

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 - Under - $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

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 for that model like 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 mgt / assessment
resource allocation.

challenges in model explainability

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

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

SHAP - Shapley, additive explanation

Dataset → avg prediction

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

$$y^* = \text{Base value} + \text{sum(Shaply value)}$$

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

How it works →

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

$$y^* = \phi_0 + \phi_1 + \phi_2 + \dots$$

$\uparrow \quad \uparrow \quad \swarrow$ feature.

model specific - optimized

Component →

model agnostic - no care which model we use.

→ Explainer → calculate SHAP value

→ output value of explainer called explanation

→ Explainer are brain of SHAP - calculate SHAP values

- TreeExplainer -

- kernelExplainer - slow

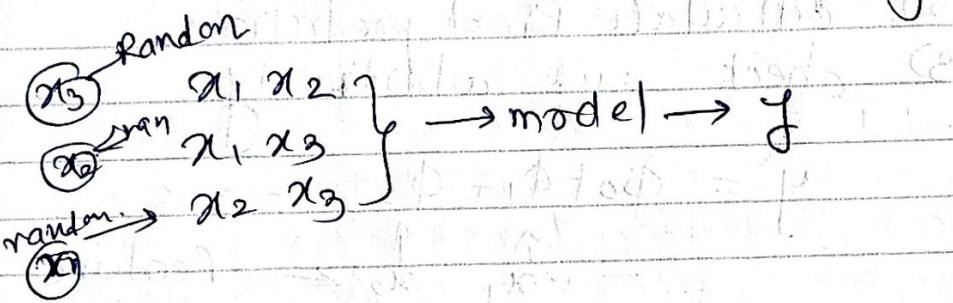
- linearExplainer -

- DeepExplainer -

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 f 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 f .



3) plot →

- summary plot

- waterfall plot

- Dependence

- Desection 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 →

$$f_1 = \text{Age}$$

$$f_2 = \text{Income}$$

Age	Income
30	60
40	80
35	70

prediction = $10 + 0.5 \times \text{Age} + 0.1 \times \text{Income}$

$$P_1 = 10 + 0.5 \times 30 + 0.1 \times 60 = 31$$

$$P_2 = 10 + 0.5 \times 40 + 0.1 \times 80 = 38$$

$$P_3 = 10 + 0.5 \times 35 + 0.1 \times 70 = 34.5$$

$$\hat{y}^{\text{test}} = 10 + 0.5(50) + 0.1(80) = 45$$

$$\hat{y}^{\text{test}} = \phi_0 + \phi_1 \text{Age} + \phi_2 \text{Income}$$

ϕ_0 = base value + Age contribution + Income contribution

calculate base

$$\phi_0 = \text{Avg Average prediction of values.}$$

$$= \frac{31 + 38 + 34.5}{3} = \frac{93.5}{3}$$

$$= 31.166666666666666 \approx 31.2$$

\hat{y}^{test} = base value + Age contribution + Income contribution

what the base value mean → 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
	35	60
	35	80
	35	70

→ model → \hat{y}^{test}

$$y_1 = 10 + 0.5(35) + 0.1(100)$$

$$= 10 + 17.5 + 10$$

$$= \underline{37.5}$$

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

$$= 45 - 37.5$$

$$\boxed{\phi_{age} = 7.5}$$

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

$$y_{pred} = 10 + 0.5(\cancel{35}) + 0.1 \times 70$$

$$= 10 + \cancel{17.5} + 7$$

$$= \cancel{35} + \cancel{17.5} \quad 42$$

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

$$= 37.5 -$$

$$= 45 - 37.5 \quad 42$$

$$\phi_{income} = 7.5 / 3$$

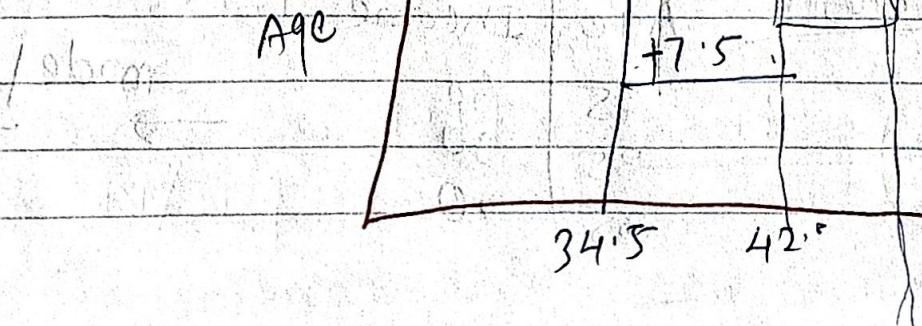
$$\boxed{\phi_{income} = 3}$$

$$y^1 = 34.5 + 7.5 + 3$$

$$\boxed{y^1_{test} = 45}$$

$$\text{age contribution} \quad \text{income contribution}$$

$$\text{Final pred} = 45$$



SHAP - Shaply additive explanation

Summary -

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

How to calculate shaply value

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

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

ϕ_i = individual contribution of feature.

$v(N)$ - final prediction.

Game theory → example kaggle competition.

Suppose we have 3 players P_1, P_2, P_3 who entered in the kaggle competition. They all have their strength & whenever they team up winning 1st prize. So

Now suppose they win \$900k then with our normal calculation everyone get \$300k but this is not shaply does. Shaply calculate the contribution of each player/value in the competition

so that's the shaply value concept so players are feature and shap values are contribution of each player.

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

$|N| \rightarrow$ 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.

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

subset of $\{P_1, P_2\}$

P_3 ← in this example

$$v(S \cup \{i\}) - v(S) = \text{sum}_{j \neq i} \{P_1, P_2\} \cup \{P_3\} - \{P_1, P_2\}$$

= will get individual contribution of P_3

$\sum_{S \subseteq N \setminus \{v\}}$ = calculation of subsets where i is not part of it.

subset S

$$\left\{ \begin{array}{l} S_1 = \{P_1, P_2\} \text{ for } \{P_2, P_3\} \text{ if } i = P_3 \\ S_2 = \{P_1\} \\ S_3 = \{P_2\} \\ S_4 = \{P_3\} \end{array} \right.$$

here S is ④

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

* Example - 2 player Problem.

player A } 10
player B } 6

} given.

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

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

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

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

$\phi_A \rightarrow$ shaply value \rightarrow for A

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

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

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

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

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

$$S_2 = \{B\}$$

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

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

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

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

$$\begin{aligned}\phi_A &= \phi_0 + \phi_1 \\ &= S_1 + S_2 \\ &= 2 + 2 \\ &= 4\end{aligned}$$

$$\text{Player B} = \frac{1}{2} \cdot (v(s \cup \{\bar{A}\}) - v(s))$$

same

$$S_1 = \emptyset$$

$$= \frac{1}{2} \cdot (v(\emptyset \cup \{\bar{B}\}) - v(\emptyset))$$

\uparrow
empty set

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

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

$$S_2 = \emptyset_A$$

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

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

$$\therefore \Phi_B = \Phi_A + \emptyset_A = 3 + 3 = 6$$

$$\emptyset_A = 4 \quad \emptyset_B = 6$$

additivity rule: $v(\emptyset) = \sum v(\emptyset_i)$

$$v(\{\bar{A} \cup \bar{B}\}) = v(\emptyset_A) + v(\emptyset_B)$$

$$10 = 4 + 6 \quad \checkmark$$

3 Player Problem

(P₁, P₂, P₃) - 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

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

$\phi_{P_1} = ?$ = Shapley value = contribution of P₁ subset length

make subset for P₁ → {P₂} = 1
{P₃} = 1
{P₂, P₃} = 2
{ } = 0

For P₂ - {P₁} = 1
{P₃} = 1
{P₁, P₃} = 2
{ } = 0

For P₃ = {P₁} = 0
{P₂} = 1
{P₂, P₃} = 2

formula = \sum weight × (with feature - without feature)

$$|N| = 3$$

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

$$\text{weight for } S_1 = \frac{0!(3-0-1)}{3!} = \frac{1(2!)}{3!} = \frac{2}{3 \times 2} = \frac{1}{3}$$

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

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

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

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

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

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

outcome(s) $v(S \cup \{P_1\})$ diff.

P_1

\emptyset

0

6

8

2

$\{P_2\}$

12

28

12

2

$\{P_3\}$

6

18

12

1

$\{P_2, P_3\}$

18

30

12

4

$$\Phi P_1 = 2 + 2 + 2 + 4 = 10$$

ΦP_2

P_2

\emptyset

0

18

12

1/3

4

$\{P_1\}$

6

24

18

1/6

2

$\{P_3\}$

6

18

12

1/6

2

$\{P_1, P_3\}$

18

30

12

1/3

4

ΦP_3

P_3

\emptyset

0

6

6

1/3

2

$\{P_2\}$

12

18

6

1/6

1

$\{P_1\}$

6

18

12

1/6

2

$\{P_1, P_2\}$

24

30

6

1/3

2

$$\Phi P_1 = 10$$

$$\Phi P_2 = 4 + 3 + 2 + 4 = 13$$

$$\Phi P_3 = 2 + 1 + 2 + 2 = 7$$

$$\text{additivity} = \Phi P_1 + \Phi P_2 + \Phi P_3$$

$$= 10 + 13 + 7$$

$$= 30 \quad \checkmark$$

total $\{P_1, P_2, P_3\}$

order no - 6678473305 - Sun no 30 to 103

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 loose ⁴⁸⁴ interpretability, so that's why its 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 sharply value - game theory.
- sharply value comes via feature - x
- so we have x and of model output
- Shaply value -
 - need background data
 - background data mostly be ztrain
- calculate base value
- shap value more base value
 - if its +ve - base value goes to higher
 - if its -ve - base value goes to lower

SHAP — ^{model Agnostic}
_{model specific}

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

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

- shap value have - magnitude — how much original prediction effected
- sign \downarrow +ve
- SHAP uses \rightarrow 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 trains own model which is simple so interpretable as well.
- SHAP — surrogate model concept.
- implementation fast
- no more complex computations required.
- < example $f_1 f_2 f_3 f_4$
10 20 30 135
so we create perturbation / disturbance.

Actual

f_1 f_2 f_3 y

Random	10	20	30	135	41
11	20	28	135	41	1
13	19	31	135	41	2
9	20	30	135	41	3
8	18	35	135	41	4

slight disturbance -

* like we create synthetic data.

x_{new} \rightarrow $f_{new} \rightarrow$ interpretable

model (example linear)

↓
train

↓
train regression model
called surrogate model

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

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

- QD we have 5 feature

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

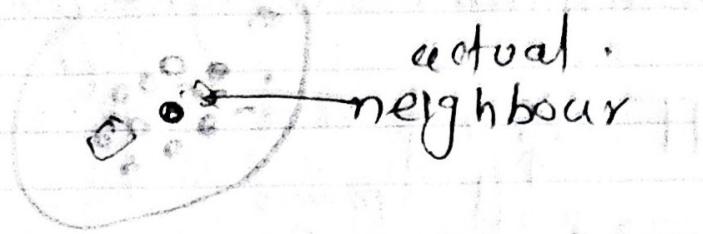
coeff - +ve or -ve

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.

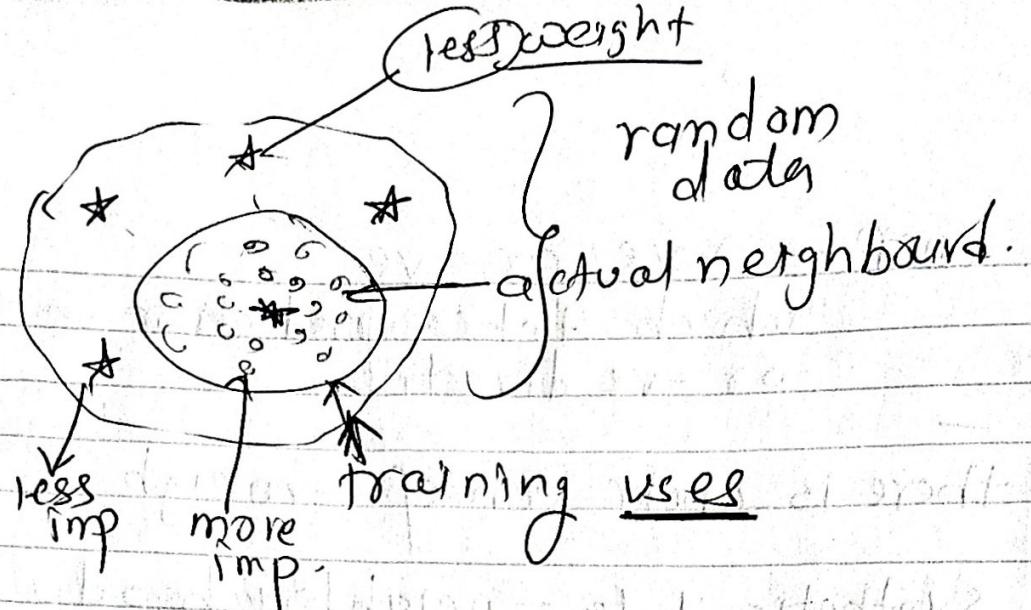
so here if we create the synthetic data or
we can say neighbour hood we need the
neighbour which don't have big disturbance
so those created far from actual point
we apply less weight so model will
give not actual learn from them but
from actual neighbours.

- increase weight for actual neighbourhood

- less on far

- weight assigned by distances

euclidean
distances



weighted mechanism - weight base

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

↑ ↑
actual synthetic

we use **kernel function** for weight assignment

example

Pt	x ₁	x ₂	x ₃	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

x_{new} ↓
↓ new

Black box model.

new vector \vec{y} distances weight - with kernel function

(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

$$\hat{y}_{\text{surrogate}}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$
$$= 2.93773$$

intercept work as base value shap.

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 from data

- 3) lime - no mathematical approach
shape - have mathematical boundary of eg one theo
- 4) surrogate - linear regression

Non linear ~~dataset~~ → not work well.

- Lime - not used in production setup
 - as use ~~random~~ synthetic data every time so not touchability
- local explanation / No global explanation
- Not more plots in lime.