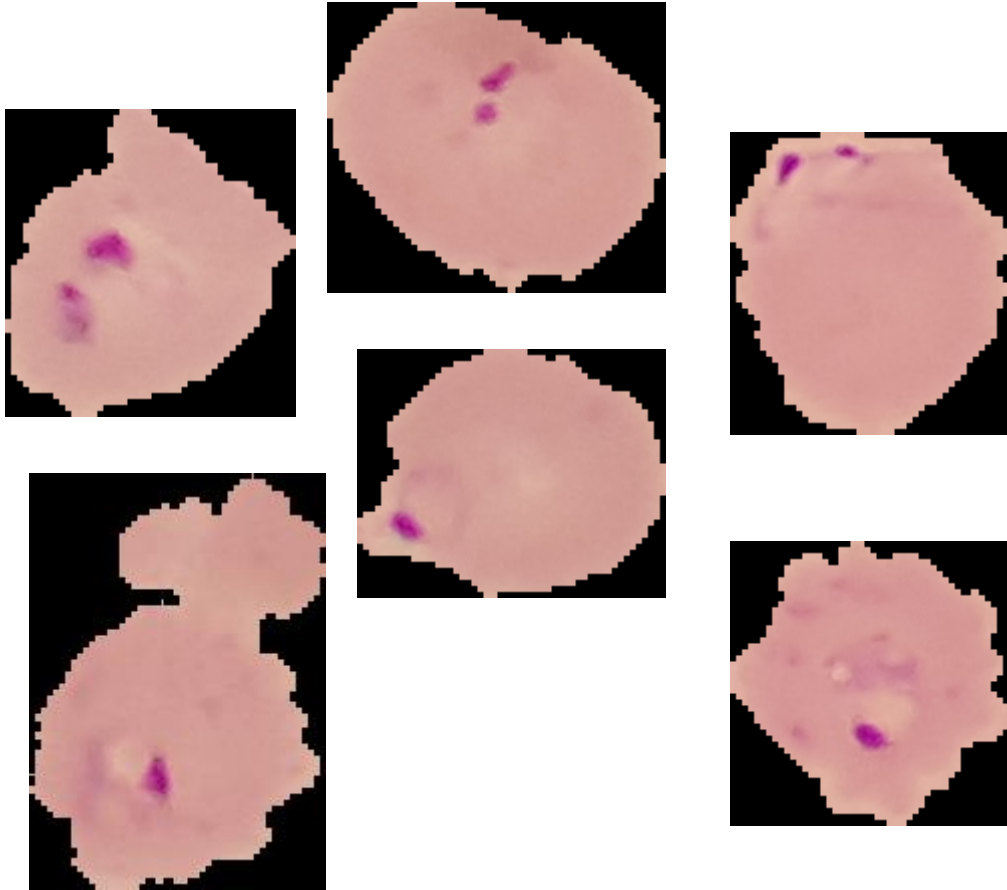# Interpretable AI

How I Learned to Stop Worrying
and Trust AI

**Ajay Thampi (PhD)**
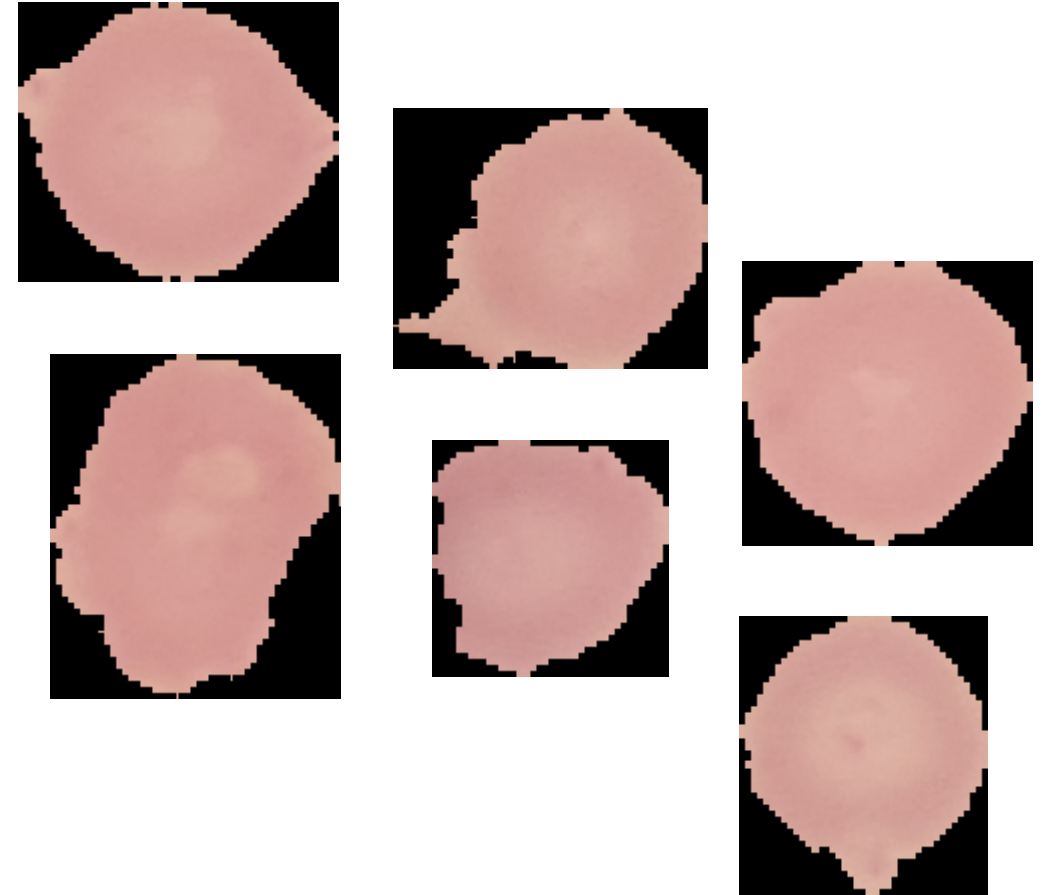Lead Data Scientist

Microsoft

# Learning – Detecting Malaria
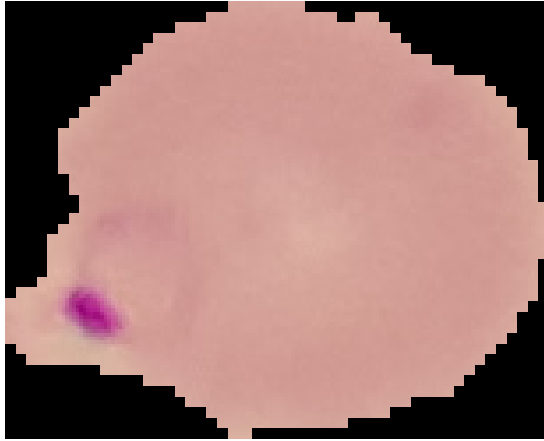


**INFECTED**

**UNINFECTED**
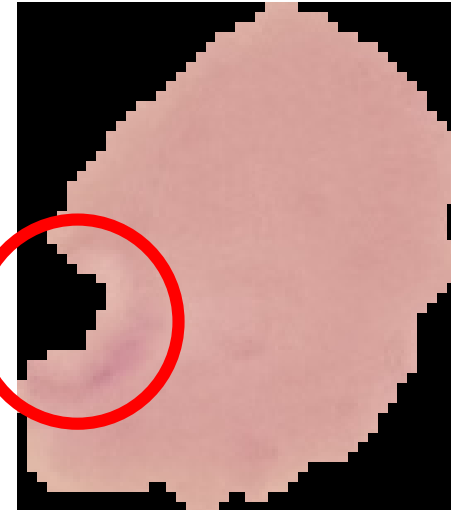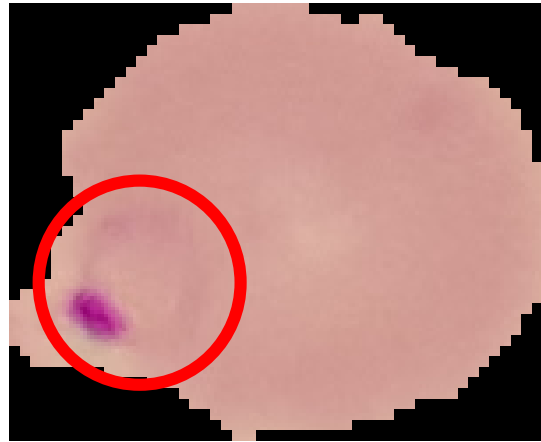
**UNINFECTED**

**INFECTED**
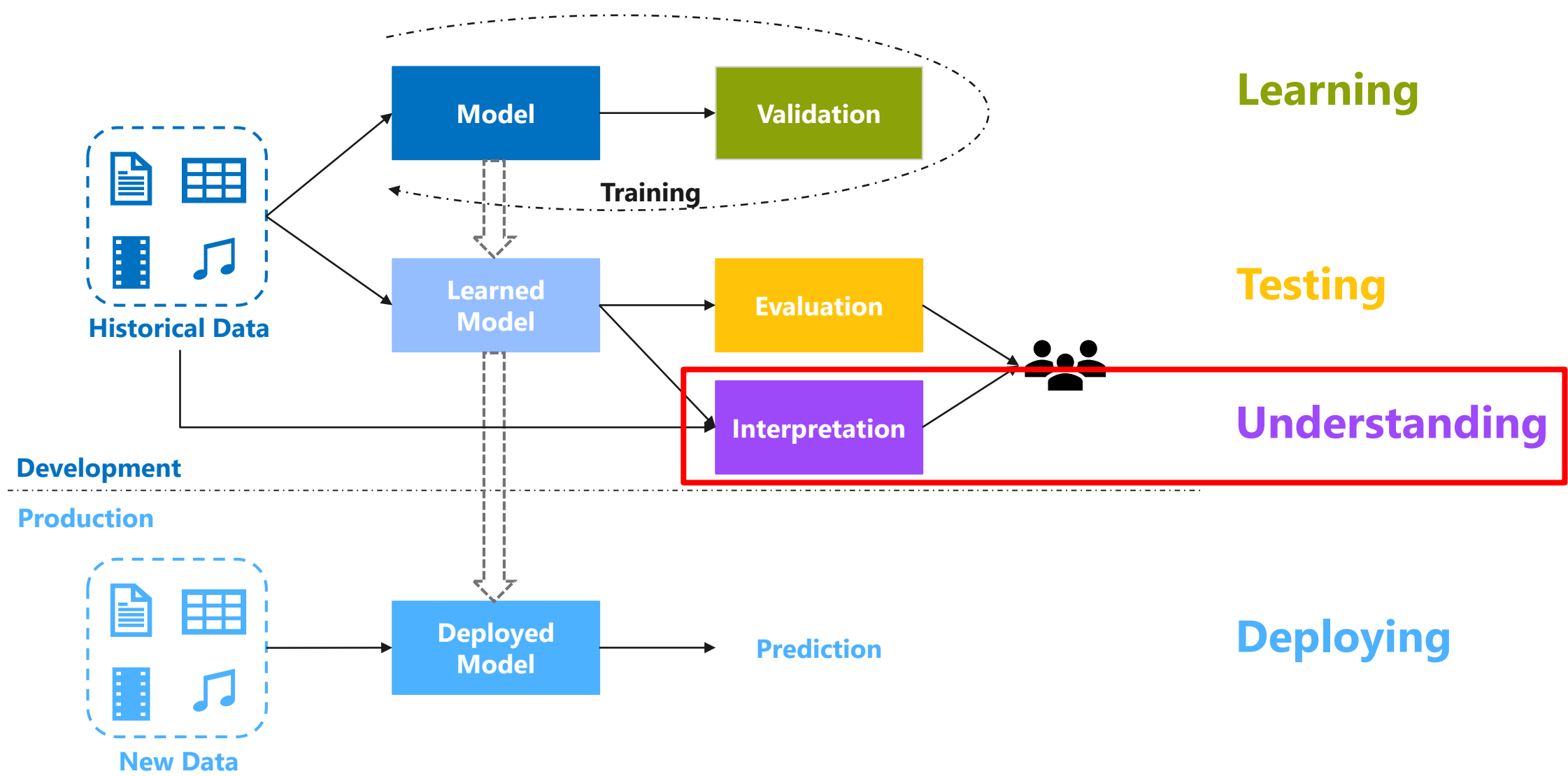
**INFECTED**

# Understanding



**UNINFECTED**

**INFECTED**

# Modelling and Deployment

# Transparent Models

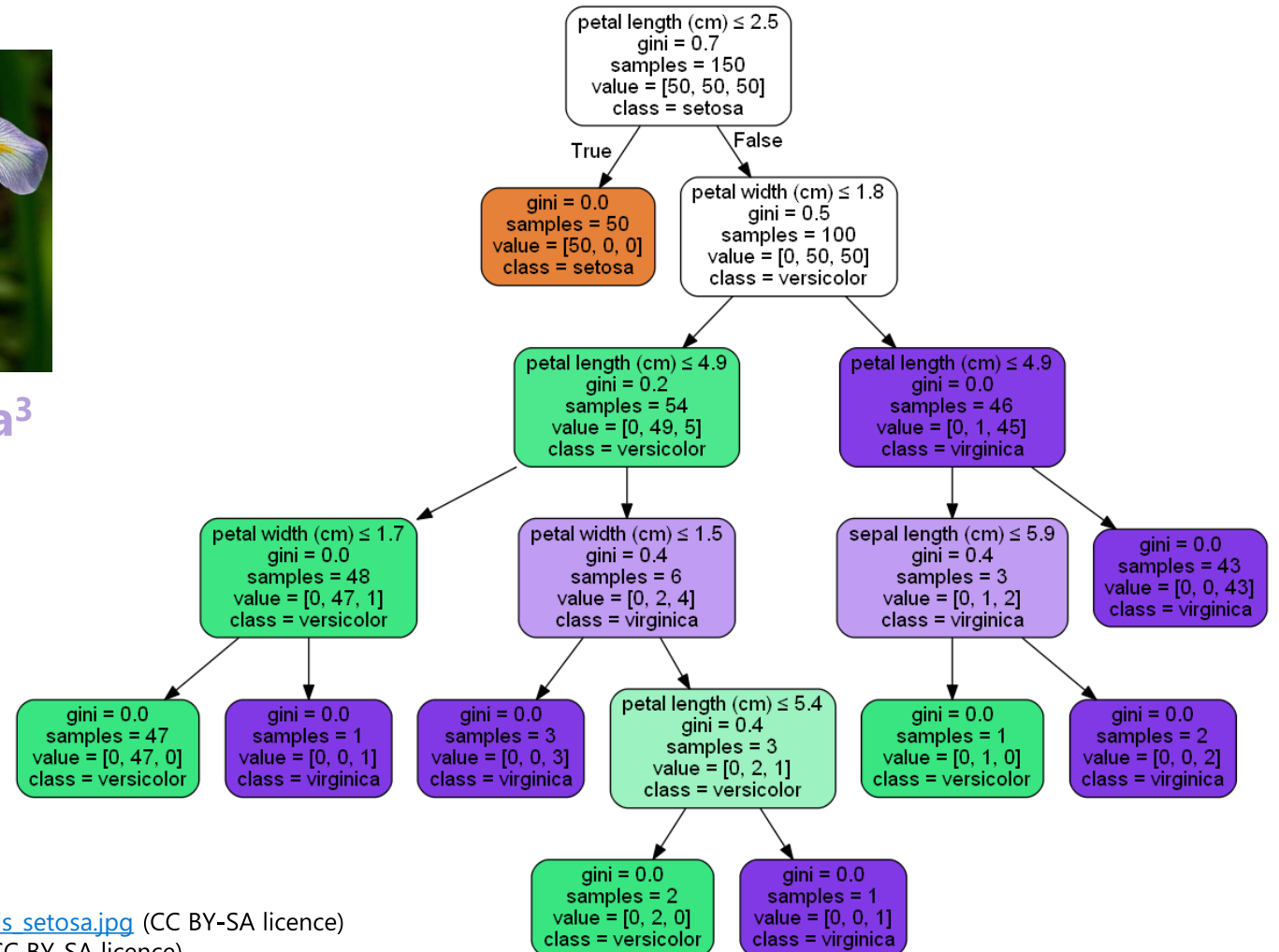# Decision Trees - Simple

## Iris Flower Dataset



**Setosa**[1]     **Versicolor**[2]     **Virginica**[3]

**Features:**
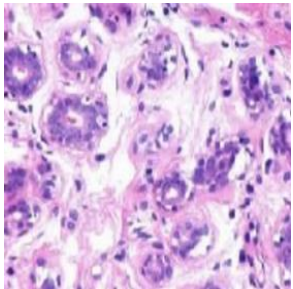- Petal Length
- Petal Width
- Sepal Length
- Sepal Width



petal length (cm) ≤ 2.5
gini = 0.7
samples = 150
value = [50, 50, 50]
class = setosa

True — False

gini = 0.0
samples = 50
value = [50, 0, 0]
class = setosa

petal width (cm) ≤ 1.8
gini = 0.5
samples = 100
value = [0, 50, 50]
class = versicolor

petal length (cm) ≤ 4.9
gini = 0.2
samples = 54
value = [0, 49, 5]
class = versicolor

petal length (cm) ≤ 4.9
gini = 0.0
samples = 46
value = [0, 1, 45]
class = virginica

petal width (cm) ≤ 1.7
gini = 0.0
samples = 48
value = [0, 47, 1]
class = versicolor

petal width (cm) ≤ 1.5
gini = 0.4
samples = 6
value = [0, 2, 4]
class = virginica

sepal length (cm) ≤ 5.9
gini = 0.4
samples = 3
value = [0, 1, 2]
class = virginica

gini = 0.0
samples = 43
value = [0, 0, 43]
class = virginica

gini = 0.0
samples = 47
value = [0, 47, 0]
class = versicolor

gini = 0.0
samples = 1
value = [0, 0, 1]
class = virginica

gini = 0.0
samples = 3
value = [0, 0, 3]
class = virginica

petal length (cm) ≤ 5.4
gini = 0.4
samples = 3
value = [0, 2, 1]
class = versicolor

gini = 0.0
samples = 1
value = [0, 1, 0]
class = versicolor

gini = 0.0
samples = 2
value = [0, 0, 2]
class = virginica

gini = 0.0
samples = 2
value = [0, 2, 0]
class = versicolor

gini = 0.0
samples = 1
value = [0, 0, 1]
class = virginica

# Decision Trees – More Complex

## Breast Cancer Dataset



**Benign[1]**   **Malignant[1]**
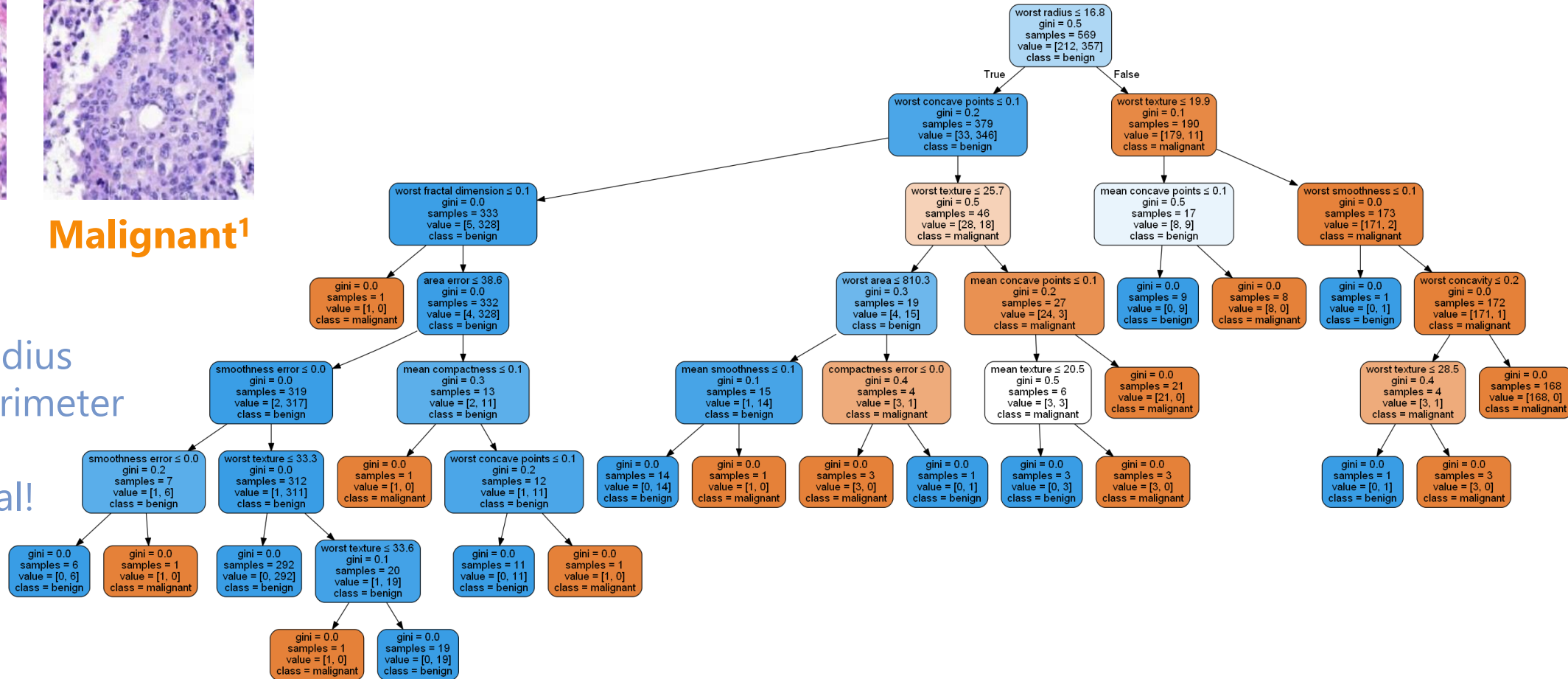
**Features:**
- Mean Radius
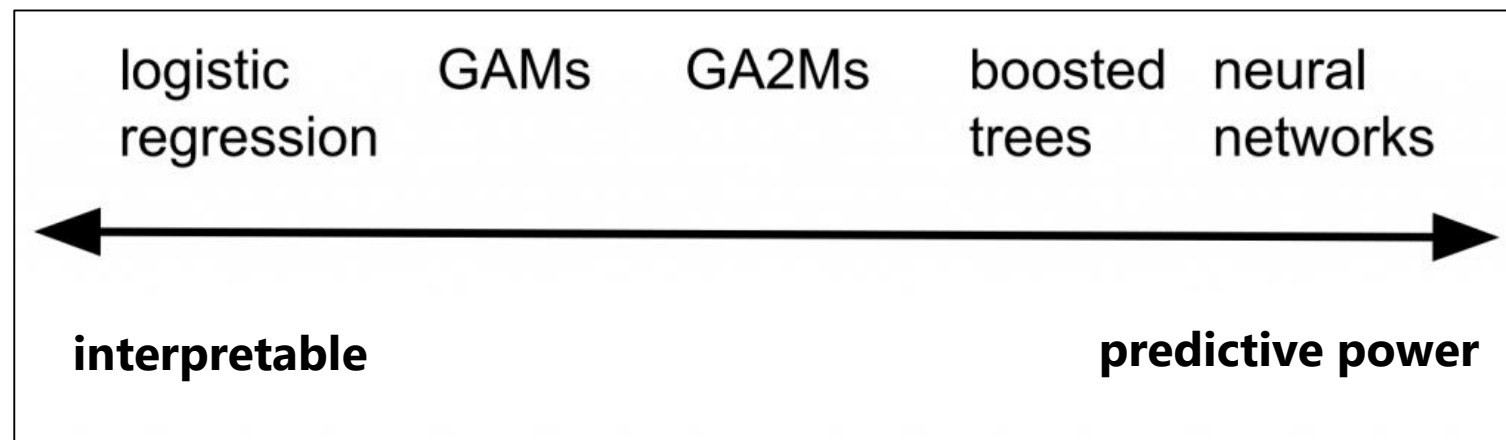- Mean Perimeter
- ...
- 30 in total!

# Other Transparent Models

- Linear Regression
- Logistic Regression
- Generalized Additive Models (GAMs)
  - GA2Ms

$$g(E[y]) = \sum f_i(x_i)$$

$$g(E[y]) = \sum f_i(x_i) + \sum f_{ij}(x_i, x_j)$$

| logistic regression | GAMs | GA2Ms | boosted trees | neural networks |
| --- | --- | --- | --- | --- |

interpretable ⟵⟶ predictive power

# Data & Features

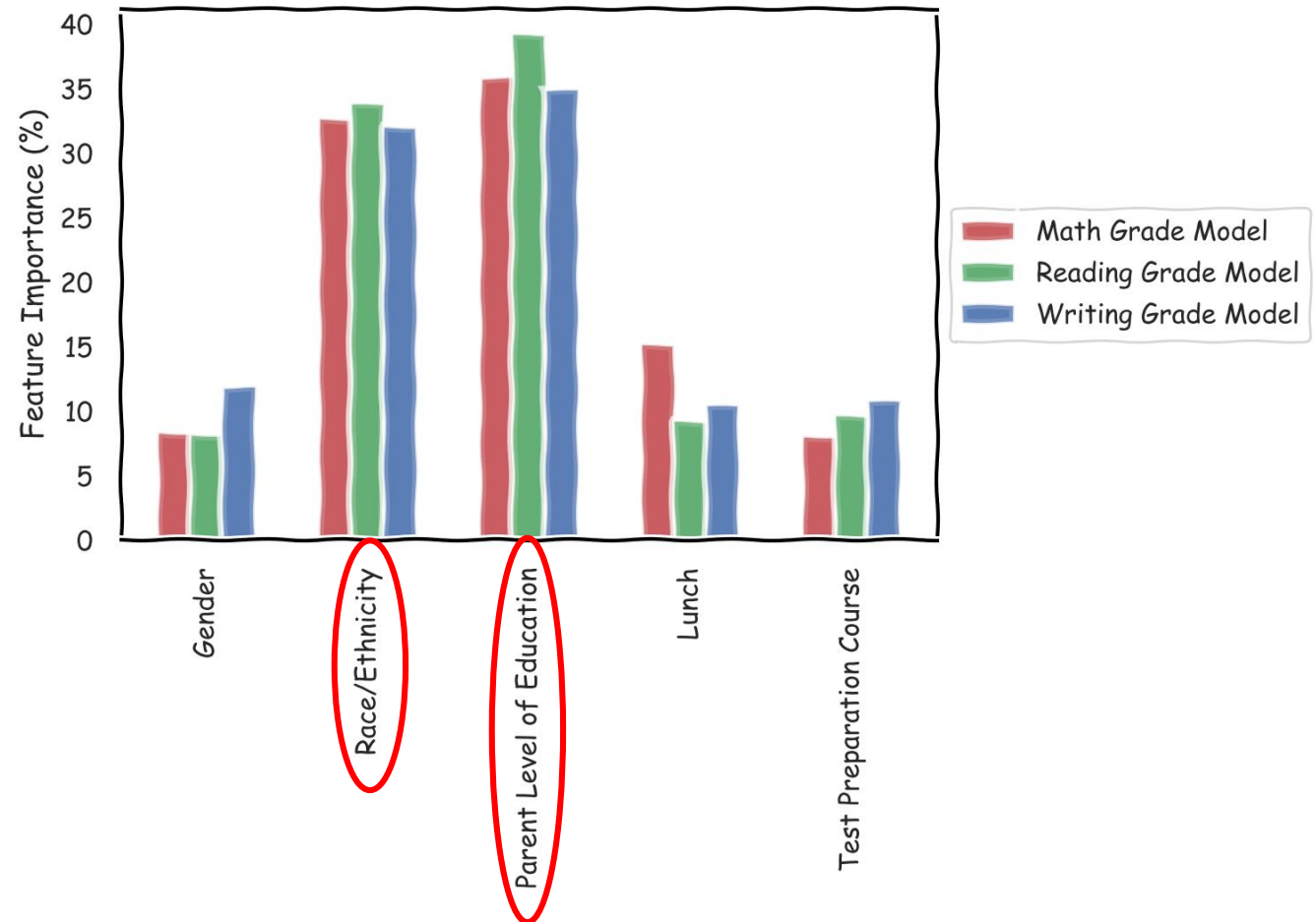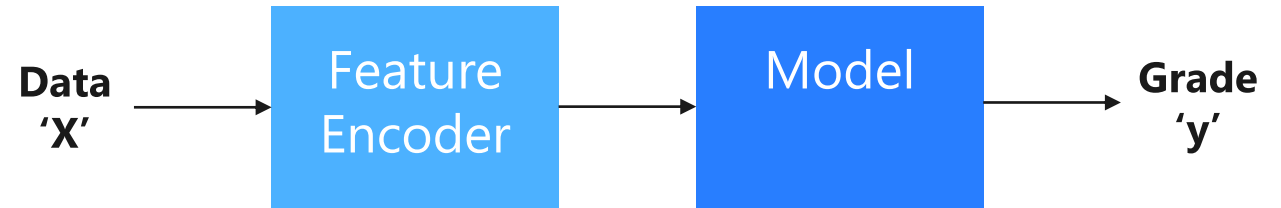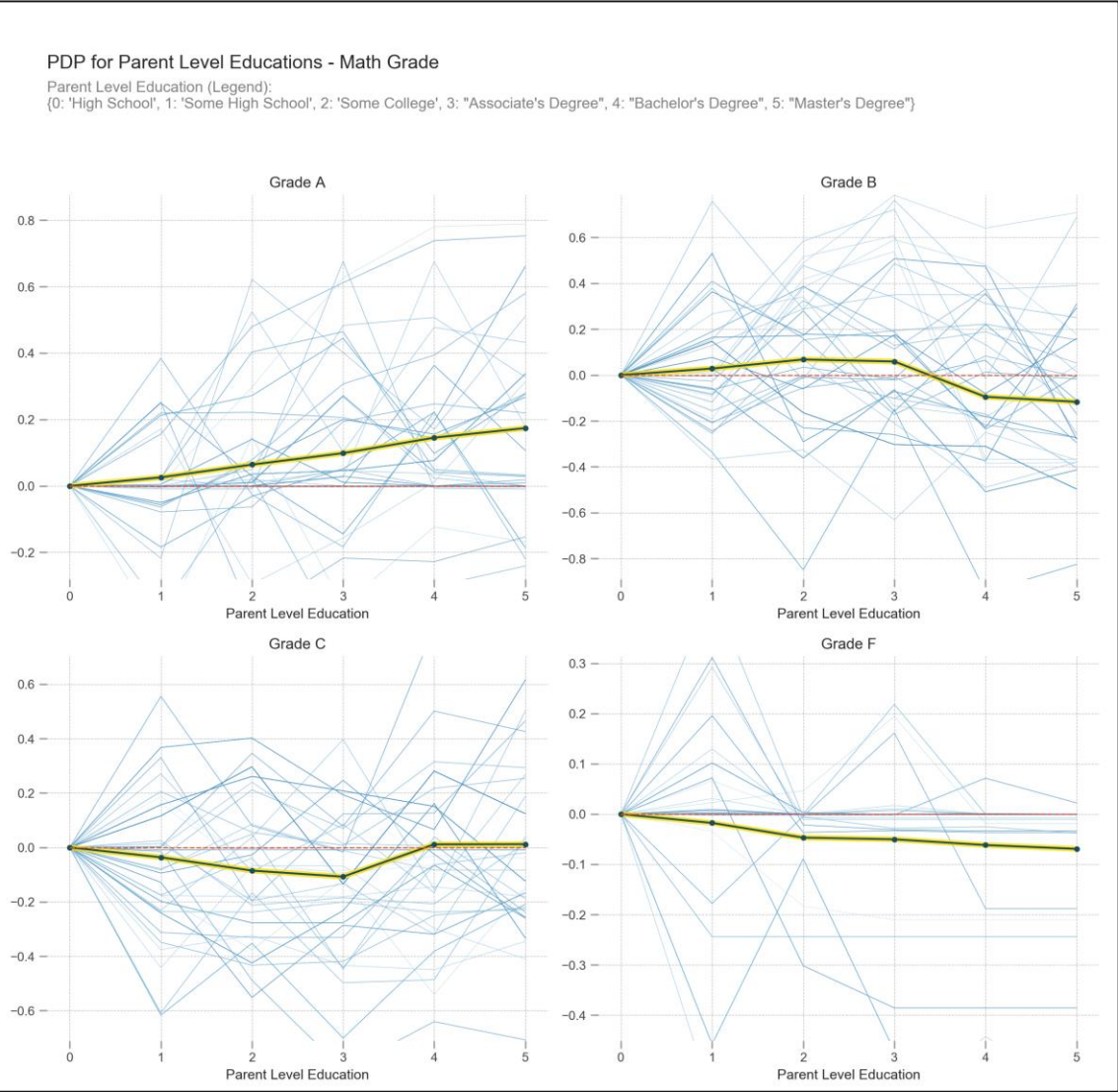Going Beyond Feature Importance

# Feature Importance

Problem: Predict high-school student grades for math, reading and writing

Features:

- Gender

- Race/Ethnicity

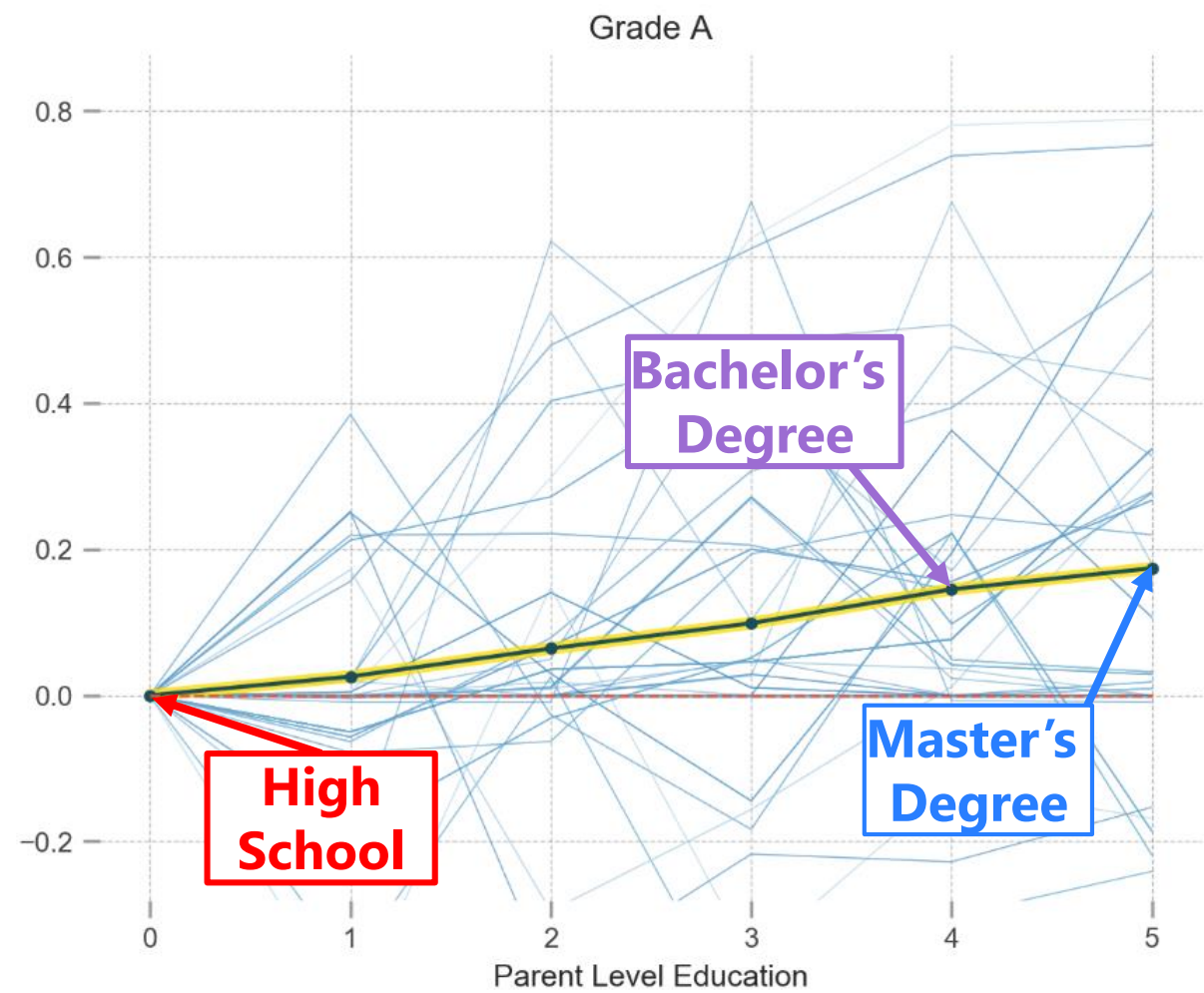- Parent Level of Education

- Lunch

- Test Preparation

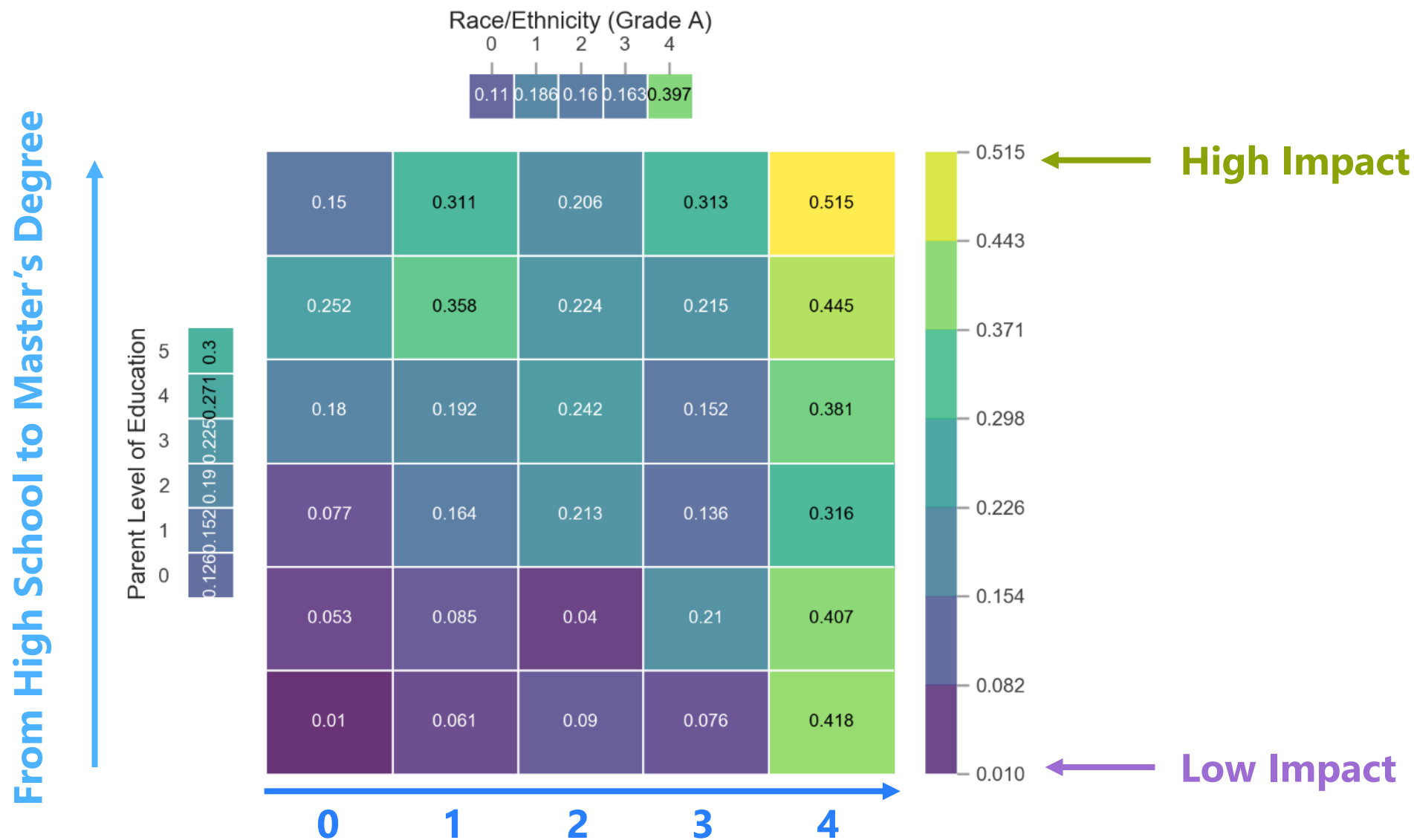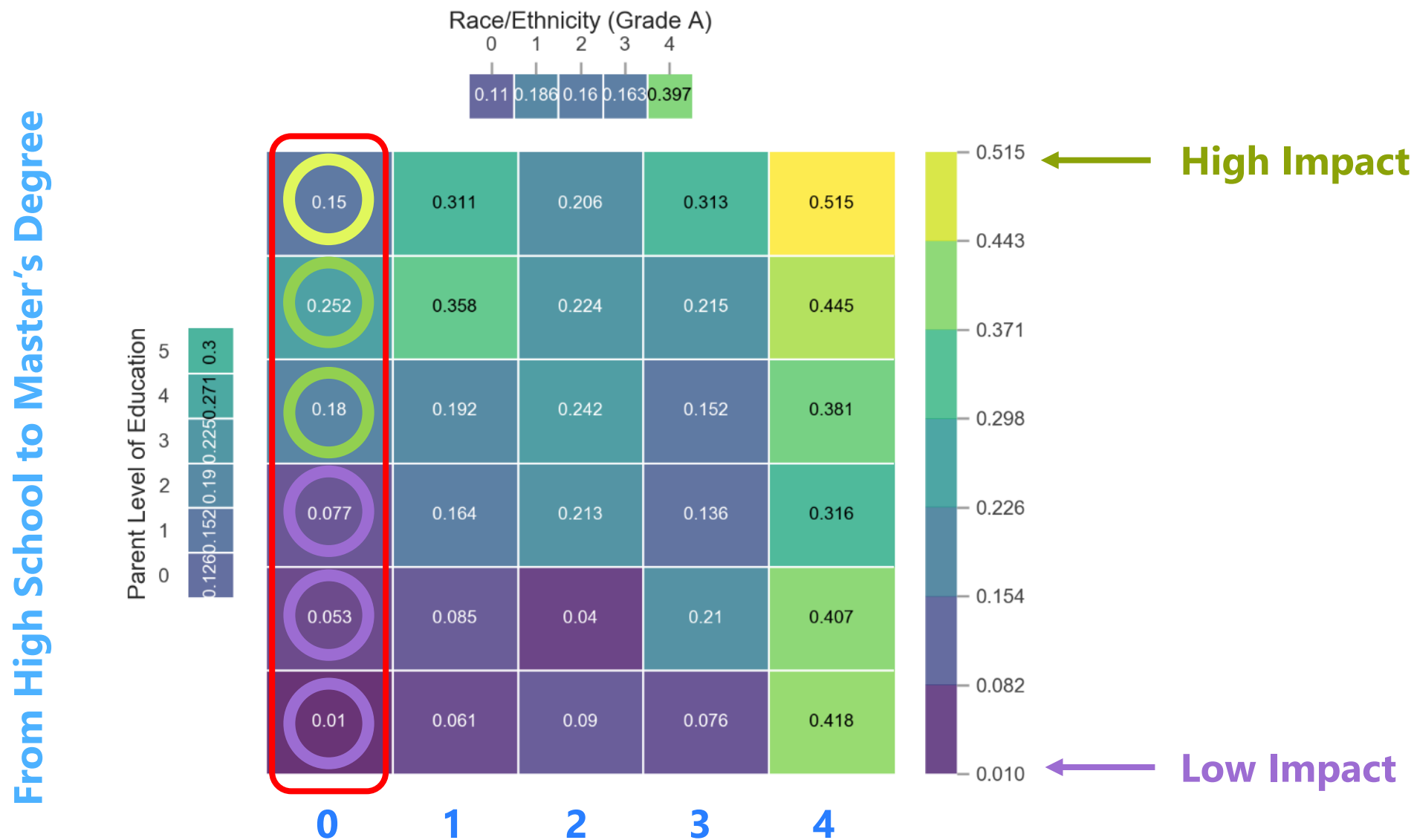# Partial Dependency Plots (PDPs)

# Parent Level Education

# Feature Interactions

# Feature Interactions

# Feature Interactions

# PDP in Python

```
from pdpbox import pdp


pdp_race = pdp.pdp_isolate(model=math_model,
                           dataset=df,
                           model_features=features,
                           feature='race')
```



```
pdp_race_parent = pdp.pdp_interact(model=math_model,
                                   dataset=df,
                                   model_features=features,
                                   features=['race', 'parent'])
```

# Black-Box Models
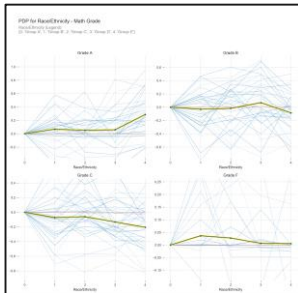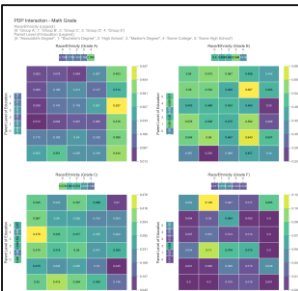
Post-Hoc Explanations

# Post-Hoc Explanations



Learned representation that is locally faithful but not globally

Pick an instance to explain

Complex decision function - hard to explain

**LIME =** **L**ocal **I**nterpretable **M**odel-agnostic **E**xplanations

**2016**

Data x

Complex Model f

Prediction f(x)

Linear Surrogate Model g(x') ≈ f(x)

Explanation ∅

$$g(x') = \emptyset_0 + \sum_{i=1}^{M} \emptyset_i x'_i \approx f(x)$$

**Explanation Parameters**

**SHAP =** **S**hapley **A**dditive ex**P**lanations

**2017**

# SHAP Tree Ensemble Explainer

## Breast Cancer Detection

**Data**
- Mean Radius
- Mean Perimeter
- Worst Concave Points
- Mean Concave Points
- …

Random Forest → **Malignant Or Benign**

SHAP TreeExplainer

# SHAP Deep Learning Explainer

# SHAP Deep Learning Explainer – Explained!

# Mitigating Bias

Algorithmic Debiasing

# Adversarial Debiasing

Input **X** → **Predictor** (Weights **W**) —Prediction **Ŷ**→ **Adversary** (Weights **U**) → Protected Variable **Z**

- Demographic Parity: $\hat{Y} \perp Z$
- Equality of Odds: $\hat{Y} \perp Z \mid Y$
- Equality of Opportunity: $\hat{Y} \perp Z \mid Y = y$

*Source*: Brian Zhang, Blake Lemoine and Margaret Mitchell. Mitigating Unwanted Biases with Adversarial Learning. AAAI Conference on AI, Ethics and Society, 2018.

# Adversarial Debiasing Demo

Word Embedding Analogy Task

## He : She :: Doctor : ?

| Word | Score |
|------|-------|
| Nurse | 0.62 |
| Her | 0.60 |
| Woman | 0.58 |
| Mother | 0.57 |
| Doctors | 0.55 |
| Physician | 0.53 |
| Pregnant | 0.51 |

# Summary

- Include model understanding in your data science process
- Be mindful of your audience – interpretability means different things to different people
- Apply interpretability techniques (like PDPs, LIME, SHAP, etc.) to improve model understanding
- Build fair models by mitigating bias

# Additional Resources (1/2)

- Source code of demos: https://github.com/thampiman/interpretability

- Blog post on interpretability: https://towardsdatascience.com/interpretable-ai-or-how-i-learned-to-stop-worrying-and-trust-ai-e61f9e8ee2c2

- Saliency Maps: https://distill.pub/2018/building-blocks/

- Representational Learning: https://www.cl.uni-heidelberg.de/courses/ws14/deepl/BengioETAL12.pdf

- t-SNE: https://lvdmaaten.github.io/tsne/

- PDP Box: https://github.com/SauceCat/PDPbox

- LIME: https://arxiv.org/pdf/1602.04938.pdf

- Kernel SHAP: http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions

# Additional Resources (2/2)

- SandDance: https://www.microsoft.com/en-us/research/project/sanddance/

- GAMut: https://www.microsoft.com/en-us/research/uploads/prod/2019/01/19_gamut_chi.pdf

- Datasheets for Datasets: https://arxiv.org/pdf/1803.09010.pdf

- Challenges for Transparency: https://arxiv.org/pdf/1708.01870.pdf

- Synthetic Data (MSR): https://arxiv.org/pdf/1810.00471.pdf

- Counterfactual Explanations: https://arxiv.org/abs/1711.00399

- Noise Audit: https://hbr.org/2016/10/noise

- Interpretable ML: https://christophm.github.io/interpretable-ml-book/limo.html

# Thank You

Q&A

@thampiman