Interpretable Machine Learning

MOVINGFROM MYTHOSTO DIAGNOSTICS

VALERIECHEN, JEFFREYLI, JOONSIKKIM, GREGORYPI UMB AMFETTAI WAI KAR

heemergenceofmachinelearningasasocietychangingtechnologyinthepastdecadehas triggeredconcernsaboutpeople'sinabilityto understandthereasoningofincreasinglycomplex models.ThefieldofIML(interpretablemachine

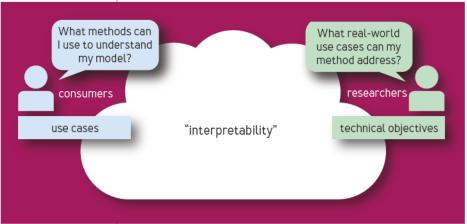
learning)grewoutoftheseconcerns,withthegoalof empoweringvariousstakeholderstotackleusecases,such asbuildingtrustinmodels,performingmodeldebugging, andgenerallyinformingrealhumandecision-making. 7,10,17

YetdespitetheflurryofIMLmethodological developmentoverthepastseveralyears,astark disconnectcharacterizesthecurrentoverallapproach. Asshowninfigure1,IMLresearchersdevelopmethods thattypicallyoptimizefordiversebutnarrowtechnical objectives,yettheirclaimedusecasesforconsumers remainbroadandoftenunderspecified. Echoingsimilar critiquesaboutthefield, ¹⁷ithasthusremaineddifficult toevaluatetheseclaimssufficientlyandtotranslate methodologicaladvancesintowidespreadpractical impact.

This article outlines apath forward for the ML



FIGURE 1: THE GAP BETWEEN IML CONSUMERS AND RESEARCHERS



community to address this disconnect and foster more widespread adoption, focusing on two key principles:

■ Embrace a "diagnostic" vision for IML. Instead of aiming to provide complete solutions for ill-defined problems, such as "debugging" and "trust," the field of IML should focus on the important, if less grandiose, goal of developing a suite of rigorously tested diagnostic tools. By treating IML methods as diagnostics, each can be viewed as providing a targeted, well-specified insight into a model's behavior. In this sense, these methods should then be used alongside and in a manner similar to more classical statistical diagnostics (e.g., error bars, hypothesis tests, methods for outlier detection), which have clearer guidelines for when and how to apply them. Under this vision, existing IML methods should be treated as potential diagnostics until they are rigorously tested.

■ RigorouslyevaluateandestablishpotentialIML diagnostics.IMLresearcherstypicallydevelopand evaluatemethodsbyfocusingonquantifiabletechnical objectives(e.g.,maximizingvariousnotionsoffaithfulness oradherencetosomedesirableaxioms 4,18,24).While theseIMLmethodsgenerallytargetseeminglyrelevant aspectsofamodel'sbehavior,itisimperativetomeasure theireffectivenessonconcreteusecasesinorderto demonstratetheirutilityaspracticaldiagnostics.

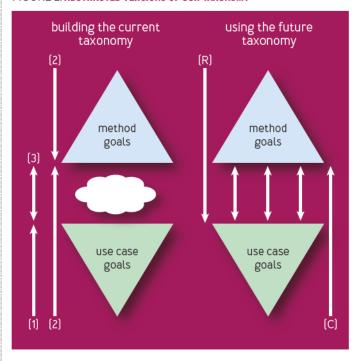
Thesetwoprinciplesmotivatedustofirstillustrate our diagnostic vision via an incomplete taxonomy that synthesizes found at ional works on IML methods and evaluation. The taxonomy (shown at an abstract level in the left side of figure 2) serves as not only a template for building an explicit mapping between potential IML diagnostics and specificuse cases, but also a too lounify studies of IML's useful ness in real-world settings. Further, the incompleteness of the current taxonomy emphasizes the need for researchers and consumers to work together to expand the coverage of the use case organization (i.e., in the "Use Case Goals"), and to establish connections between methods and use case by following the proposed work flow below.

- (1) Problemdefinition, where researchers work with consumer stodefine a well-specified target usecase.
- (2) Methodselection, wherethey identify potential IML methods for a targetuse case by navigating the methods part of the taxonomy and lor leveraging previously established connections between similar use cases and methods.
- [3] Methodevaluation.whereresearchersworkwith

ai 4oF29







consumers to test whether selected methods can meet target use cases.

Then, the latter part of this article includes an extensive discussion about best practices for this IML workflow to flesh out the taxonomy and deliver rigorously tested diagnostics to consumers. Ultimately, there could be an increasingly complete taxonomy that allows consumers (C) to find suitable IML methods for their use cases and helps researchers (R) to ground their technical work in real applications (as seen on the right side of figure 2). For instance, Table 1 highlights concrete examples

TABLE1: **EXAMPLEUSECASES**

COMPUTERVISION: CLASSIFIERTODETECTOBJECTSINIMAGES	
UseCase:	Debugthemodelbyidentifyingifitusespositivespurious correlations(i.e.,reliesonobjectYtodetectobjectX).
DiagnosticInsight:	Whenfeatures(i.e.,spuriousobjects)arepresentormissing, howdoesthisaffectaspecificprediction?
BANKLENDING: CLASSIFIERTOGRANT/DENYLOANSTOCLIENTS	
UseCase:	Recommendactionablerecourseforanindividualtoget aloanaftertheyhavebeenpreviouslydenied.
DiagnosticInsight:	What(low-cost)changescananindividualmaketoachieve adesiredoutcome?
COMPUTATIONALBIOLOGY: CLUSTERINGTOANALYZESINGLE-CELLRNASEQUENCES	
UseCase:	Verifywhetherdifferencesbetweenclusterscorroborate knownscientificknowledge(e.g.,differentcelltypes).
DiagnosticInsight:	Whatfeaturechanges(i.e.,togeneexpression)canbemade toagroupofpointstoachieveadesiredoutcome?

ofhowthreedifferentpotentialdiagnostics, each corresponding to different types of IML methods (local feature attribution, local counterfactual, and global counterfactual, respectively), may provide useful in sights for three usecases. In particular, the computer vision use case from Table 1 is expanded upon a sarunning example.

BACKGROUND

Anincreasinglydiversesetofmethodshasbeenrecently proposed and broadly classified as part of IML. Multiple concernshave been expressed, however, in light of this rapid development, focused on IML's underlying foundations and the gap between research and practice.

Critiquesofthefield'sfoundations

ZacharyC.Liptonprovidedanearlycritique,highlighting thatthestatedmotivationsofIMLwerebothhighly variableandpotentiallydiscordantwithproposed methods.¹⁷MayaKrishnanaddedtothesearguments fromaphilosophicalangle,positingthatinterpretability asaunifyingconceptisbothunclearandofquestionable usefulness.¹⁵Instead,morefocusshouldbeplacedonthe actualendgoals,forwhichIMLisonepossiblesolution.

Gapsbetweenresearchandpractice

Multipleworkshavealsohighlightedimportant gapsbetweenexistingmethodsandtheirclaimed practicalusefulness.Somehavedemonstratedalackof stability/robustnessinpopularapproaches. ^{1,2,16}Others, meanwhile,discusshowcommonIMLmethodscanfailto helphumansintherealworld,boththroughpointingout hiddenassumptionsanddangers, ^{6,21}aswellasconducting casestudieswithusers. ^{5,14}

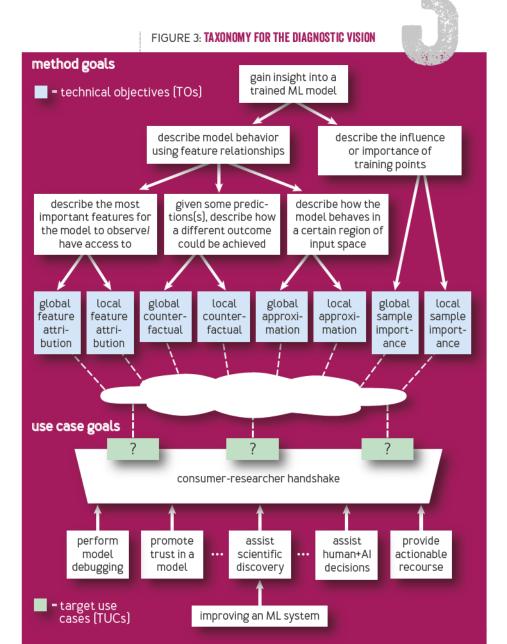
Morerecently,manyreviewpapers 3,10,19,20 have attemptedtocleanupandorganizeaspectsofIMLbut largelydonotaddresstheseissuesheadon.Incontrast, thereframingofIMLmethodsasdiagnostictoolsproposed herefollowsnaturallyfromtheseconcerns.Notably, thisarticleembracestheseemingshortcomingsof IMLmethodsasprovidingmerely "facts" 15 or "summary statistics" 21 aboutamodel, and instead focuses on the practical questions of when and how the semethods can be useful

ADIAGNOSTICVISIONFORIML Inourvision, adiagnostic is a tool that provides some

actionableinsightaboutamodel.Asananalogy,consider thesuiteofdiagnostictoolsatadoctor'sdisposalthat similarlyprovidesvariousinsightsaboutapatient.An X-raycouldbeusefulforidentifyingbonefractures,while aheart-ratemonitorwouldbehelpfulforidentifyingan irregularrhythm.Importantly,neithertoolenablesthe doctortobroadly "understand" aperson'shealth,buteach canbeusefulifappliedproperlytoawell-scopedproblem. Asimilarlyrigorousapproachtoestablishingconnections betweenIMLmethodsandwell-definedusecasesis imperativefortheIMLcommunity.

Tobeginsuchapursuit, let's identify and reconcile the manymethodgoalsandusecasegoalsthatyoumight currently encounter. Based on contemporary practices anddiscourse,let'sconsiderataxonomythatorganizesseparate hierarchiesforthemethodgoalsatthetopendandusecase goalsatthebottomend(asillustratedinfigure3).Whilethe diagnosticvisionforthefieldideallyinvolvesaclearlydefined setofusecasesandarobustsetofconnectionsbetween thesetwosides.acloudisusedtoillustratethecurrent. overalllackofwell-establisheddiagnostics. Moving forward, thegoalforresearchersandconsumersistoconduct principledstudiesfocusedonfillinginbothgaps.First,they shouldworktorefinethecurrentorganizationofusecases, consisting of an incomplete list of common ly discussed broadgoals, by defining more well-specified target usecases (showningreen)viatheconsumer-researcherhandshake. Second, they should aim to establish explicit connections betweenthesetargetsandtechnicalobjectives(shownin bluel.

ai 8 o F 2 9



Methodgoals

EachIMLmethodprovidesaspecifictypeofinsightinto agivenmodel. The forms of these insights help provide a hierarchical organization that divides the set of existing IML methods into eight method clusters. In the diagnostic vision, each method cluster is thought of broadly as a class of diagnostics that addresses a TO (technical objective). Later, each TO is described in a way that allows individual method goals to be specified.

Hierarchicalorganization

Thetopendofthetaxonomyaimstodifferentiatebetween thevariousperspectivesthatexplanationsprovide, basedonthreefactorscommonlydiscussedinexisting literature: 3.9.11

- **Explanationrepresentation.** Modelexplanations are typically given interms of either *feature relationships* between inputs and outputs or *training examples*.
- Typesoffeaturerelationships.Inthecontextof explanationsbasedonfeaturerelationships,thereare threedistinctapproachesforexplainingdifferentaspects ofthemodel'sreasoning:(1) featureattribution ;(2) counterfactual;and(3) approximation .Notethatbecause the IML community focuses lessongenerating example-based explanations, we consider one main grouping along that branch: sample importance explanations .
- **Explanationscale.** Explanationsvaryintermsofthe scaleofthedesiredinsights, with their scoperanging from *local*(i.e., for an individual instance) to *global* (i.e., for a well-defined region of the input space).

AttheleafnodesaretheTOs,classesofgoalsthat

arepreciseenoughtobegenerallylinkedtoa*method* cluster thatmostdirectlyaddressesthem.Intotal, thereareeightTOs/methodclustersthatcapturea largeportionofthegoalsofexistingIMLmethods. Thereareacoupleofimportantnuancesregardingthe characterizationofTOs.

- First, although TOs and method clusters are one-to-one in the proposed taxonomy, it is important to explicitly distinguish these two concepts because of the potential for cross-cluster adaptation. This notionarises because it is frequently possible for a method to, in an adhoc fashion, be adapted to address a different TO.
- → Second, each TO should be thought of as defining a class of related goals. Indeed, for a given TO, we hypothesize what some of the key technical detail (s) are that must be considered toward fully parametrizing meaning fully different instantiations of the same broader goal. These important technical details, taken to gether with the TO, allowy out ode fine individual proxymetrics that reflect the desired properties of your explanations. Proxymetrics can then serve a stractable objective functions for individual method stooptimize for, as well as measures of how well any method addresses a particular instantiation of the TO.

Technicalobjectives

Thefollowingisanoverviewofthe TOs (and their technical details) that correspond to various method clusters. Because of the overlaps in content, local and global versions of the same general method type lobjective are grouped to gether. (For more details and examples of specific methods for each, see our longer-form paper,

- "InterpretableMachineLearning:MovingfromMythosto Diagnostics,"byChen,etal. 8).
- → Featureattributionexplanations addresshowthe model'sprediction(s)areaffectedwhenfeaturesare present(ormissing),i.e.,how"important"eachfeatureisto themodel'sprediction(s).Often,measuresofimportance aredefinedbasedonhowthemodel'sprediction(s) changesrelativetoitspredictionforsomebaselineinput. Thebaselineinputissometimesimplicitanddomain specific(e.g.,allblackpixelsforgrayscaleimagesorthe meaninputintabulardata).Thus,thetechnicaldetailsare boththeprecisenotionof" importance" andthechoice ofthebaselineinput .Relevantproxymetricstypically measurehowmuchthemodelpredictionchangesfor differenttypesofperturbationsappliedtotheindividual (orthetrainingdata)accordingtothe "importance" values ascomputedbyeachmethod.
- → Counterfactualexplanations addresswhat "low-cost" modificationcanbeappliedtodatapoint(s)toachieve adesiredprediction. Themost commonte chnical detail is the specific measure of cost, and the most common proxymetric is how of tenthe counterfactual changes the model's prediction(s).
- Approximationmethods addresshowtosummarizethe modelbyapproximatingitspredictionsinaregion,either locallyaroundadtapoint,globallyaroundasmanypoints aspossible,oracrossaspecificregionoftheinputspace. Thesemethodsrequirethetechnicaldetailsofboththe definitionoftheregion andthesimplefunction's model family. For local approximation, a canonical metric is local fidelity, which measures how well the method predicts

- withinacertainneighborhoodofadatapoint. Forglobal approximation, aproxymetric is coverage, which measures how many datapoints the explanation applies to.
- Sampleimportancemethods address which training points most influence amodel's prediction for either an individual point or the model as a whole. Technical details differ from method to method, so it is difficult to identify a uniformaxis of variation. The semethods can be evaluated with proxymetrics that represent the useful ness of the provided explanations through simulated experiments of finding corrupted data points, detecting points responsible for data distributions hifts, and recovering high accuracy with the samples considered important.

Howdoby-designmethodsfitin?

Whiletheydonothaveacorrespondingmethodcluster inthistaxonomy, it is important to discuss another family of IML methods that propose models that are themselves interpretable by design. ²¹ The differentiating property of the semodels from the post-hoc methods referenced in the above section is that the TO(s) of the seap proaches is intrinsically tied to the model family itself; hence, the models are interpretable by design only in that they satisfy said TO(s). That said, by-design methods also fit in to this framework and should be viewed as a different way to answer the same TO sin the taxonomy. When by-design methods are proposed or used, they should clearly specify which TO sthey intend to address.

Usecasegoals

MuchofthecurrentdiscourseonIMLusecasessurrounds

differentiatingfairlybroadgoals,suchasdebugging models,gainingtrustofvariousstakeholders,and providingactionablerecoursetousers[figure3]. While thislevelofcategorizationrepresentsagoodstart, it is of limitedutilitybecause ittreats each of these categories as monolithic problems for IML to solve. For one, these problems are complex and should not be assumed to be completely, nor solely, solvable by IML itself. Rather, IML is but one potential set of tools that must be proven to be useful. That is, to show that an IML method is an effective diagnostic, specificuse cases must be identified and demonstrated. ¹⁵

Secondly,eachbroadgoalreallyincludesmultiple separatetechnicalproblems,crossedwithmanypossible practicalsettingsandconstraints.ltisunlikelythatagiven IMLmethodwillbeequallyusefulacrosstheboardforall ofthesesubproblemsanddomains.

Thus, claims of practical usefulness should ideally be specified down to the level of an adequately defined TUC (targetuse case). Like TO son the methods side, TUCs correspond to learning aspecific relevant characteristic about the underlying model (e.g., a certain property or notion of model behavior). Unlike a TO, however, they represent real-world problems that, while they can be evaluated, of ten might not be a menable to direct optimization.

For example, you can set up evaluation sto determine whether an IML method is useful for identifying a particular kind of bugin the model (e.g., positive spurious correlations), but it is not so obvious how to optimize an IML method that will succeed on those evaluations.

AWORKFLOWFORESTABLISHINGDIAGNOSTICS
Let'sturnnowtohowadiagnosticvisionforIMLcan
bemorefullyrealized,discussinghowmethodscanbe
establishedasdiagnostics,thusfillinggapsintheexisting
taxonomy.Specifically,anidealworkflowisdefinedfor
consumer-researcherteamstoconductfuturestudies
aboutIMLmethods.Itdescribeshowthetaxonomycan
guidebestpracticesforeachofthethreekeysteps:[1]
problemdefinition;[2]methodselection;and[3]method
evaluation.Thisworkflowappliesbothtoteamswhowish
tostudyexistingIMLmethodsandtothoseproposing
newones.

Arunningexamplehelpscontextualizethisdiscussion, buildingonthecomputervisionmodeldebuggingexample fromtable1.Modeldebuggingisnotonlyacommon consumerusecase, ^{7,13}butalsoawell-groundedone.Itisa naturalstartingpointbecauseoftheversatilenatureof itsassumedconsumer,datascientists,whotypicallyhave bothsubstantialMLknowledgeanddomainexpertise, minimizingthecommunicationgapbetweenthedata scientistandthelMI researcher.

Step1:ProblemDefinition

Animportantfirststepforanyprincipledstudyistodefine awell-specifiedTUC. This process is called the *consumer-researcherhandshake* (figure 3), where researchers work with consumer stoprogressively refine the latter's real-world problems into relevant TUCs. In this process, some helpful pieces of information include: the data available, the ML pipeline used, and the domain knowled gerequired to perform evaluations. Ultimately, amore fleshed-out

taxonomy will help researchers have more concrete use cases at hand to motivate their method development, and consumers will have more realistic guidance on what IML can and cannot do for them.

Running example: Consider a data scientist who wants to debug her image-based object detection model. She hopes to leverage the expertise of an IML researcher, but as shown in a hypothetical version of the use cases part of

FIGURE 4: CONSUMER-RESEARCHER HANDSHAKE FOR OUR RUNNING EXAMPLE

target use case: target use case: positive negative correlations correlations

use case goal: correlations between image objects use case goal: use case goal: spurious correlation bad edge case detection detection use case goal: perform model debugging

thetaxonomy(figure4), theumbrellaofmodeldebugging includesseveral subproblems, such as detecting spurious correlations and identifying badedge-case behavior. Thus, the team of researcher and datascient is tneed stoidentify a TUC that is more specific than "perform modeldebugging" by identifying exactly what notion of "bug" the IML method should detect. Through the consumer-researcher hands hake, it arises that the datascient is tis concerned that the model might not be making correct decisions based on the actual target objects, but rather is relying on correlated objects that also happen to be present. For example, the model might be using the presence of a person as an indicator that the reisatennis racket in the image, in stead of the racket itself.

Thisinformationallowstheteamtonavigatethe relevantbranchesofthetaxonomy. Here, by considering the datascientist's concern, they first narrow the goal from model debugging to detecting spurious correlations. Then, by also taking into account the specific setting (i.e., the presence of the tennis racket at the same time as the tennis player), they are able to arrive at a further specified use case of detecting spurious correlations between two positively correlated objects (marked by the white border in figure 4). In this case, the team takes care to differentiate this from the analogous problem of detecting reliance on negatively correlated objects, reasoning that the latter is fundamentally different (i.e., it is harder to tell whether the output depends on an object if the co-occurrence sare rare in the first place).

Step2:MethodSelection

AfteraTUChasbeenproperlydefined,thenextstepisto considerwhichIMLmethodsmightbeappropriate. This doesassumethatIMLmethodsarenecessary—thatis, the teamshouldhavedemonstrated that the TUC presents challenges to more "trivial" or conventional diagnostics. For example, Bansal, et al. found model confidence to be a competitive baseline against dedicated interpretability approaches for AII humandecision—making teams. 5

Ifnon-IMLdiagnosticsareunsuccessful, the taxonomy canbeusedintwowaystoselectmethods.First, researchersandconsumerscan.asadefault.traversethe methodspartofthetaxonomytoidentifytheTOs(and thus,respectivemethodclusters)thatmightbestalign with the TUC. Doings os hould rely on the researcher's bestjudgmentinapplyingpriorknowledgeandintuitionabout variousmethodtypestotrytonarrowdownthesetof potentialTOs.Ifamethodisbeingproposed,itshouldbe mappedtotheappropriatemethodcluster, and the same selection process should follow. Second, the team can alsonavigatestartingfromtheusecasespart, leveraging and expandingonconnectionsestablishedbypreviousstudies. Naturally, if some methods have already been shown toworkwellonaTUC,thenthose(orsimilar)methods provideimmediatebaselineswhenstudyingthesame(or similarlusecases.

Ineithercase,animportant—yetsubtle—choicemust thenbemadeforeachmethod:exactlyhowitsresulting explanationsshouldbeinterpreted(i.e.,whichTOisbeing addressed). As discussed in the section about method goals, amethod belonging to a specific cluster may most naturally

addresstheassociatedTO,butitisalsopossible,andindeed commonplace,toattemptcross-clusteradaptation for addressingotherTOs.Unfortunately,whilesuchadaptations maybeusefulattimes,theyareoftenperformedinan adhocfashion.Specifically,thedifferencesbetweenthe technicaldetailsofeachTOareoftenoverlookedinthe adaptationprocess,asillustratedviathefollowingtwo examples(andinmoredepthinChen,etal. 8).

First, youmighttry touse "feature importance weights," via SHAP (Shapleyadditive explanations), ¹⁸ as linear coefficients in a local approximation. Such an adaptation assumes that the notion of local "importance" also can reflect linear interactions with features on the desired approximation region. This is not necessarily guaranteed by SHAP, however, which in steaden forces a different set of game-theoretic desiderata on the importance values and may be set up to consider a quite disparate set of perturbations compared to the target approximation region.

Conversely, you can think of saliency maps via vanilla gradients ²³ as an adaptation in the opposite direction. These saliency maps, alocal approximation where the effective neighborhood region is extremely small, are more popularly used to address local feature attribution objectives, such as identifying which parts of the image are affecting the prediction the most. This adaptation, however, carries an underlying assumption that the pixels with the largest gradients are also the most "important." This approximation may not be accurate because the local shape measured by the gradient is not necessarily indicative of the model's behavior near abase line in put that is farther away.

Runningexample: Inthisscenario, suppose that there have been no previously established results for detecting positive spurious correlations. The team follows the methods part of the taxonomy to generate hypotheses for which types of local explanations be stsuit their needs for under standing individual images. They decide against approximation-based objectives, because as the inputs vary in pixel space, simple approximations are unlikely to hold or be semantically meaning ful across continuous local neighborhoods. They choose feature attribution because they be lieve that visualizing the features that the model deems most important would be useful for detecting the se types of spurious correlations.

Theteamproposesamethodinthelocalcounterfactual methodclusterthatidentifiesthesuper-pixelsthatmust changeinordertoflipthepredictionfrom "tennisracket" to "notennisracket." By "visualizing" the counterfactual explanation likeas aliency map, the teamper forms across-cluster adaptation to interpret the counterfactual as a feature attribution explanation. To do so, they are assuming that the most changed features are also the most important for detecting the tennisracket. They reason that a feature attribution explanation would be a more intuitive format for the data scientist for this TUC. Interms of comparison, the feature attribution method that the teams elects for comparison is Grad-CAM (gradient-weighted class activation mapping), 22 which also produces as a liency map.

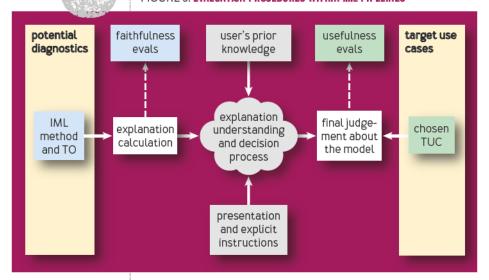
Step3:MethodEvaluation

Onceappropriatemethod(s)havebeenchosen,thelast stepistoevaluatethem. Evaluation is the crucial step

of testing whether proposed methods can actually help address the specified TUC. Yet despite its importance, this step is often carried out in manners incongruent with the properties it claims to test. One common mistake is that the evaluation of an explanation's faithfulness (i.e., ability to meet a specified TO) is often problematically conflated with the evaluation of its usefulness (i.e., applicability for addressing practical TUCs). While both may play important roles, they target fundamentally different claims.

Considering these evaluations within the overall pipeline of an IML application (as shown in figure 5) addresses this kind of mistake. It first highlights differences in goals of these evaluations by connecting back to the taxonomy presented in this article; faithfulness

FIGURE 5: EVALUATION PROCEDURES WITHIN IML PIPELINES



correspondstomeetingobjectivesofaspecificTOinthe methodspart,andusefulnesscorrespondstomeetingthe TUCintheusecasespart.Then,italsolaysoutthevarious movingcomponentsthataffecteachtype,withgrayboxes denotingcomponentsthatrequiremorecarefulstudy.This servestogroundhoweachmaybecarriedout,whichwe discussingreaterdetailnext.

Faithfulnessevaluations are performed with respect to a proxymetric specified using the relevant technical details from the target TO class. For example, if the goal were to show the usefulness of an approximation-based explanation adapted as a counterfactual, the faithfulness evaluation should be with respect to a counterfactual proxymetric. Referring to the terminology from Doshi-Velezand Kim, the setypes of evaluations are called functionally grounded—that is, involving automated proxy tasks and no humans. While such evaluations are easiest to carry out, they come with key limitations.

Ingeneral, you should expect that a method would perform well at least on a proxy for its selected TO, and, naturally, those methods that do not directly target this specific proxy will likely not perform as well. An explanation's performance can also be faultily compared with another's as a result of unfair or biased settings of technical details. As an example, although GAMs (generalized additive models) ¹² and linear models both provide local approximations, comparing the semethods only in the context of fidelity ignores the fact that GAMs potentially generate more "complicated" explanations.

Further, while faithfulness evaluations can act as a first-steps anity check before running more costly

usefulnessevaluations, showing that amethod is faithful to the model alone is not conclusive of the method's real-world usefulness until a direct link is established between the corresponding proxy and TUC. Once these links are established, these proxies can then be used more confidently to help rule out bad set ups before performing expensive usefulnesse valuations.

Usefulnessevaluations, incontrast to faithfulness, measureauser's successina pplying explanations to the specified TUC. Since they are ultimately an evaluation of what auser does with an explanation, usefulness depends crucially on factors, such as the user's prior knowledge—for example, their domain and ML/IML experience. Again, using terminology from Doshi-Velezand Kim, ⁹ users' perspectives can be incorporated through studies on real human sperforming simplified or actual tasks (i.e., human-grounded or application-grounded evaluations, respectively). In particular, as part of conducting usefulness studies, you would need to consider how users might act differently depending on the presentation of the explanation and explicit in structions that are provided.

Ashighlightedbythecloudinfigure5, exactlyhow userstranslateexplanationcalculations (intheirminds) to their final judgments remains murky. This motivates further research relating to better understanding what users understand explanations to tell the mandhow they actupon these understandings. Then, when establishing new diagnostics, these assumptions / limitations should be clearly spelled outfor when researchers use the method in a future study and when consumers deploy the method.

Motivatedbythesechallenges, researchersmightwant toalsoconsideranothertypeofusefulnessevaluation: simulationevaluation. This is an algorithmic evaluation on a simulated version of the real task where success and failure are distilled by a domain expertinto a measurable quantity (a sillustrated in the running example). This type of evaluation is still based on the real task but is easier and potentially more reliable to run than user studies.

Bysimulatingtheusersandtheirdecision-making processalgorithmically, thus controlling some noisier aspects of usefulnessevaluation, researchers may be able to better understand why their methods are "failing": is it because of the algorithmit selfortheusers' actual decision-making process?

Overall, successon these various levels of evaluations provides evidence for establishing a connection between the method in question and the TUC. Specifically, the team should check to see if the proxymetrics considered earlier were correlated to successon the TUC. If so, this would provide evidence for whether the proxymetrics considered should be used again in future studies, connecting faithfulness and usefulness evaluations.

Runningexample: Theteamfirstperformsseparate localfeatureattribution faithfulness evaluations for both methods using the respective notions of importance that each defines. For example, for the proposed method, the teamensures that each generated explanation faithfully carries out its intended TO of identifying the effect of the presence or absence of a super-pixel. Goodperformance on any proxymetric, however, does not conclusively imply

goodperformanceontheactual TUC, so the team turns to usefulnesse valuation.

Theteamfirstconductsasimulationevaluation, where datasets are created that contained ther (artificially induced) positive correlation between a pair of objects or no such correlations. By carefully controlling the training and validation distributions, they can automatically verify whether a model has learned the problematic behavior they want to detect. Then they can define as coring function for the explanations (i.e., how much attention they pay to the spurious object) and measure how well that score correlates with the ground truth for each explanation.

Second, the team runs a human study with multiple models where they know the ground truth of which ones uses purious correlations. They scored a tascientists based on whether they are able to use each explanation generated by the counterfactual versus Grad-CAM to identify models that uses purious correlations. If the methods are successful on the human studies, the team has demonstrated the connection between the mand the TUC of detecting positively correlated objects.

CONCLUSION

AssumingadiagnosticvisionforIML, the taxonomy presented here is a way to clarify and begin bridging the gap between methods and use cases. Further, this article discusses be st practices for how the taxonomy can be used and refined over time by researchers and consumers to establish which methods are useful for which use cases. As the taxonomy is fleshed out via more studies by consumer researcher teams, our vision is that it will be increasingly

usefulforbothpartiesindividually(figure2,right).Overall, thegoalistopromotebetterpracticesindiscovering, testing,andapplyingnewandexistingIMLmethodsmoving forward.

References

- Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., Kim, B. 2018. Sanitychecks for saliency maps. In Proceedings of the 32 nd International Conference on Neural Information Processing Systems, 9525-9536; https://dl.acm.org/doi/10.5555/3327546.3327621.
- Alvarez-Melis, D., Jaakkola, T. 2018. On the robustness of interpretability methods. arXiv:1806.08049; https://arxiv.org/abs/1806.08049.
- 3. Arya,V.,Bellamy,R.K.,Chen,P.-Y.,Dhurandhar,A., Hind,M.,Hoffman,S.C.,Houde,S.,Liao,Q.V.,Luss, R.,Mojsilovi´c,A.,etal.2019.Oneexplanationdoes notfitall:atoolkitandtaxonomyofAlexplainability techniques.arXiv:1909.03012; https://arxiv.org/pdf/1909.03012.pdf.
- Bach,S.,Binder,A.,Montavon,G.,Klauschen,F.,Müller, K.-R.,Samek,W.2015.Onpixelwiseexplanationsfor non-linearclassifierdecisionsbylayer-wiserelevance propagation. *PloSONE* 10(7):e0130140;https:// journals.plos.org/plosone/article?id=10.1371/journal. pone.0130140.
- Bansal,G.,Wu,T.,Zhu,J.,Fok,R.,Nushi,B.,Kamar, E.,Ribeiro,M.T.,Weld,DS.2020.Doesthewhole exceeditsparts?TheeffectofAlexplanationson complementaryteamperformance.arXiv:2006.14779; https://arxiv.org/pdf/2006.14779.pdf.

- Barocas,S.,Selbst,A.D.,Raghavan,M.2020.Thehidden assumptionsbehindcounterfactualexplanationsand principalreasons.InProceedingsoftheConference onFairness,Accountability,andTransparency ,80-89; https://dl.acm.org/doi/abs/10.1145/3351095.3372830.
- Bhatt,U.,Xiang,A.,Sharma,S.,Weller,A.,Taly,A., Jia,Y.,Ghosh,J.,Puri,R.,Moura,J.M.F.,Eckersley,P. 2020.Explainablemachinelearningindeployment. In ProceedingsoftheConferenceonFairness, Accountability,andTransparency ,648-657;https://dl.acm.org/doi/abs/10.1145/3351095.3375624.
- 8. Chen, V., Li, J., Kim, J.S., Plumb, G., Talwalkar, A.2021. Interpretable Machine Learning: Moving from Mythos to Diagnostics. arXiv:2103.06254.
- Doshi-Velez, F., Kim, B. 2017. Towards a rigorous science of interpretable machine learning. arXiv:1702.08608; https://arxiv.org/pdf/1702.08608.pdf.
- Gilpin, L.H., Bau, D., Yuan, B.Z., Bajwa, A., Specter, M., Kagal, L. 2018. Explaining explanations: an overview of interpretability of machine learning. In Fifth IEEE International Conference on Data Science and Advanced Analytics; https://lieeexplore.ieee.org/document/8631448.
- Guidotti,R.,Monreale,A.,Ruggieri,S.,Turini, F.,Giannotti,F.,Pedreschi,D.2018.Asurveyof methodsforexplainingblackboxmodels.ACM ComputingSurveys 51(5),1–42;https://dl.acm.org/ doi/10.1145/3236009.
- 12. Hastie, T.J., Tibshirani, R.J. 1990. Generalized Additive Models. *Monographson Statistics and Applied Probability*, 43. Chapman and Hall ICRC.

- 13. Hong, S.R., Hullman, J., Bertini, E.2020. Human factors inmodel interpretability: industry practices, challenges, and needs. In *Proceedings of the ACM on Human-Computer Interaction* 4 (CSCW1), 1–26; https://dl.acm.org/doi/10.1145/3392878.
- Kaur, H., Nori, H., Jenkins, S., Caruana, R., Wallach, H., Wortman Vaughan, J. 2020. Interpreting interpretability: understanding datascientists' use of interpretability tools for machine learning. In Proceedings of the CHI Conference on Human Factors in Computing Systems, 1-14; https://dl.acm.org/doi/ abs/10.1145/3313831.3376219.
- Krishnan, M. 2020. Against interpretability: a critical examination of the interpretability problem in machine learning. *Philosophy& Technology* 33,487–502; https://link.springer.com/article/10.1007/s13347-019-00372-9.
- Laugel, T., Lesot, M.-J., Marsala, C., Detyniecki, M. 2019. Issueswithpost-hoccounterfactual explanations: adiscussion.arXiv:1906.04774; https://arxiv.org/pdf/1906.04774.pdf.
- 17. Lipton, Z.C. 2018. The mythosof model interpretability. *ACMQueu* e16(3), 31–57; https://queue.acm.org/detail.cfm?id=3241340.
- Lundberg,S.M.,Lee,S.-I.2017.Aunifiedapproachto interpretingmodelpredictions.In AdvancesinNeural InformationProcessingSystems 30; https://papers.nips.cc/paper/2017/ hash/8a20a8621978632d76c43dfd28b67767-Abstract. html.
- 19. Mohseni, S., Zarei, N., Ragan, E. 2020. Amultidisciplinary surveyand framework for designand evaluation

- ofexplainableAlsystems. *ACMTransactionson InteractiveIntelligenceSystems* 1(1); https://arxiv.org/pdf/1811.11839.pdf.
- 20. Murdoch, W.J., Singh, C., Kumbier, K., Abbasi-Asl, R., Yu, B.2019. Interpretable machine learning: definitions, methods, and applications. In *Proceedings of the National Academy of Sciences* 116 (44), 22071-22080; https://www.pnas.org/content/116/44/22071.
- 21. Rudin, C.2019. Stopexplaining blackboxmachine learning models for high stakes decisions and use interpretable models in stead. *Nature Machine Intelligence* 1,206–215; https://www.nature.com/articles/s42256-019-0048-x.
- 22. Selvaraju,R.R.,Cogswell,M.,Das,A.,Vedantam,R., Parikh,D.,Batra,D.2017.Grad-CAM:visualexplanations fromdeepnetworksviagradient-basedlocalization.In *IEEEInternationalConferenceonComputerVision*,618-626; https://iieeexplore.ieee.org/document/8237336.
- 23. Simonyan, K., Vedaldi, A., Zisserman, A. 2013. Deepinside convolutional networks: visualising image classification models and saliency maps. arXiv:1312.6034; https://arxiv.org/abs/1312.6034.
- 24. Sundararajan, M., Taly, A., Yan, Q. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning*; http://proceedings.mlr.press/v70/sundararajan17a.html.

ValerieChen is a Ph.D. student in the Machine Learning Department at Carnegie Mellon University. Her research is on interpretability techniques to aid human decision-making and, more broadly, as a way to study the societal impact of machine learning.

JeffreyLi is a Computer Science Ph.D. student at the University of Washington. He is interested in topics that address challenges limiting the deployment of machine learning in practice, including learning from weak sources of supervision and interpretable machine learning.

JoonSikKim is a Ph.D. student in the Machine Learning Department at Carnegie Mellon University. He is interested in methods that can facilitate the understanding of complex machine learning models and their implications for model interpretability and fairness.

GregoryPlumb is a Ph.D. student in the Machine Learning Department at Carnegie Mellon University. His research focuses on explainable machine learning, with an emphasis on developing novel techniques for model debugging.

AmeetTalwalkar is an assistant professor in the Machine Learning Department at Carnegie Mellon University. His current work is motivated by the goal of democratizing machine learning, with a focus on topics related to automation, interpretability, and distributed learning.

Copyright ${\hbox{$\mathbb Q$}}$ 2021 held by owner/author. Publication rights licensed to ACM.

