# A Bayesian's approach to predict the probability of getting diabetes
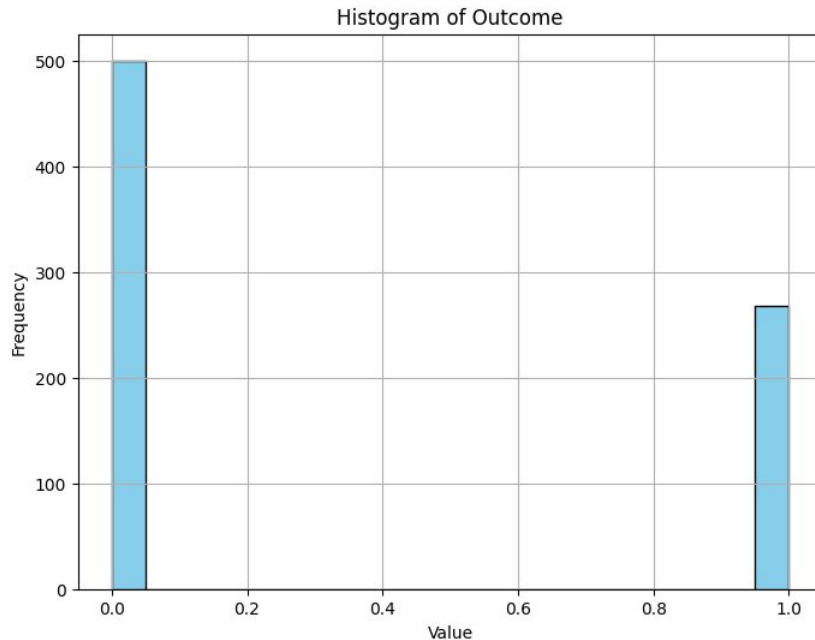
Saksham Malik, Alex Yang, Feifan Liu

# Dependent variable  - Outcome

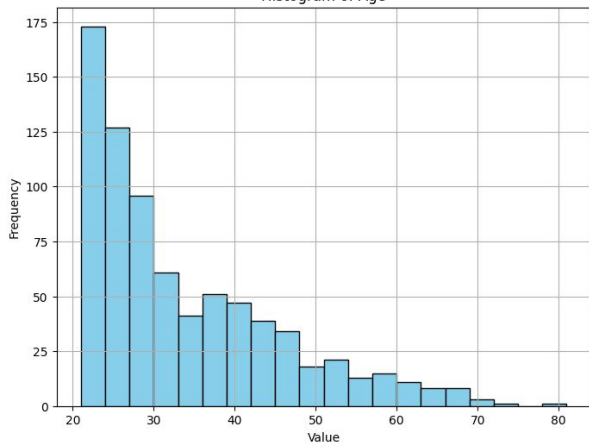Outcome: Bernoulli RV
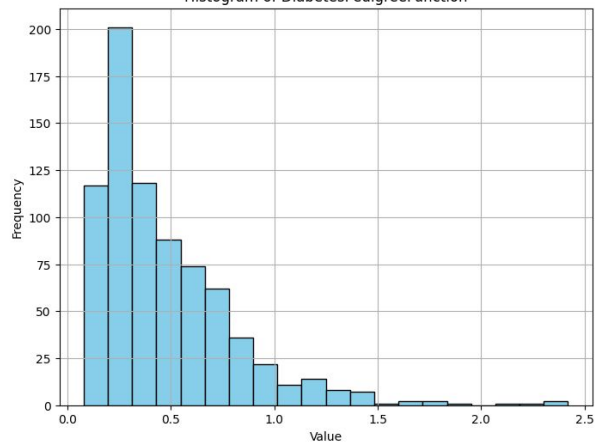
whether the person have a diabetes or not.

P ≈ 0.34



Histogram of Outcome
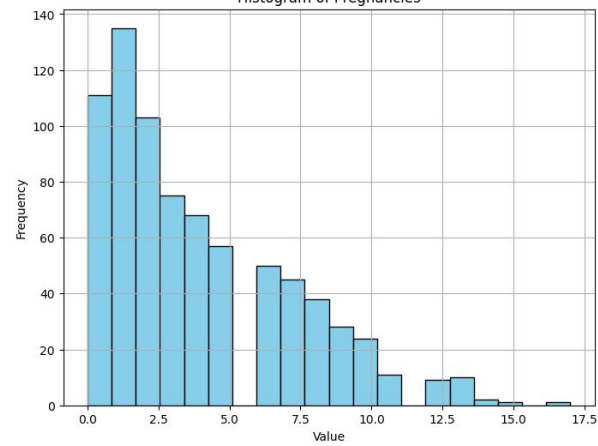
# Features
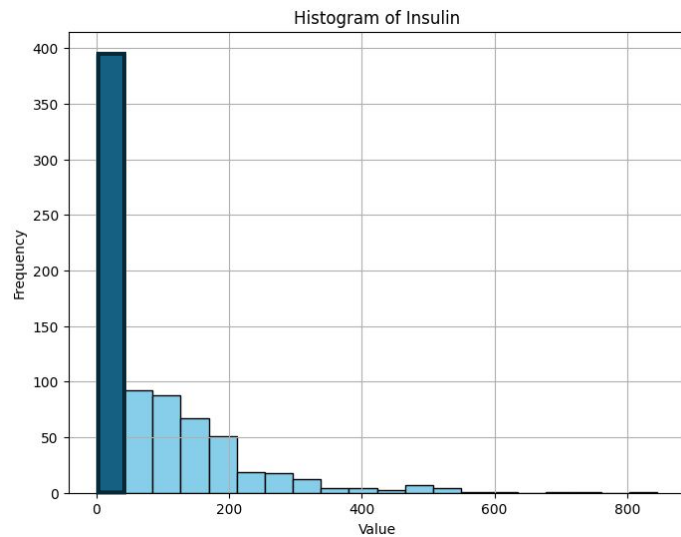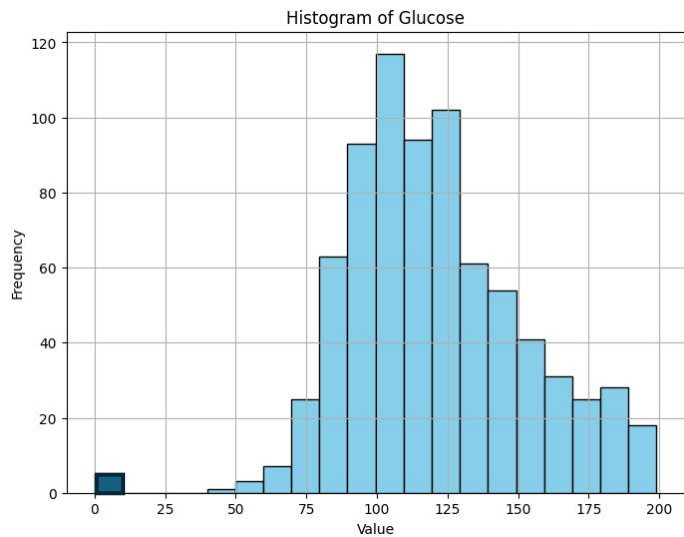
# Features

# Features

# Data Imputation

Multiple Imputation by Chained Equations (MICE) algorithm

1. Replace missing values with the column mean
2. Fix a column with missing values and fit a regression model on the remaining columns to predict the fixed column values
3. Predict missing value (which was set to mean initially) based on other columns and replace it.
4. Repeat for other columns
5. Repeat steps 2-4 for k iterations (often use k=5 in practice)

Assumption: Data is missing at random (MAR) or missing completely at random (MCAR). Missing not at random (MNAR) can only be imputed properly via Bayesian methods (equivalent algorithm is BICE)

# Bayesian Additive Regressive Trees (BART) model

- Non-parametric model regression approach
- Performs bayesian model averaging (ensembling) on a large number of shallow(low depth) and sparse(low number of splits per level) decision trees.
- The hyperparameters $\alpha$ and $\beta$ parametrize the probability that a node at depth $d(=0,1,2,...)$ is non-terminal, given by $\alpha(1+d)^{-\beta}$.
- The default values $\alpha=0.95$ and $\beta=2$ ensure the trees are shallow.

Example of decision tree

# BART algorithm

1. Recursive partitioning:
   a. Iterate through each feature (X) at each step.
   b. Choose the split point that minimizes an impurity measure (like variance) for the target variable (y). This creates two child nodes.
   c. Repeat splitting on the child nodes until a stopping criteria is met (e.g., minimum number of data points in a node). This creates a single decision tree.
2. Assign a constant value (average target variable) to each terminal node (leaf) of the tree.
3. Repeat Steps 1-2 *m* times (often use 50, 100, 200 in practice)
4. Apply regularization based on prior node probabilities
5. Prediction for new data point:
   a. Obtain a prediction from each tree (the value assigned to the terminal node where the data point lands).
   b. Use BMA on the predictions from all trees in the ensemble to get the final BART prediction for the new data point.
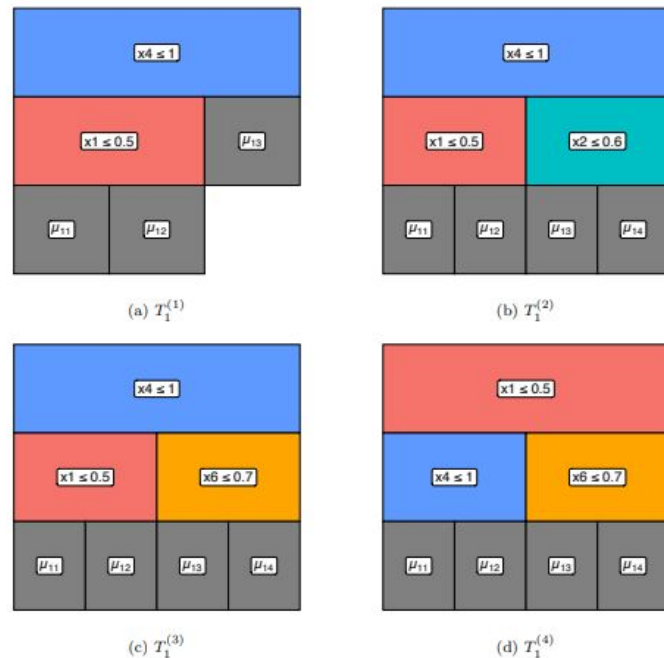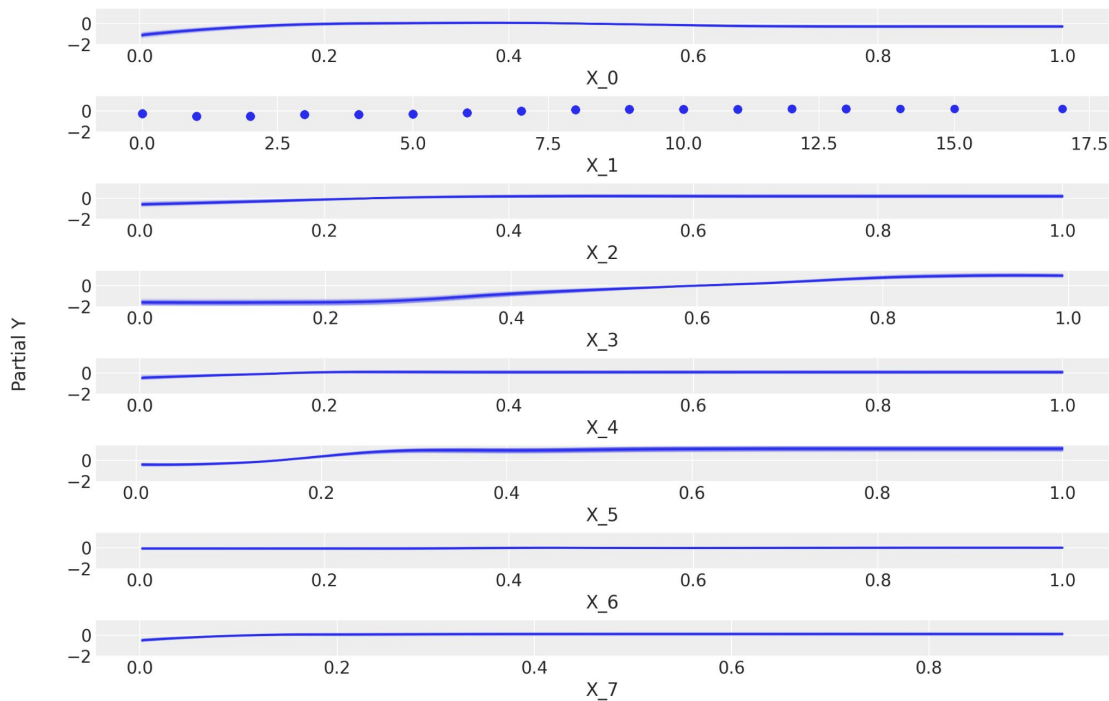


Figure1. - Inglis, A., Parnell, A., & Hurley, C. (2022). Visualizations for Bayesian Additive Regression Trees [arXiv:2208.08966]
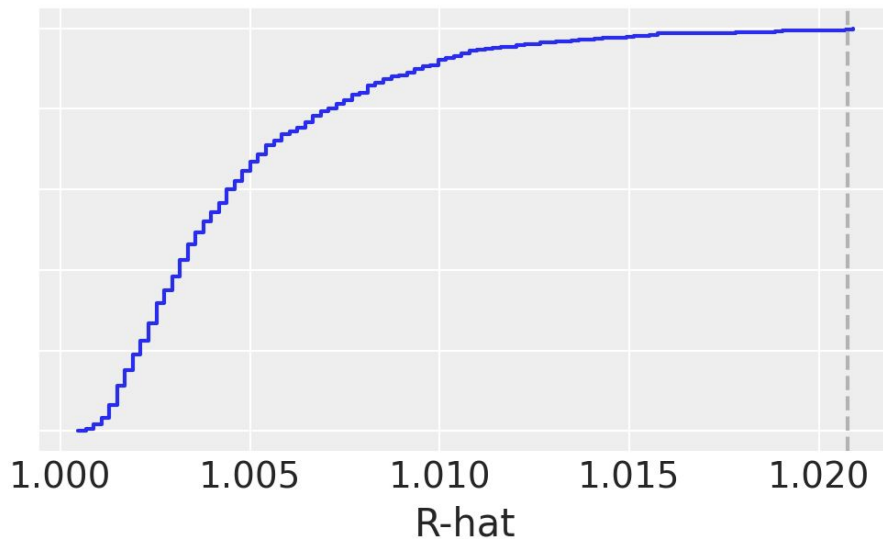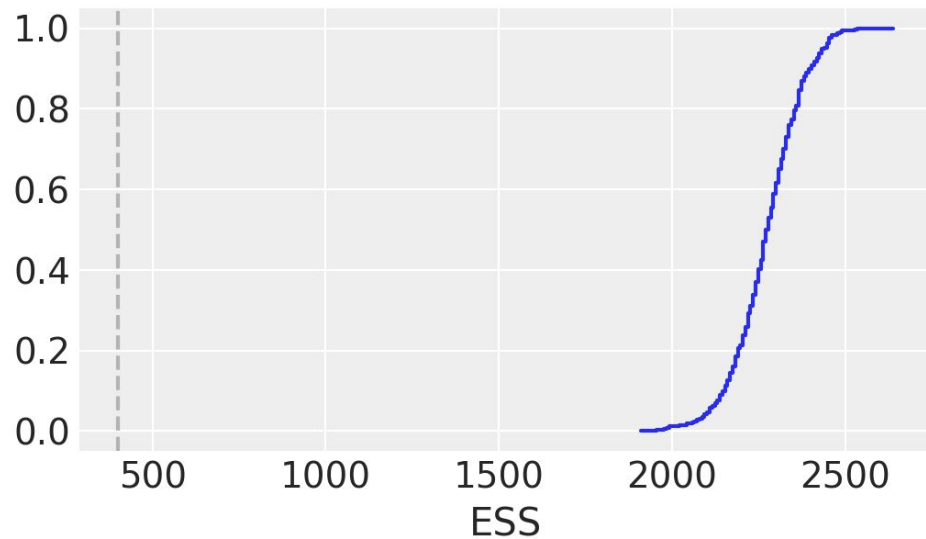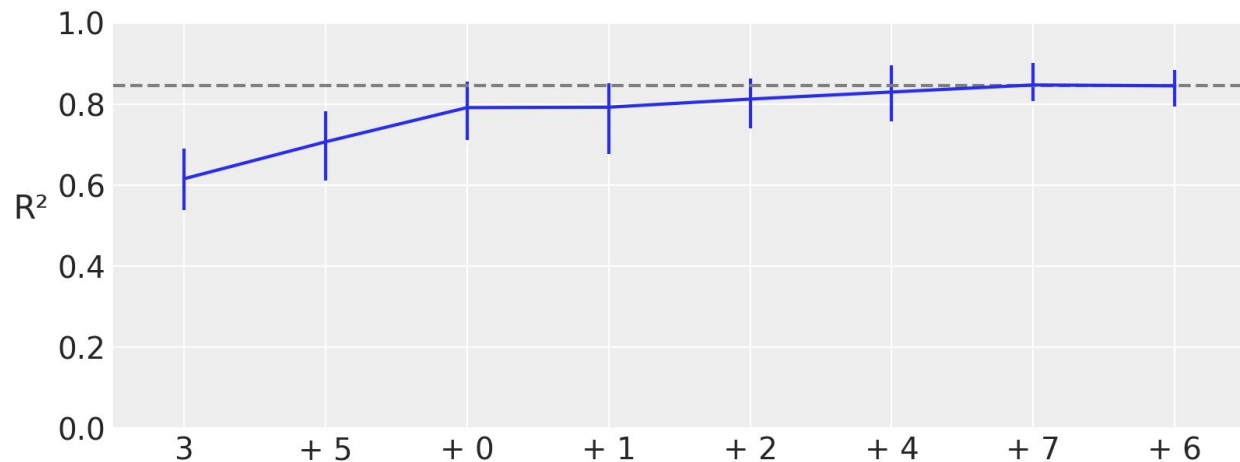
# Individual Conditional Expectance Plots



Age:X_0,

Pregnancies:X_1,

DiabetesPedigreeFunction:X_2,

Glucose:X_3,

Insulin:X_4,

BMI:X_5,

BloodPressure:X_6,
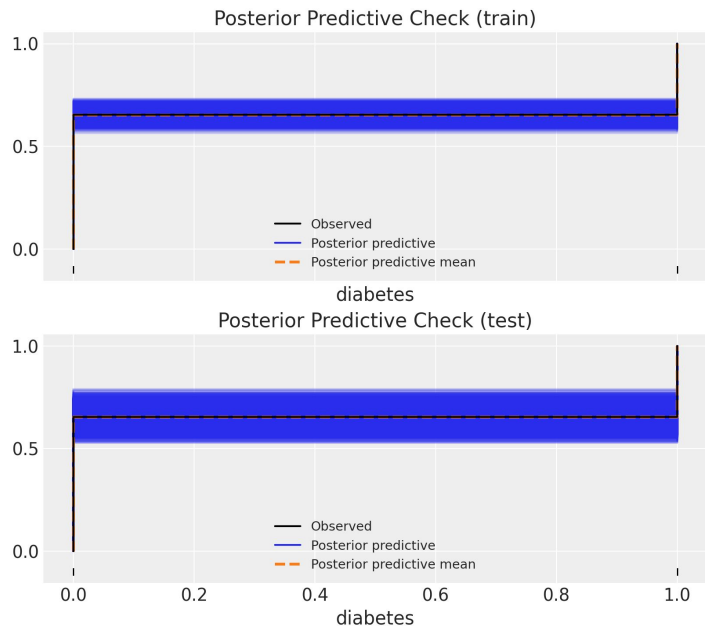
Glucose_Insulin_Product:X_7
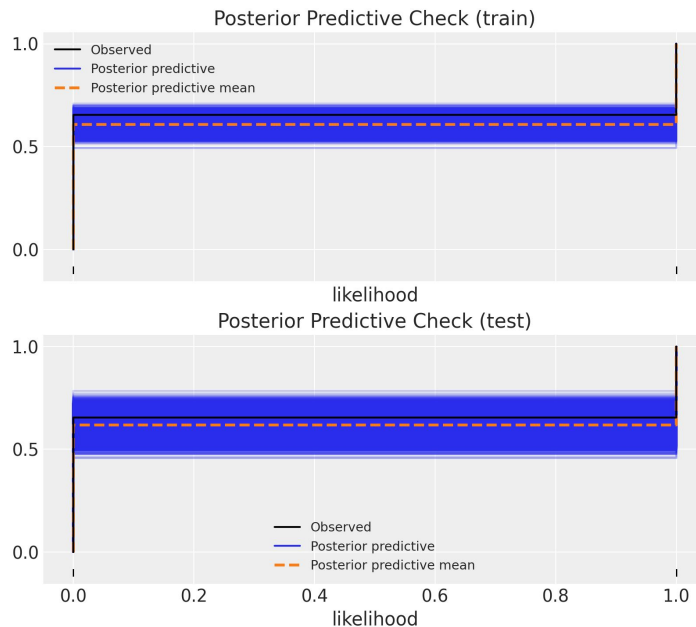
# Convergence diagnostics

# Variable Importance Plots



Age:X_0,
Pregnancies:X_1,
DiabetesPedigreeFunction:X_2,
Glucose:X_3,
Insulin:X_4,
BMI:X_5,
BloodPressure:X_6,
Glucose_Insulin_Product:X_7

# References

- PIMA Indians Diabetes Database. (2016, October 6). Kaggle. https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database

- Chipman, H. A., George, E. I., and McCulloch, R. E. (2010). "BART: Bayesian additive regression trees." The Annals of Applied Statistics, 4(1): 266–298 https://arxiv.org/pdf/0806.3286.pdf

- Inglis, A., Parnell, A., & Hurley, C. (2022). Visualizations for Bayesian Additive Regression Trees https://arxiv.org/pdf/2208.08966.pdf

- Decision Tree example image -Chris J Maddison STA314 Lecture 2 2021 https://www.cs.toronto.edu/~cmaddis/courses/sta314_f21/slides/lec02.pdf

# Thanks For Listening