

Paper					
Ontology explanationExperience.rdf aimodel.rdf aimodevaluation.rdf explainer.rdf user.rdf behaviour_tree.rdf userevaluation.rdf	ExplanationExperience hasDescription	Description hasAIModel	AIModel trained on	Dataset hasDataType	DataType number of features number of instances
		hasExplainer (needExplainer?)	Explainer	<i>hasOutputType</i> hasPortability hasConcurrentness hasPresentation <i>hasExplanationScope</i> targetType	<i>Explanation</i> Portability ExplainerConcurrentness InformationContentEntity <i>Explanation Scope</i> Explanation Target
		hasUser	User asks	UserQuestion hasTarget hasType Intent hasType Technical Facilities can handle process Domain Knowledge level of domain knowledge level of AI knowledge	UserQuestionTarget QuestionType IntentType ExplanationModality KnowledgeLevel KnowledgeLevel
	hasSolution	Solution hasExplainer	Explainer utilises	ExplainabilityTechnique hasType hasOutputType hasPortability hasConcurrentness hasPresentation <i>hasExplanationScope</i> targetType isCompatiblewithFeatureTypes hasComplexity	ExplainabilityTechniqueType Explanation Portability ExplainerConcurrentness InformationContentEntity <i>Explanation Scope</i> Explanation Target DataType ComputationalComplexity
	hasDone	UserEvaluation basedOn	Metric measures	Dimensions	
Code Availability					

Make it personal: a social explanation system applied to group recommendations

PSIE

PsieGroupRecommendationExplanationExperience

Recommend movies to groups based on social knowledge

HappyMovieRecommenderSystem

HappyMovieDataset

Tabular

N/A

N/A

GroupRecommendation

Recommendation

Recommend movies to groups based on social knowledge

SocialGroupRecommender

Knowledge based Recommender

3.86, 3.69, 3.56, 3.89

Likert scale (5 points - being 5 strongly agree)

Usefulness, Decision process, Reusability, Usability

Content Based Explanation

model-specific

post-hoc

Any

local

prediction

Moviegoer

Why does the system recommend movie Y for group X?

System Recommendation

Why Question

Understand system recommendation

Trust, Satisfaction

ScreenDisplay

Any

User profile

high

low

PsieRuleBasedTechnique

Knowledge Extraction

ContentBasedExplanation

model-specific

post-hoc

Visual/Textual

local

prediction

Tabular

N/A

Questionnaire

Usefulness/Helpfulness

N/A

DisCERN: Discovering Counterfactual Explanations using Relevance Features from Neighbourhoods

DisCERN

LungCancerRiskExplanationExperience

predict lung cancer risk given clinical data of patients

LungCancerRiskPredictionModel

LungCancerRiskDataset

Tabular

12

427

CancerRiskPrediction

Multi-class Classification

predict lung cancer risk given clinical data of patients

CancerRiskRandomForest

RandomForest

87.94

Accuracy

Performance

Counterfactual Explanation

model-agnostic

post-hoc

Computational Entity

local

prediction

Patient/Clinician

How can patient X reduce cancer risk Y to lower

System Recommendation

How/What-if Question

Reducing cancer risk

Education, Taking Action

ScreenDisplay

Any

Any/ Clinical knowledge

low/high

low

DisCERNCancerRiskExplainer

DisCERN

Feature Relevance+Example based

Counterfactual Explanation

model-agnostic

post-hoc

Computational Entity

local

prediction

Tabular

N/A

N/A

N/A

N/A

<https://github.com/RGU-Computing/DisCERN-XAI>

Evaluating Explainability Methods Intended for Multiple Stakeholders

BTTelecom

BTTelecomRecommenderExplanationExperience

recommend engineering notes to desk support staff to help on-site engineers

EngineerNoteRecommender

BTEngineeringNotes

Text

300 tf-idf features

5352

EngineerNoteRecommendation

Recommendation,classification

Predict next scenario based on description in the engineering notes

EngineerNoteRecommender

Content based Recommender, Machine Learning / term frequency, unsupervised

50.88%, 99.10%(in lab), completeness - 70% (in practice)

Accuracy with and without (NNR) class, top K Accuracy, confidence score

Performance, efficiency of scenario organisation

Neighbourhood Explanation, numerical, textual

model-agnostic

post-hoc

Any

local

prediction

Desk Agent

Why task Y is recommended as next task?

System Recommendation

Why question

why a recommendation has been made?

Transparency, Taking Action, Education

ScreenDisplay

Text, Image

BTNetworkPlannerDomainExpert, BTFieldEngineerDomainExpert, BTDeskAgentDomainNovice

high, high, low

low

BTRecommenderExplainer

BTContentSimilarityBasedTechnique: confidence score, feature-importance, summarisation of sim/difs

Knowledge Extraction + Feature Relevance

Neighbourhood Explanation

model-agnostic

post-hoc

Content + Similarity

local

prediction

Text

N/A

see Notes

Question to get feedback

Usefulness/Educatingness/Efficiency

N/A

Directing exploratory search: Reinforcement learning from user interactions with keywords

SciNet

DocumentSearchExplanationExperience

Determine the most related documents given a set of keywords

SciNetSearchEngine

WebOfScienceDataset

Documents

7 things, title, abstract, author names, publication year, publication forum, article, keywords

50million

DocumentRetrieval

InformationRetrieval

Determine the most related documents given a set of keywords

SciNetReinforcementLearning

Reinforcement Learning

0.71

Kappa

agreement between expert and system

Neighbourhood Explanation

model-specific

ante-hoc

visual

local

prediction

Scientist

Why was this result X retrieved for this query Y?

System Recommendation

Why Question

Understand system prediction

Effectiveness, Satisfaction

ScreenDisplay

Any

Search Domain

high

low

SciNetReinforcementLearning

Knowledge Extraction

Neighbourhood Explanation

model-specific

ante-hoc

Interactive visual

local

prediction

Documents

N/A

Questionnaire

Usability/Quality of user experience

N/A

Visualizing Recommendations to Support Exploration, Transparency and Controllability

TalkExplorer

TalkExplorerExplanationExperience
Recommend papers based on content and social connections
ConferenceNavigator3RecommenderSystem
ConferenceNavigator3Dataset
Tabular
N/A
N/A
TalkPaperRecommendation
Recommendation
Recommend papers based on content and social connections
CN3ContentBasedRecommender
Tf-idf + kNN
N/A
N/A
N/A

Neighbourhood Explanation
model-specific
ante-hoc
Image
local
prediction
Conference attendee

Why does the system recommend paper Y to user X?
System Recommendation
Why Question
Understand system prediction
Effectiveness, Transparency, Scrutability
ScreenDisplay
Any

Conference Topic Domain
high
low

TalkExplorerKNNTechnique
k-nearest Neighbour
Neighbourhood Explanation
model-specific
ante-hoc
visual
local
prediction
Tabular
N/A
Questions about explanation visualisation knowledge + tasks with TalkExplorer + Likert Scale questions about their needs
Think Aloud + Likert Scale
Effectiveness
N/A

Axiomatic Attribution for Deep Networks

IntGradImage

IGImageClassificationExplanationExperience
predict the category of a given image
IGImageClassificationModel
ILSVRC-2014
Image
89401 pixels
456182
ImageClassification
Multi-class Classification
predict the category of a given image
GoogleNet
Convolutional Neural Network
6.67%
top-5 error
Performance

Saliency Map Explanation
model-specific
post-hoc
Image
local
prediction
Any User

Why does the system predict category Y for image X?
System Recommendation
Why Question
Understand how system works
Transparency, Trust
ScreenDisplay
Any

Public Domain
Any
Any

IntegratedGradientTechnique
IntegratedGradient
Saliency Map Explanation
model-specific
post-hoc
Annotated Computational Entity
local
prediction
Image
N/A
N/A
N/A
N/A
<https://github.com/ankurtaly/Integrated-Gradients>

Axiomatic Attribution for Deep Networks

IntGradRetinopathy

DiabeticRetinopathyDetectionExplanationExperience

predict if a given medical image contains diabetic retinopathy

DiabeticRetinopathyDetectionModel

EyePACS

Image

1382400 pixels

128175

DiabeticRetinopathyDetection

Binary Classification

predict if a given medical image contains diabetic retinopathy

Fine-tubed InceptionV3

Convolutional Neural Network

90.3%, 98.1%, 99.1%

FOCP Sensitivity, FOCP Specificity, AUROC for EyePACS

Performance

Saliency Map Explanation

model-specific

post-hoc

Image

local

prediction

Clinician, Optomologist

Why does the system predict RDR for image X?

System Recommendation

Why Question

Understand how system works

Transparency, Trust, Education

ScreenDisplay

Any

Clinical Knowledge

High

Low

IntegratedGradientTechnique

IntegratedGradient

Saliency Map Explanation

model-specific

post-hoc

Annotated Computational Entity

local

prediction

Image

N/A

N/A

N/A

N/A

<https://github.com/ankurtaly/Integrated-Gradients>

Axiomatic Attribution for Deep Networks

IntGradTextClassification

QuestionCategoryExplanationExperience

predict the question category based on question text

KimQuestionCategoryPredictionModel

WikiTableQuestions dataset

text

N/A

22033

QuestionCategoryPrediction

Multi-class Classification

predict the question category based on question text

KimCNN-multichannel

Convolutional Neural Network

N/A

N/A

N/A

Saliency Map Explanation

model-specific

post-hoc

Text

local

prediction

Any User

Why does the system predict category Y for question text X?

System Recommendation

Why Question

Understand how system works

Transparency, Trust

ScreenDisplay

Any

Public Domain

low

low

IntegratedGradientTechnique

IntegratedGradient

Saliency Map Explanation

model-specific

post-hoc

Annotated Computational Entity

local

prediction

Text

N/A

N/A

N/A

N/A

<https://github.com/ankurtaly/Integrated-Gradients>

Textual Explanations for Self-Driving Vehicles

KimEtAlMethod

SelfDrivingExplanationExperience

make acceleration or change course decisions in a self-driving car based on video

Self-drivingDecisionMakingModel

BerkeleyDeepDriveDataset

Video

40 seconds (frame rate not known)

6984

Self-drivingVehicalControl (acceleration and course)

Regression

make acceleration or change course decisions in a self-driving car based on video

Deep Neural Networks with Attention

NeuralNetwork

[2.29, 0.82], [6.06, 0.47]

[Mean of absolute error, Mean of distance correlation] of Acceleration and Course

self-drivingVehicalControl performance

Introspective Explanation

model-agnostic

post-hoc

text

local

prediction

Driver

Why does the vehical system make decision X?

System Recommendation

Why Question

Understand system decision

User acceptance, Trust, Understanding, Effective communication

ScreenDisplay

text

Public Domain

high

low

LSTMTextGeneratorExplanationTechnique

Data-driven Explanation Generation

Introspective Explanation

model-agnostic

post-hoc

text

local

prediction

Any

N/A

N/A

N/A

N/A

<https://github.com/pair-code/saliency>

iBCM: Interactive Bayesian Case Model Empowering Humans via Intuitive Interaction

iBCM

iBCMGradingExplanationExperience

Cluster student assignment submissions to design grading rubric or to compose feedback

iBCMClusteringModel

iBCMAssessmentsDataset

Code

N/A

N/A

AssessmentsClustering

Clustering

iBCMClusteringMethod

BCM

N/A

N/A

N/A

Prototype Explanation

model-specific

ante-hoc

Computational Entity

cohort

prediction

Lecturer

Why does the system assign certian assessments in to one cluster? How does the system assign clusters?

Model

Why/How Question

Understand system/Understand cohort of predictions

Education/Transparency

ScreenDisplay

Any

Lecturer Knowledge

High

Low

iBCMTechnique

Interactive Bayesian Case Model

Prototype Explanation ,Feature Importance Explanation

model-specific

ante-hoc

Visual, textual

cohort

prediction

Tabular,Image,Text

N/A

N/A

Questionnaire

Usefulness/Efficiency

N/A

A case-based reasoning system for aiding detection and classification of nosocomial infections

InNoCBR

InnoInfectionExplanationExperience

Predict patient's infection based on a clinical, laboratory, and medico administrative based data

InnoInfectionDiagnosisModel

InnHostpitalDataset

Tabular

6 tuple {S, B, V, S, C, E}

5385

InfectionDiagnosis

Multi-class Classification

Predict patient's infection based on a clinical, laboratory, and medico administrative based data

InnoHybridMethod

Rules+PARTRules+NLP(NB)

70.21%, 0.62, 55.75%, 19.18%

Accuracy, Kappa, false-positive rate (before and after modified data)

performance

Reasoning Path Explanation

model-specific

ante-hoc

text

local

prediction

Spanish NHS Doctor

Why does the system predict infection Y for patient X?

System Recommendation

Why Question

Validate system prediction

Trust/Transparency

ScreenDisplay

Any

Infection Detection/Prevention Knowledge

high

low

InnoDecisionPathTechnique

Decision tree

Reasoning Path Explanation

model-specific

anti-hoc

Text

local

prediction

Tabular

N/A

N/A

N/A

N/A

N/A

Explaining Models by Propagating Shapley Values

DeepSHAPGlobal

MortalityPredictionExplanationExperience

Predict patient mortality base on clinical, nutritional and behaviorial factors

MLPMortalityPredictionModel

NHANESDataset

Tabular

79

14407

MortalityPrediction

Binary Classification

Predict patient mortality base on clinical, nutritional and behaviorial factors

MoralityMLP

Neural Network

82.56%

Accuracy

Performance

Feature Importance Explanation

model-agnostic

post-hoc

Any

global

model

Clinician

What/How features contributed to predicting mortality Y for patient X?

Model

What/How Question

Understand how model make decisions

Transparency

ScreenDisplay

Any

Clinical Knowledge

high

low

DeepSHAPExplanationTechnique

SHAP/Game-theory

Feature Importance Explanation

model-agnostic

post-hoc

violin plot chart

global

model

Tabular

N/A

N/A

N/A

N/A

<https://github.com/lrjball/shap>

Explaining Models by Propagating Shapley Values

DeepSHAPLocal

MortalityExplanationExperience
Predict patient mortality base on clinical, nutritional and behaviorial factors
MLPMortalityPredictionModel
NHANESDataset
Tabular
79
14407
MortalityPrediction
Binary Classification
Predict patient mortality base on clinical, nutritional and behaviorial factors
MoralityMLP
Neural Network
82.56%
Accuracy
Performance

Feature Importance Explanation
model-agnostic
post-hoc
Any
local
prediction
Clinician

What/How features contributed to predicting mortality Y for patient X?
System Recommendation
What/How Question
Understand why model made a decision
Transparency/Education
ScreenDisplay
Any

Clinical Knowledge
high
low

DeepSHAPExplanationTechnique
SHAP/Game-theory
Feature Importance Explanation
model-agnostic
post-hoc
bar chart
local
prediction
Tabular
N/A
N/A
N/A
N/A
<https://github.com/lrjball/shap>