

Explainable AI

- 1) Global - Explainability
- 2) Local Explainability
- 3) Post-hoc Explainability
- 4) Intrinsic -

1) Global - overall explainability
- once we build the model with help of the final model parameter we explain the why?

example - Linear - $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

where $\beta_0, \beta_1, \beta_2$ are the coefficients which get calculated & we explain what are the values/weighted of those features to get final value.

2) Local Explainability

- with the single row if we are predicting the output then we explain why per that model

like if x_1, x_2, x_3 are feature which produces value y_1 , then why y_1 & how get that value.

Post - a-hoc →

1) After training done we get some values so we need to use SHAP or LIME to explain the why?

Intrinsic → With help of this we can explain the features those play important role in calculate output value.

- Linear - we use coefficient
Random forest and decision tree we use calculate important feature.

Real world use cases

- 1) Banking - credit scoring, fraud detection
- 2) Health care - Patient Risk Scoring
- 3) Retail - Recommendation
- 4) Public sector - risk mgmt / assessment resource allocation.

challenges in model explainability

- Accuracy vs Interpretability
- more accurate models harder to explain
- Domain knowledge req.
- Communicating to non-tech audience
- Risk of misinterpretation.

Explainability builds trust, ensure fairness & meets regulations & improve model.



Low Blue Light
(Hardware Solution)
www.tuv.com
ID: 01100000

SHAP - shapley, additive explanation

Dataset \rightarrow avg prediction

$$\text{Base value} \rightarrow \phi = \frac{1}{n} \sum_{i=1}^n f(x_i)$$

$$\hat{y} = \text{Base value} + \text{sum (shapley value)}$$

$$y = \phi_0 + \sum_{i=1}^n \phi_i$$

How it works \rightarrow

- 1) calculate Base value
- 2) calculate Final prediction
- 3) check out additivity

$$\hat{y} = \phi_0 + \phi_1 + \phi_2 + \dots$$

$\uparrow \quad \uparrow \quad \leftarrow \text{feature}$

model specific - optimized

Component \rightarrow

model agnostic - no care which model we use.

1) Explainer \rightarrow calculate SHAP value

\rightarrow output value of explainer called Explanation

\rightarrow Explainer are brain of SHAP - calculate SHAP values

- Tree Explainer -

- Kernel Explainer - slow and non-differentiable

- Linear Explainer -

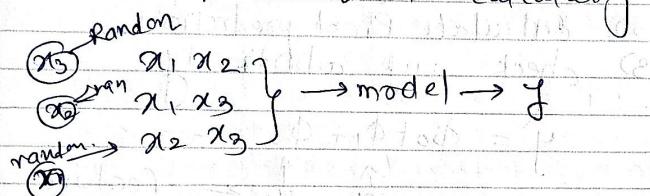
- Deep Explainer - forward pass





- 2) markers → control how SHAP wides or replace features when calculating their contribution
- They has background dataset

- suppose we have three feature x_1, x_2, x_3 & we are calculating y if we want to get contribution of x_1, x_2 without x_3 we can either replace with random value or constant or avg value so that doesn't matter while calculating y .



- 3) plot →
- summary plot
 - waterfall plot
 - dependence
 - decision plot

- 4) Background data → Reference Point

when SHAP calculates the contribution of a feature it need to simulate - What would the prediction be if we didn't know feature value?

Don't know - leave it blank, but model can't handle it. & SHAP replace with some reference point. The reference value come from Background data



Example →

$$f_1 = \text{Age} \cdot 0.1 + 2 \cdot 0.5 + 10 = 1$$

$$f_2 = \text{Income} \cdot 0.1 + 2 \cdot 0.5 + 10 = 1$$

Age	Income
30	60
40	80
35	70

$$p_1 = 10 + 0.5 \times 30 + 0.1 \times 60 = 31$$

$$p_2 = 10 + 0.5 \times 40 + 0.1 \times 80 = 38$$

$$p_3 = 10 + 0.5 \times 35 + 0.1 \times 70 = 34.5$$

$$\hat{y}_{\text{test}} = 10 + 0.5(30) + 0.1(60) = 45$$

$$\hat{y}_{\text{test}} = \phi_1 + \phi_{\text{age}} + \phi_{\text{income}}$$

$$45 = \text{basevalue} + \text{Age contribution} + \text{Income contribution}$$

calculate base

$$\begin{aligned}\phi_1 &= \text{Avg Average prediction of values.} \\ &= \frac{31 + 38 + 34.5}{3} = \frac{103.5}{3} \\ &= 34.5 \quad (\text{basevalue})\end{aligned}$$

what the base value means? if we don't have any feature value so our model provide prediction as base value.

② calculate contribution of Age column.

$$\text{Age} = 30 + 40 + 35 / 3 = 105 / 3 = 35$$

New dataset	Age	Income	model	\hat{y}
	30	60		
	40	80		
	35	70		

$$\begin{aligned}y_1 &= 10 + 0.5(35) + 0.1(100) \\&= 10 + 17.5 + 10 \\&= \underline{\underline{37.5}}\end{aligned}$$

$\phi_{age} = \text{final prediction} - \text{prediction without age}$

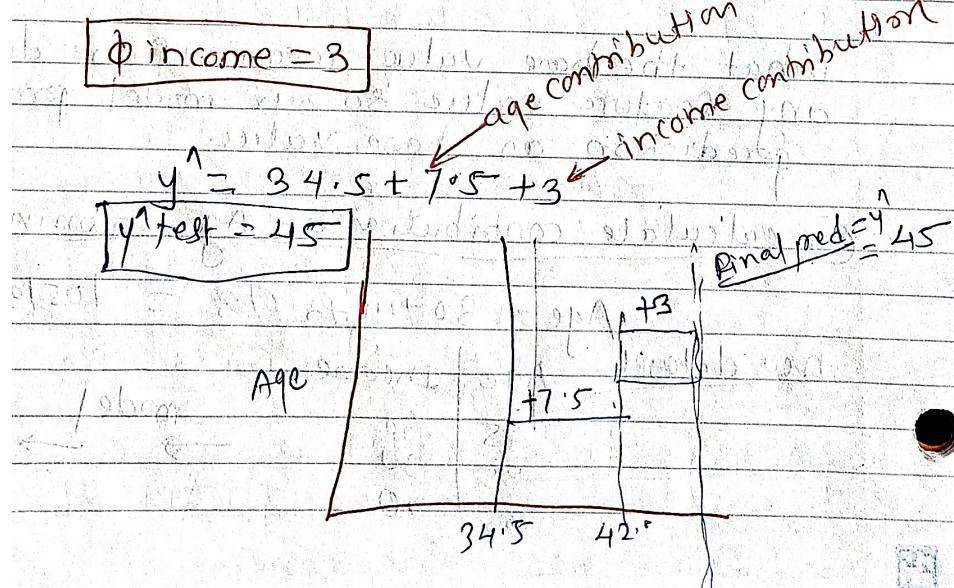
$$\begin{aligned}&= 45 - 37.5 \\&= 7.5\end{aligned}$$

$$\phi_{income} = 60 + 80 + 70 / 3 = 210 / 3 = 70$$

$$\begin{aligned}\phi_{income} &= 10 + 0.5(\cancel{35}) + 0.1 \times 70 \\&= 10 + 18 \cancel{.5} + 7 \\&= \cancel{35} + 8 \cancel{.5} + 7 \\&= 42\end{aligned}$$

$\phi_{income} = \text{final prediction} - \text{prediction without income}$

$$\begin{aligned}&= 45 - 37.75 \\&= 7.25\end{aligned}$$





SHAP - Shaply additive explanation

Summary -

- 1) calculate base value - Avg of all pred
- 2) calculate contribution of features
- 3) check which feature's contribution is more.

How to calculate shaply value

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|! (|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S))$$

$$\sum_{i \in N} \phi_i(v) = v(N) \quad \leftarrow \text{checks for additivity}$$

ϕ_i = individual contribution of feature.

$v(N)$ = final prediction.

Game theory → example of kaggle competition.

Suppose we have 3 players P₁, P₂, P₃ who entered in the kaggle competition. They all have their strength & whenever they team up, win 1st prize. So now suppose they win \$900k then with our normal calculation everyone get \$300k but this is not shaply. Shaply does. Shaply calculate the contribution of each player/value for the competition so that the shaply value concept & 3 players are feature and shap values are contribution of each player.



Low Blue Light
(Hardware Solution)
www.tuv.com
ID: 0217000905



$$\phi(v) = \sum_{S \subseteq N \setminus \{v\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (v(S \cup \{v\}) - v(S))$$

$|N| \rightarrow$ length of set - no of feature/player.

$|N| \rightarrow 3$ (in this example)

$S \rightarrow$ subset $\{P_1, P_2\}$ example.

$|S| \rightarrow$ length of subset.

$v(S) = F(S) \rightarrow$ model or mathematical function or prediction,

$v(\{P_1, P_2\}) = \$5000k$ - outcome / payoff value.

$v(S) + v(S \cup \{v\}) \rightarrow$ total $\rightarrow \$9000k$ in this example.

subset of $\{P_1, P_2\}$ \leftarrow in this example

$$v(S \cup \{v\}) - v(S) = \sum_{i \in S} \{P_1, P_2\} \setminus \{P_3\} - \{P_1, P_2\}$$

\rightarrow will get individual contribution of P_3

$\sum_{S \subseteq N \setminus \{v\}}$ = calculation of all subsets where $v(S)$ summation is not part of it

$N = \{P_1, P_2, P_3\}$

subset $\left\{ \begin{array}{l} S_1 = \{P_1, P_2\} \text{ for } \{P_1, P_2\} \text{ if } i = P_3 \\ S_2 = \{P_1\} \end{array} \right.$

$S_3 = \{P_2\}$

$S_4 = \emptyset$

here S is ④

$$= \sum_{A \subseteq N \setminus \{P_3\}} \frac{12!}{(3-2-1)!} \cdot \left(\sum_{i \in A} \{P_1, P_2\} \right) \cdot \left(\sum_{i \in A} \{P_1, P_2\} \right)$$



*Example - 2 player Problem.

player A } 10
player B }

} given.

$$1) v(\emptyset) = 0$$

$$2) v(\{A\}) = 4$$

$$3) v(\{B\}) = 6$$

$$4) v(\{A, B\}) = 10$$

$\phi_A \rightarrow$ shaply values \rightarrow for A

$$\phi_n = \frac{|S_1|! (N-s-1)!}{N!} (v(S \cup \{i\}) - v(S))$$

$$\text{emp set } S_1 = \emptyset, N=2, s=0 \quad \frac{0! (2-0-1)!}{2!} = \frac{1 \times 1}{2} = \frac{1}{2}$$

$$= \frac{1}{2} (v(S \cup \{i\}) - v(\emptyset))$$

$$= \frac{1}{2} (v(S \cup \{i\}) - 0)$$

$$= \frac{1}{2} (v(\emptyset \cup \{A\})) = \frac{1}{2} \times 4 = 2$$

$S_2 = \{B\}$

$$= \frac{1}{!} \frac{(2-1-1)}{2!} = \frac{1}{2}$$

$$= \frac{1}{2} (v(S \cup \{i\}) - v(\{B\}))$$

$$= \frac{1}{2} (v(\{B\} \cup \{A\}) - 6)$$

$$= \frac{1}{2} \times (10 - 6) = \frac{1}{2} \times 4 = 2$$

$$\phi_A = \phi_0 +$$

$$= S_1 + S$$

$$= 2 +$$

$$= 4$$

$$\text{Player B} = \frac{1}{2} \circ (v(S \cup \{i\}) - v(S))$$

↑ same

$$S_1 = \emptyset = \frac{1}{2} \circ (v(\emptyset \cup \{B\}) - v(\emptyset))$$

$$= \frac{1}{2} \circ (v(\{B\}) - v(\emptyset))$$

$$= \frac{1}{2} \times (6 - 0) = \frac{1}{2} \times 6 = 3$$

$$S_2 = \emptyset^A$$

$$= \frac{1}{2} \times (v(\emptyset^A \cup \{A\}) - v(\emptyset^A))$$

$$= \frac{1}{2} \times (10 - 4) = \frac{1}{2} \times 6 = 3$$

$$\phi_B = \phi_A + \phi_B = 3 + 3 = 6$$

$$\phi_A = 4, \phi_B = 6$$

additivity rule = $\sum v(\phi_i)$

$$v(\{\phi_A, \phi_B\}) = v(\phi_A) + v(\phi_B)$$

$$10 = 4 + 6$$

$$10 = ((\phi_A \cup \phi_B) \vee)$$

$$10 = (\phi_A \vee \phi_B)$$

$$10 = ((\phi_A \vee \phi_B) \vee \phi_C)$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C))$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee \phi_D$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee ((\phi_A \vee \phi_B) \vee \phi_D)$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee ((\phi_A \vee \phi_B) \vee \phi_D) \vee \phi_E$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee ((\phi_A \vee \phi_B) \vee \phi_D) \vee ((\phi_A \vee \phi_B) \vee \phi_E)$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee ((\phi_A \vee \phi_B) \vee \phi_D) \vee ((\phi_A \vee \phi_B) \vee \phi_E) \vee \phi_F$$

$$10 = ((\phi_A \vee \phi_B) \vee ((\phi_A \vee \phi_B) \vee \phi_C)) \vee ((\phi_A \vee \phi_B) \vee \phi_D) \vee ((\phi_A \vee \phi_B) \vee \phi_E) \vee ((\phi_A \vee \phi_B) \vee \phi_F)$$



3 Player Problem (P_1, P_2, P_3) - Player

$\{P_1\}$	6
$\{P_2\}$	12
$\{P_3\}$	6
$\{P_1, P_2\}$	24
$\{P_1, P_3\}$	18
$\{P_2, P_3\}$	18
$\{P_1, P_2, P_3\}$	30

$$\phi_{P_1} \phi_{P_2} \phi_{P_3} = ?$$

$\phi_{P_1} = ?$ = shaply value = contribution of P_1 subset length

make subset for $P_1 \rightarrow \{P_2\} = 1$

$$\{P_3\} = 1$$

$$\{P_2, P_3\} = 2$$

$$\{ \} = 0$$

For P_2 - $\{P_1\} = 1$

$$\{P_3\} = 1$$

$$\{P_1, P_3\} = 2$$

$$\{ \} = 0$$

For P_3 - $\{P_1\} = 0$

$$\{P_2\} = 1$$

$$\{P_2, P_3\} = 2$$

formula = \sum weight \times (with feature - without fe

$$|N| = 3$$

$$\text{weight} = \frac{s! (|N| - s - 1)!}{|N|}$$

$$\text{weight for } P_1 = \frac{0! (3-0-1)}{3!} = \frac{1(2!)}{4(3!)} = \frac{2}{3 \times 2} = \frac{1}{3}$$

$$\rightarrow S_2 \text{ for } P_2 = \frac{1! (3-1-1)}{3!} = \frac{1}{3} = \frac{1}{6}$$

$$\rightarrow S_3 \text{ for } P_3 = \frac{1}{6}$$

$$\rightarrow S_4 \text{ for } (P_2, P_3) = \frac{2! (3-2-1)}{3!} = \frac{2! \times 0!}{3!} = \frac{2}{3 \times 2} = \frac{1}{3}$$

$$\rightarrow S_5 \text{ for } (P_1, P_2) = \frac{2! (3-2-1)}{3!} = \frac{2! \times 0!}{3!} = \frac{1}{3}$$

$$\rightarrow S_6 \text{ for } (P_1, P_3) = \frac{2! (3-2-1)}{3!} = \frac{1}{3}$$

$$\Phi_{P_1} = \{\}, \{P_2\}, \{P_3\}, \{P_2, P_3\}$$



	outcome(s)	$v(S \cup \{P_1\})$	diff.			
$\emptyset P_1$	0	6	$1/3 \cdot 2$			
$\emptyset P_2 \}$	12	28	$12 \cdot 1/6 \cdot 2$			
$\emptyset P_3 \}$	6	18	$6 \cdot 1/6 \cdot 1 = 1/6$			
$\emptyset P_1 P_2 P_3 \}$	18	30	$12 \cdot 1/3 \cdot 4 = 1/3 \cdot 4$			
$\emptyset P_1 =$	$2+2+2+4=10$					
$\emptyset P_2$	outcome	$v(S \cup \{P_2\})$	diff. weigh out			
$\emptyset P_2$	0	12	$1/3 \cdot 4$			
$\emptyset P_1 \cup P_2 \}$	6	24	$12 \cdot 1/6 \cdot 2$			
$\emptyset P_3 \}$	6	18	$6 \cdot 1/6 \cdot 2$			
$\emptyset P_1 P_2 \}$	18	30	$12 \cdot 1/3 \cdot 4$			
$\emptyset P_3$	P_3	outcome	$v(S \cup \{P_3\})$	diff	wei	dir
$\emptyset P_3$	0	6	12	6	$1/3$	2
$\emptyset P_2 \}$	12	18	6	$1/6$	1	
$\emptyset P_1 \cup P_3 \}$	6	18	12	$1/6$	2	
$\emptyset P_1 P_2 \cup P_3 \}$	24	30	6	$1/3$	2	
$\emptyset P_1 = 10$						
$\emptyset P_2 = 4+3+2+4 = 13$						
$\emptyset P_3 = 2+1+2+2 = 7$						
additivity	$= \emptyset P_1 + \emptyset P_2 + \emptyset P_3$					
	$= 10+13+7$					
	$= 30$	✓				
			total	$\{P_1 P_2 P_3\}$		



order no - 6678473305 - Sun no 30 to 103

SHAP (SHapley Additive exPlanations)

- model agnostic is slow compare to other explainer but we can use for any model
- simple data - simple model - good explainability
- complex data need complex model
 - so they works fine but loses interpretability, so that's why it's called black box model as if they gives good prediction & high accuracy but have no interpretability.
 - so we need some model which provide to explain the model + why?
- with explainable AI we can find the answer if model provide certain prediction then why its that so.

SHAP Working -

- it based on Shapley value — game theory.
- Shapley value comes via feature - x
- so we have x and y model output
- Shapley value — $\text{Avg of } f(x) \text{ or base value}$
 - need background data
 - background data mostly be train
- calculate base value
- Shap value move base value
 - if its +ve - base value goes to higher
 - if its -ve - base value goes to lower



SHAP - $\begin{cases} \text{model Agnostic} \\ \text{model specific} \end{cases}$

$$\phi_0 = \arg \text{ prediction}$$

$$\phi_0 + \sum_{i=1}^d \phi_i = \hat{y}$$

- shap value have - magnitude - how much original prediction effected
- sign - tve
- SHAP uses
 - mathematical concept - game theory
 - calculate individual contribution of each feature for prediction.

LIME - local interpretable model Agnostic explanation

- local explanation only.
- interpretable - explainable
- model agnostic - dont care what model we use.
- LIME learns own model which is simple so interpretable as well.
- SHAP - surrogate model concept.
- implementation fast
- no more complex computations required.

example $f_1 f_2 f_3 f$

10 20 30 135

so we create perturbation / disturbance.

Added	f_1	f_2	f_3	y
	→ 10	20	30	135 ↑
Random	11	28	28	y_{n^1}
	13	19	31	y_{n^2}
	20	30	35	y_{n^3}
	3	18	35	y_n^4

slight disturbances —

* Like we create synthetic data.

x_{new} → new → interpretable
model (example linear)

↓
train

↓
train regression model

called surrogate model

suppose we have linear regression so model
trying to find coefficient values

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$$

- so we have 5 feature

once we train the model we get our coef and
intercept.



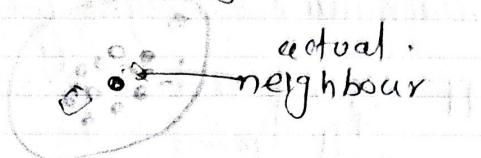
coeff - +ve or -ve

which takes prediction in +ve direction
or -ve direction.

- there is another concept - weight

synthetic data - neighbourhood

example



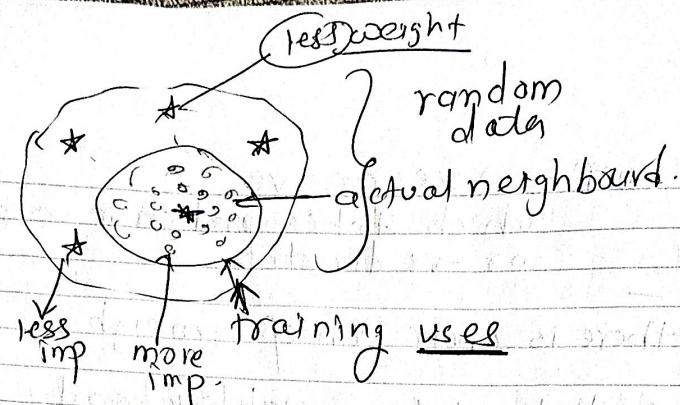
why synthetic data - because it shows that
the disturbance in the actual value
how to affect the prediction.

so here if we create the synthetic data or
we can say neighbour hood we need the
neighbour which don't have big disturbance
so those created far from actual point
we apply less weight so model will
give not actual learn from them but
from actual neighbours

- increase weight for actual neighbourhood
- less weight for

- weight assigned by instances

euclidean
distance



Weighted mechanism - weight base

$$TF(x, s) = \exp\left(-\frac{D(x, s)}{\sigma^2}\right)$$

↑
actual synthetic

we use kernel Function for weight assignment

example

Pt	x_1	x_2	x_3	x_{new}	\rightarrow new	model pred
1	2.2	4.8	1.1			10.4
2	1.8	5.2	0.9			9.4
3	2.5	5.5	1.0			11.2
4	1.5	4.5	1.5			8.7
5	2.0	6.6	0.5			10.8
6	3.0	5.0	1.0			12.0

Black box model.

in new vector \vec{y} distances weight - with parallel

(2.2, 4.8, 1.1)	10.4	0.300	0.91
(2.0, 6.0, 0.5)	10.8	1.1120	0.2865

↓
Φ-shap

$$\hat{y}_{\text{surrogate}}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$
$$= 3.93723$$

intercept work as base value shape.

If we feed above value to our surrogate model then we gets the coef - $\beta_0, \beta_1, \beta_2, \beta_3$

so coef values are your explanation

- create synthetic data
- calculate prediction using black box model
- we calculate distances from actual point using euclidean distance
- assign weights \rightarrow distance less - weight more
distance more - weight less
- so model train on the new data.
and get intercept and coef value
which explain model easily.

Advantage - fast to implement

- Hardware light

Disadvantage - coef values not stable.

linear model \leftarrow synthetic data

random \rightarrow change data

- 3) lime - no mathematical approach
shape - have mathematical grounding of game theory
- 4) surrog ate - linear regression

Non linear dataset → not work well.

- Lime - not used in production setup
 - do use random synthetic data
every time so not too doofy.
- local explanation / No global explanation
- Mat mosse plots in LIME.