# Trustworthy AI: Explainability

Destiny Pounds

M.S. Student, Biomedical Data Science

Advised by Vibhuti Gupta, Ph.D.

Assistant Professor, Computer Science and Data Science

School of Applied Computational Sciences
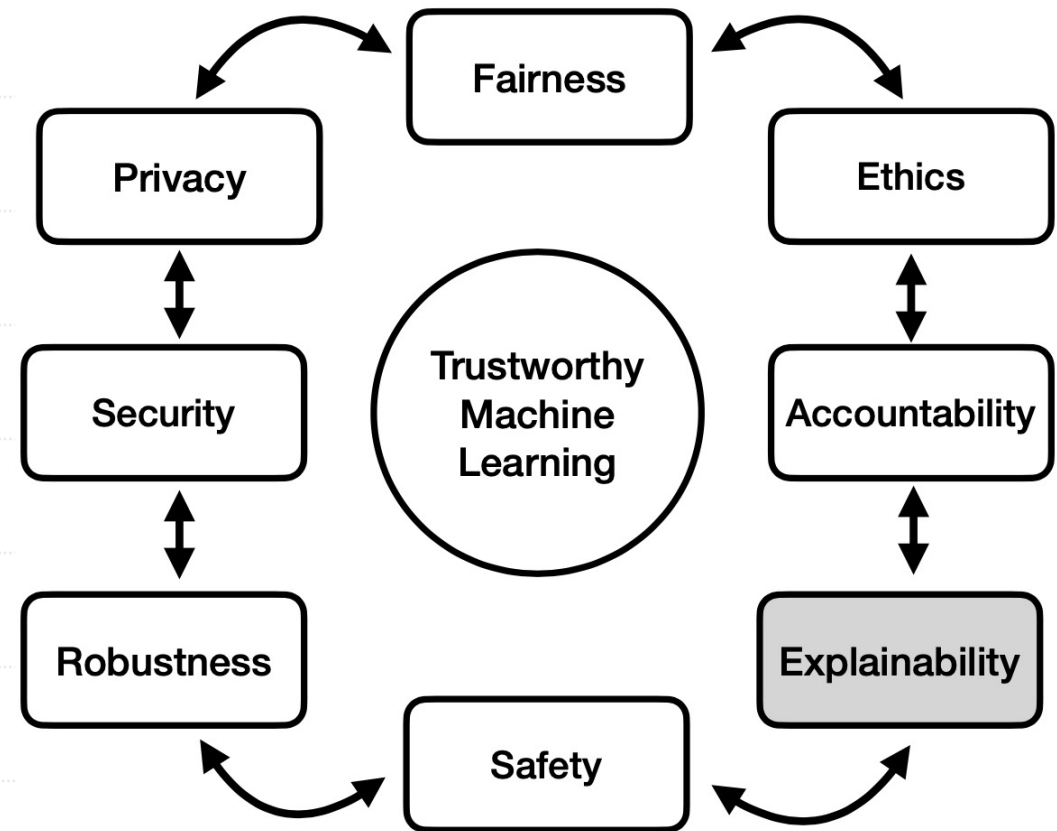
Meharry Medical College

# Overview

- Defining Explainability
- Model Understanding
  - Examples
  - Benefits

  Approaches to Model Understanding
  - Interpretable Models
  - Post-hoc Explainability
    - Local Explanations
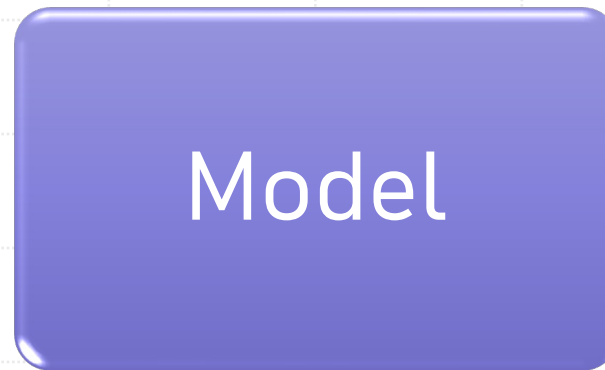    - Global Explanations

# Explainable Artificial Intelligence (XAI)

- Explainability of an AI model describes the extent to which human–users can comprehend and trust the results and output created by the model.

# Overview of Predictive Modeling Process

Input
(Data)

Model

Output
(Prediction)

Explainable AI requires model understanding.

# Example: Why Model Understanding?

Input



Predictive Model 🐛

Prediction = Siberian Husky

# Example: Why Model Understanding?

Input

Model Understanding

This model is relying on incorrect features to make this prediction!! Let me fix the model

Predictive Model

Prediction = Siberian Husky

Adapted from Farnadi, G. (2022). *Trustworthy Machine Learning* [Slides]. HEC Montréal

Adapted from Farnadi, G. (2022). *Trustworthy Machine Learning* [Slides]. HEC Montréal

# Example: Why Model Understanding?

Defendant Details

Predictive Model 🐛 ——— Prediction = Risky to Release

# Example: Why Model Understanding?



Defendant Detai...

This prediction is biased. Race and ...eing ...the ...!

**Model understanding facilitates bias detection.**

Gender

Predictive Model

Prediction = Risky to Release

# Example: Why Model Understanding?

Loan Applicant Details



Predictive Model

Prediction = Denied Loan

Loan Applicant

Adapted from Farnadi, G. (2022). *Trustworthy Machine Learning* [Slides]. HEC Montréal

# Example: Why Model Understanding?



Model understanding helps provide recourse to individuals who are adversely affected by model predictions.

Loan Applicant Details

FILE

Predictive Model

Prediction = Denied Loan

Loan Applicant

I have some means... Let me ... my ... pay ... ne.

... to get a loan

# Example: Why Model Understanding?



Patient Data

| | |
|---|---|
| 25, Female, Cold | |
| 32, Male, No | |
| 31, Male, Cough | |
| . | |
| . | |
| . | |
| . | |

Model Understanding

Predictions

| | |
|---|---|
| Healthy | |
| Sick | |
| Sick | |
| . | |
| . | |
| Healthy | |
| Healthy | |
| Sick | |

Predictive Model

# Example: Why Model Understanding?



Patient Data

25, Fema|
32, Male,
31, Male,
.
.
.
.

Model Understanding

This model is using irrelevant features when ... on female ... I should ... edictions ...oup.

Predictive Model

Sick
Sick
.
.
Healthy
Healthy
Sick

**Model understanding helps assess if and when to trust model predictions when making decisions.**

# Example: Why Model Understanding?



Patient Data

25, Femal[e]
32, Male,
31, Male,
.
.
.

Predictive Model

Model Understanding

Healthy
Sick
Sick
.
.
Healthy
Healthy
Sick

Model understanding allows us to vet models to determine if they are suitable for deployment in real world.

This model is using irrelevant features when ... female ... on. This ... proved!

AUTHORITY

# Summary: Benefits of Model Understanding

- Debugging

- Bias Detection

- Recourse

- Assess Trustworthiness of Predictions

- Vet Models for Deployment

# Achieving Model Understanding

Approach #1:

Build inherently interpretable (white box) models.
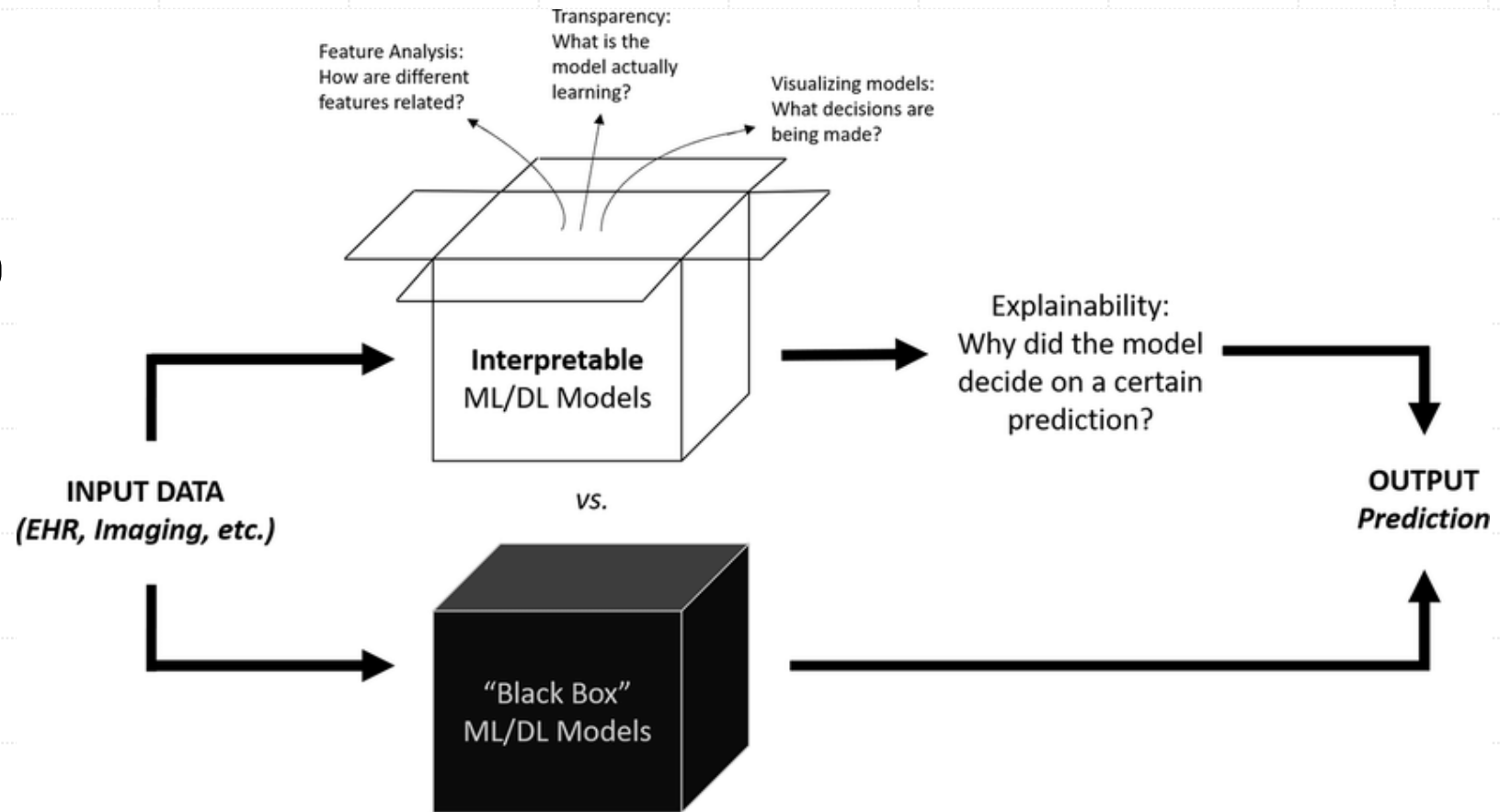
Approach #2:

Explain pre-built (black-box) models in a post-hoc manner.



Feature Analysis: How are different features related?

Transparency: What is the model actually learning?

Visualizing models: What decisions are being made?

**Interpretable** ML/DL Models

Explainability: Why did the model decide on a certain prediction?

INPUT DATA *(EHR, Imaging, etc.)*

vs.

"Black Box" ML/DL Models

OUTPUT *Prediction*

# Inherently Interpretable Models

■ These models are interpretable by design. Understanding of the model is clear before receiving results.



## Examples

▪ Linear Models

▪ Shallow Decision Tree

▪ Rule Based Models

▪ Risk Scores



| The HEART Score for Chest Pain Patients in the ED | | |
|---|---|---|
| **History** | • Highly Suspicious<br>• Moderately Suspicious<br>• Slightly or Non-Suspicious | • 2 points<br>• 1 point<br>• 0 points |
| **ECG** | • Significant ST-Depression<br>• Nonspecific Repolarization<br>• Normal | • 2 points<br>• 1 point<br>• 0 points |
| **Age** | • ≥ 65 years<br>• > 45 - < 65 years<br>• ≤ 45 years | • 2 points<br>• 1 point<br>• 0 points |
| **Risk Factors** | • ≥ 3 Risk Factors or History of CAD<br>• 1 or 2 Risk Factors<br>• No Risk Factors | • 2 points<br>• 1 point<br>• 0 points |
| **Troponin** | • ≥ 3 x Normal Limit<br>• > 1 - < 3 x Normal Limit<br>• ≤ Normal Limit | • 2 points<br>• 1 point<br>• 0 points |

**Risk Factors:** DM, current or recent (<one month) smoker, HTN, HLP, family history of CAD, & obesity

**Score 0 – 3:** 2.5% MACE over next 6 weeks → Discharge Home
**Score 4 – 6:** 20.3% MACE over next 6 weeks → Admit for Clinical Observation
**Score 7 – 10:** 72.7% MACE over next 6 weeks → Early Invasive Strategies

# Post-hoc Explainability
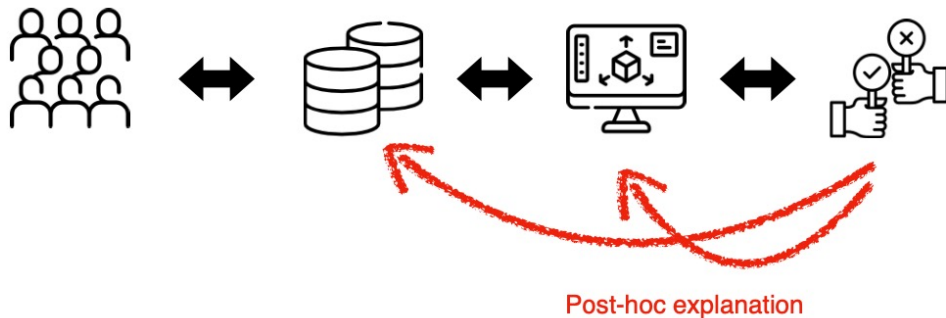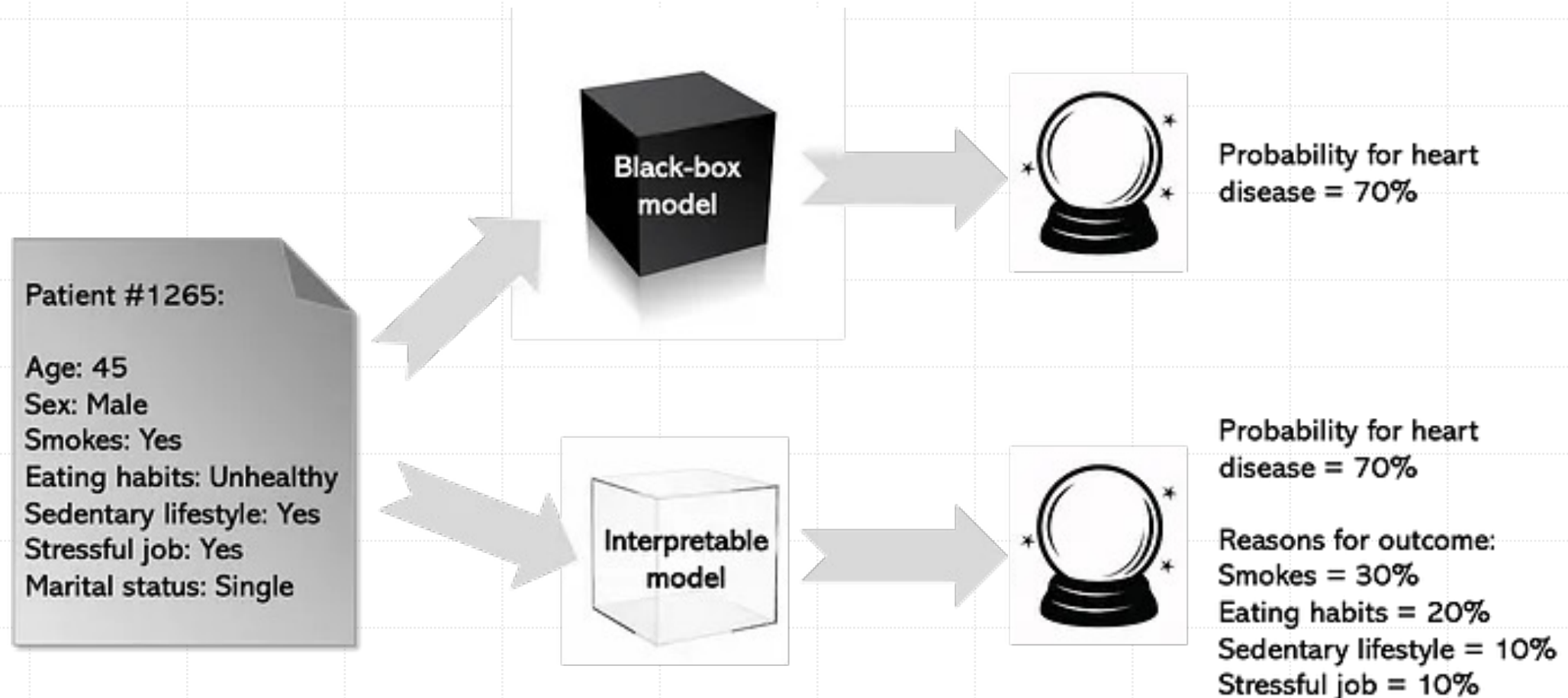
Building an inherently interpretable model is not always possible and we are left with a black box. In that case, models can be interpreted post-hoc.

Post-hoc explanation occurs after execution of the model.

Post-hoc explanation

# White Box vs. Black Box Example



Patient #1265:

Age: 45
Sex: Male
Smokes: Yes
Eating habits: Unhealthy
Sedentary lifestyle: Yes
Stressful job: Yes
Marital status: Single

Black-box model

Probability for heart disease = 70%

Interpretable model

Probability for heart disease = 70%

Reasons for outcome:
Smokes = 30%
Eating habits = 20%
Sedentary lifestyle = 10%
Stressful job = 10%

# Approaches for Post-hoc Explainability

## Local Explanations

Local explanations focus on the data and provide explanations for individual outcomes. Therefore, they provide trust for model outcomes.

## Global Explanations

Global explanations focus on the model and provide an understanding of the decision process. Therefore, they provide a sense of understanding to how the model works.

# Thank You

**Please send us your questions at:**
**vgupta@mmc.edu and**
**dpounds24@email.mmc.edu**