

XAI : SHAP $\xrightarrow{(2017)}$

Applied AI Course

Shapely values

- Co-operative Game Theory
- Lloyd Shapley 1951; Nobel Prize (2012)
- Study in ML Context
- Lots of mathematical properties & research (~70 yrs)

Game - Theory

Players: $f_1, f_2, f_3, \dots, f_d$ (features)

Game → co-operatively work together

→ generate a model output

[GAIN/[↑]profit]

Shapely value

- measures how much has each player (feature) contributed to the profit (output)
- Mathematically rigorous
- Intuitive

Intuition : Regression Problem

Model

Predict Salary

$$\bar{y} = \text{avg of } y_i's$$

features

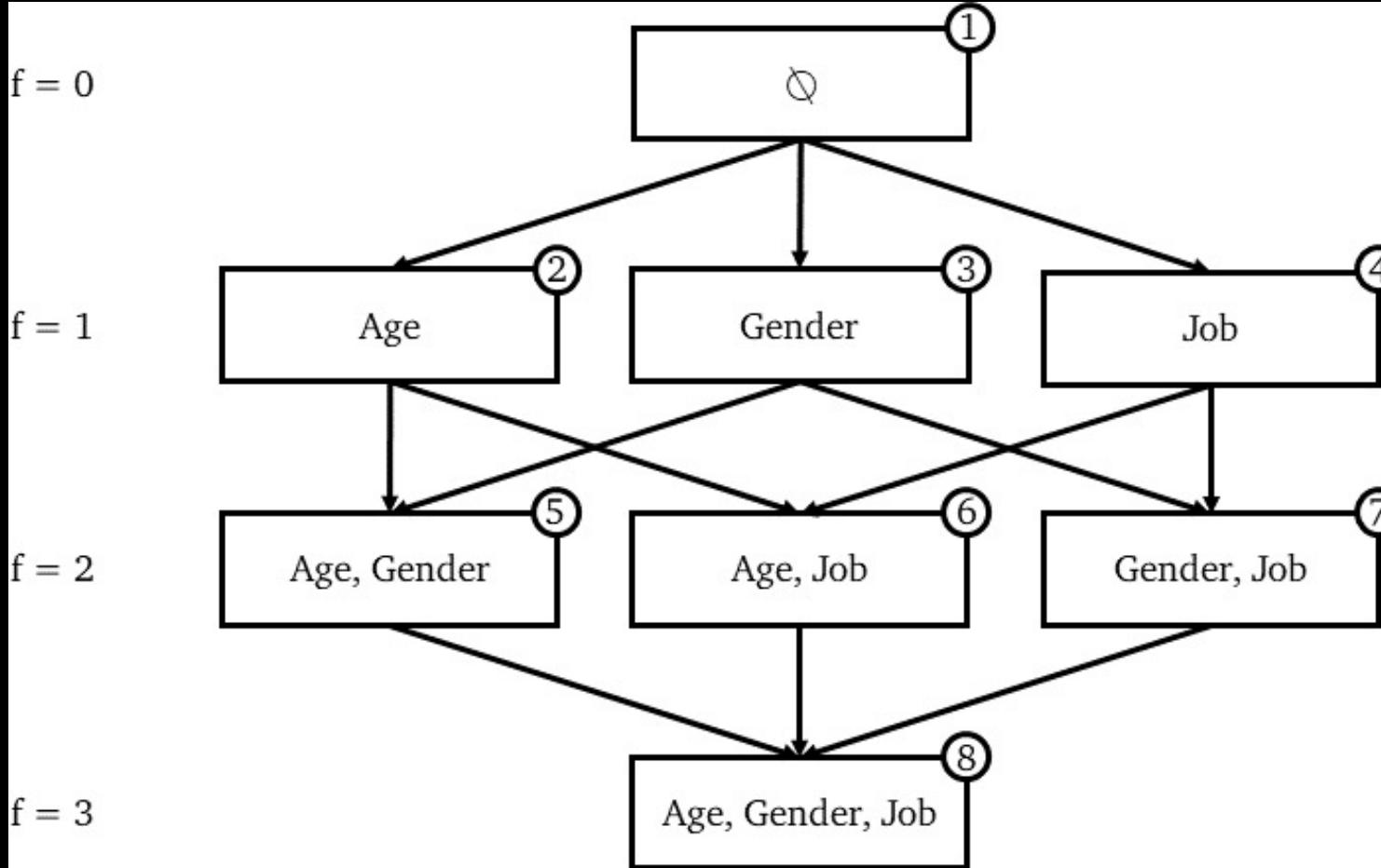
Age, gender, Job

x_0 : Query point

output

y_i

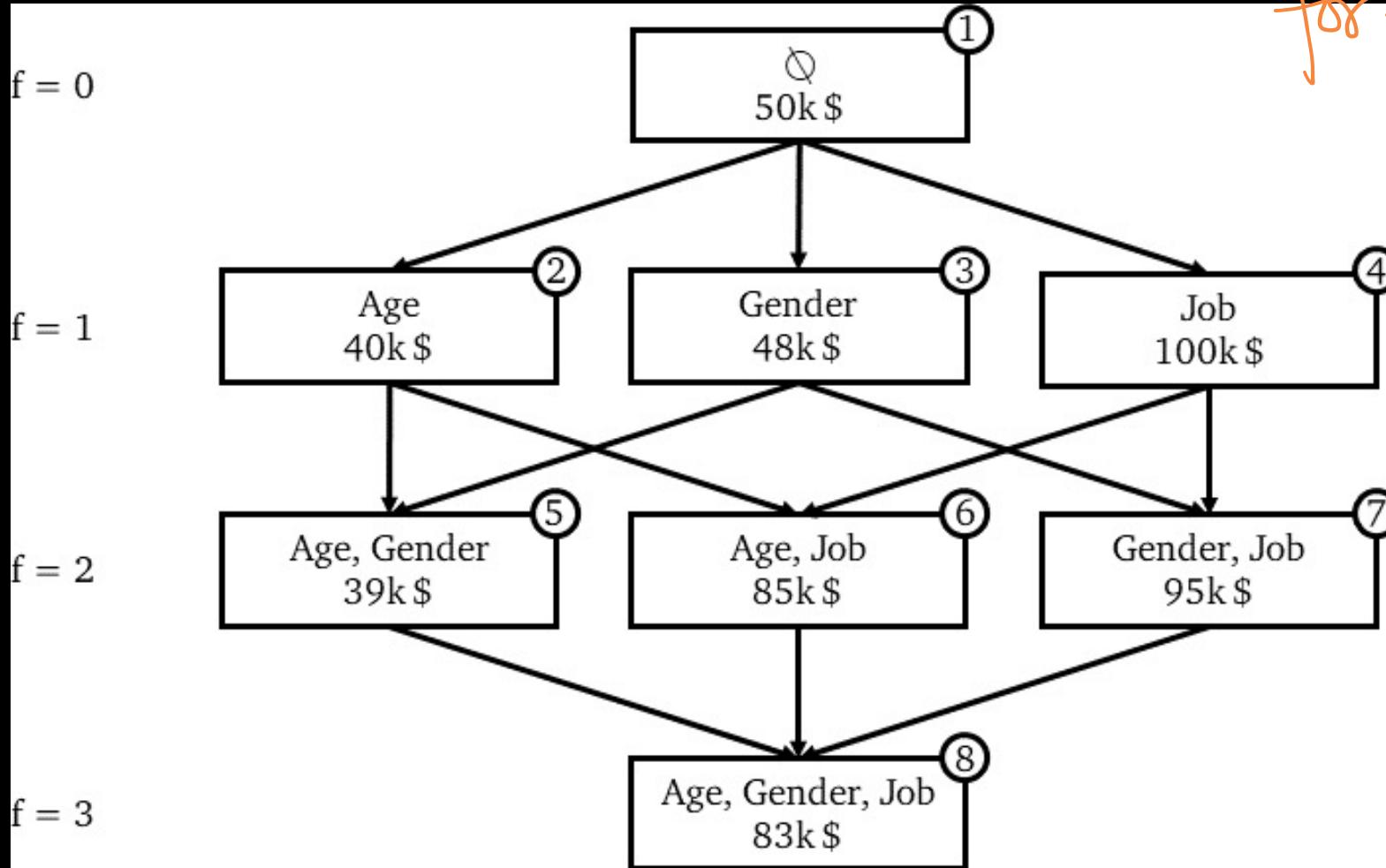
\hat{y}_0 : predicted value for x_0



d - features

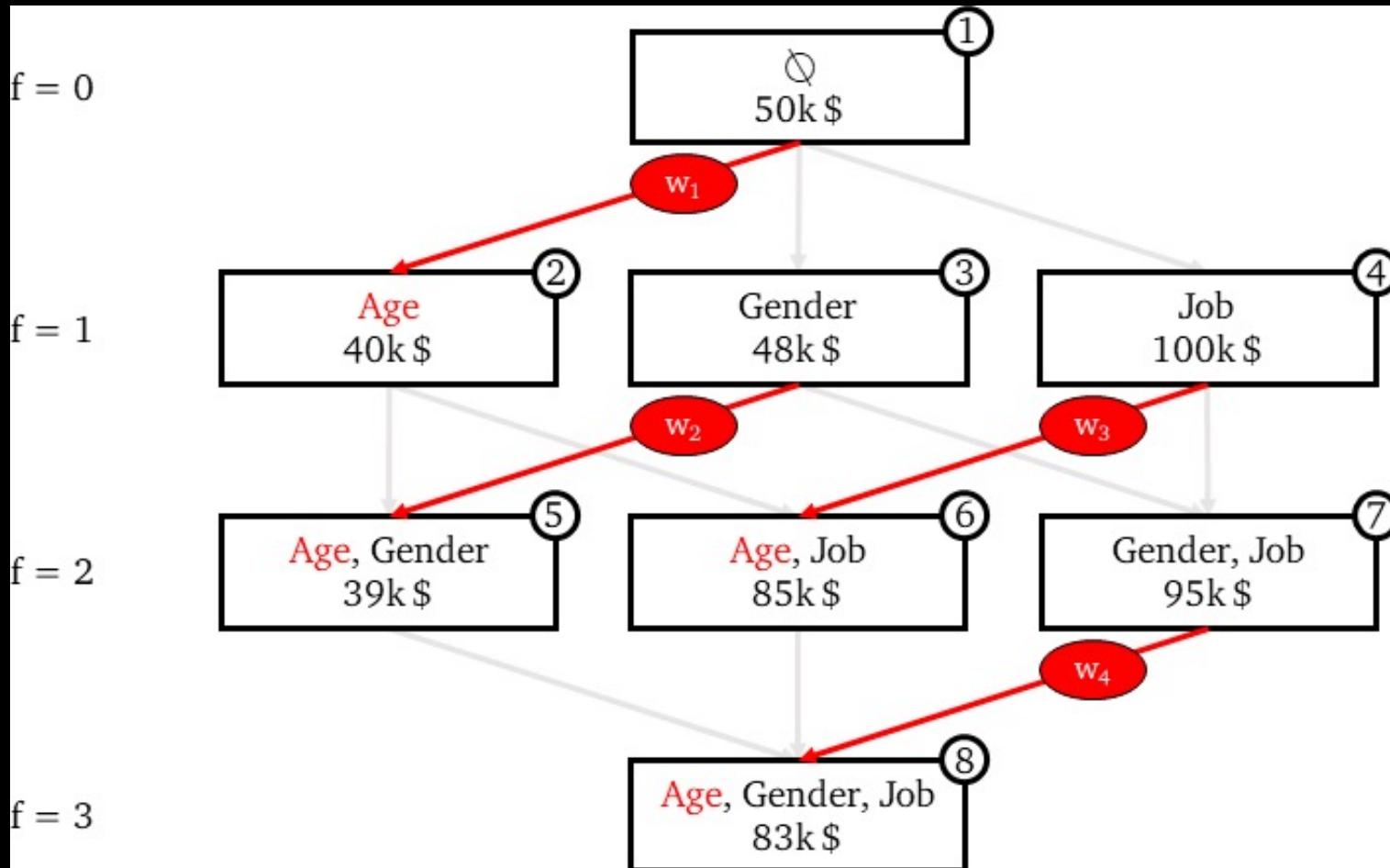
↓

2^d - models



y_a given $x_{\bar{a}}$
 for each of the
 2 models

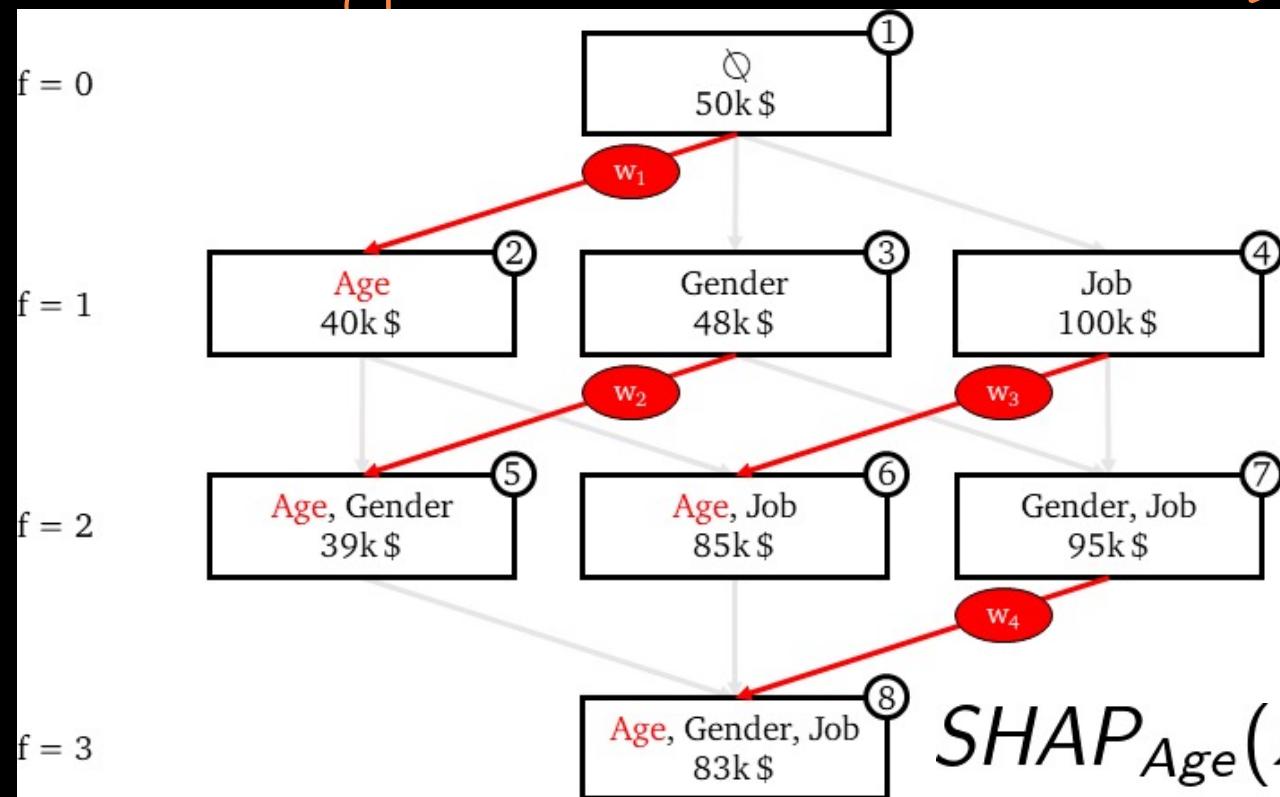
Marginal Contribution of a feature (say Age)



$$\begin{aligned} \text{MC}_{\text{Age}, \{\text{Age}\}}(x_0) \\ = -10k \end{aligned}$$

$$\begin{aligned} \text{MC}_{\text{Age}, \{\text{Age, Gender}\}}(x_0) = -9k \end{aligned}$$

Marginal contribution of a feature (say Age)



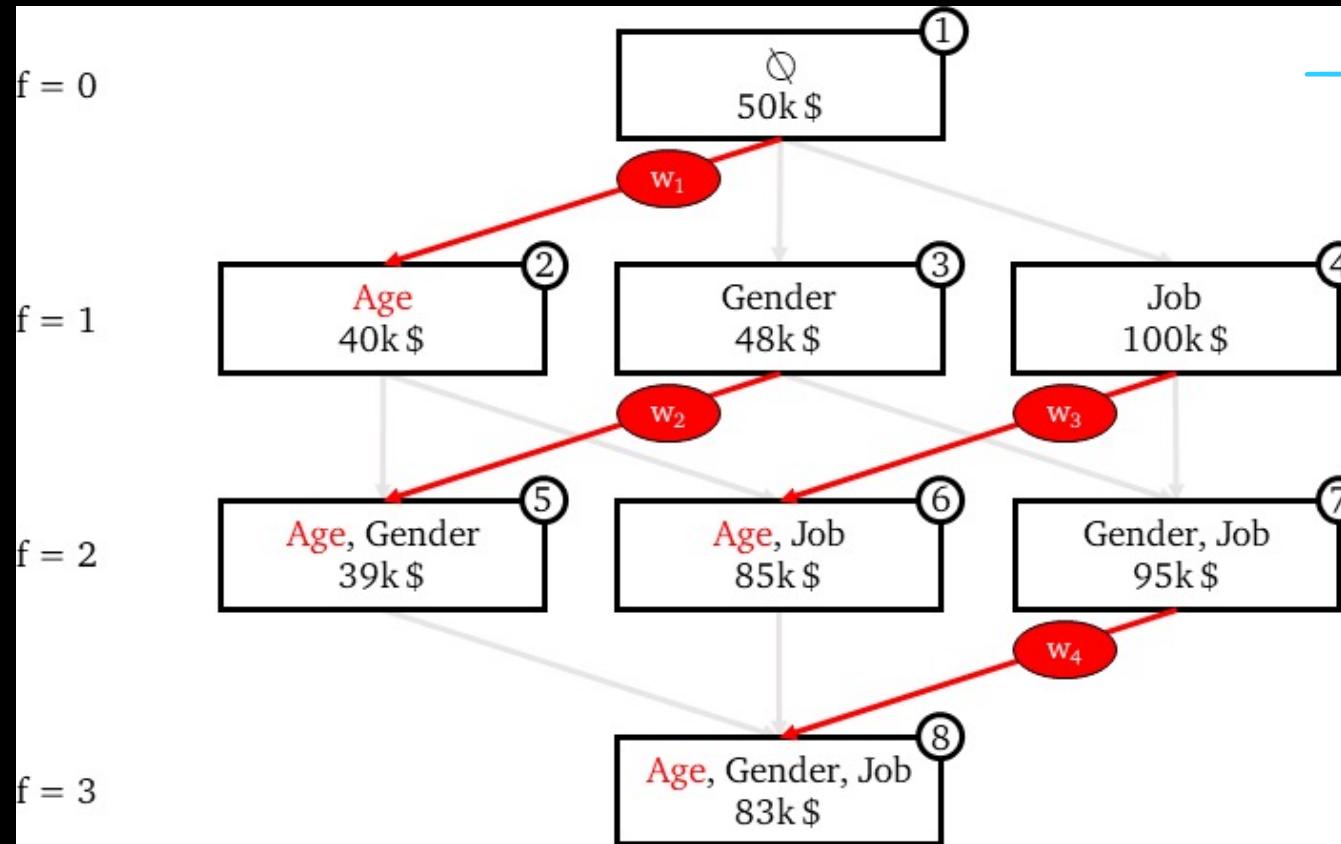
$$SHAP_{\text{Age}}(x_0) = w_1 \times MC_{\text{Age}, \{\text{Age}\}}(x_0) + \\ w_2 \times MC_{\text{Age}, \{\text{Age, Gender}\}}(x_0) + \\ w_3 \times MC_{\text{Age}, \{\text{Age, Job}\}}(x_0) + \\ w_4 \times MC_{\text{Age}, \{\text{Age, Gender, Job}\}}(x_0)$$

$$SHAP_{Age}(x_0) = w_1 \times MC_{Age, \{Age\}}(x_0) + \\ w_2 \times MC_{Age, \{Age, Gender\}}(x_0) + \\ w_3 \times MC_{Age, \{Age, Job\}}(x_0) + \\ w_4 \times MC_{Age, \{Age, Gender, Job\}}(x_0)$$

$$w_1 + w_2 + w_3 + w_4 = 1$$

How to compute the weights?

$$w_1 + w_2 + w_3 + w_4 = 1$$



w_1

$w_2 + w_3$

w_4

w_1

w_4

$$w_2 + w_3 = w_4$$

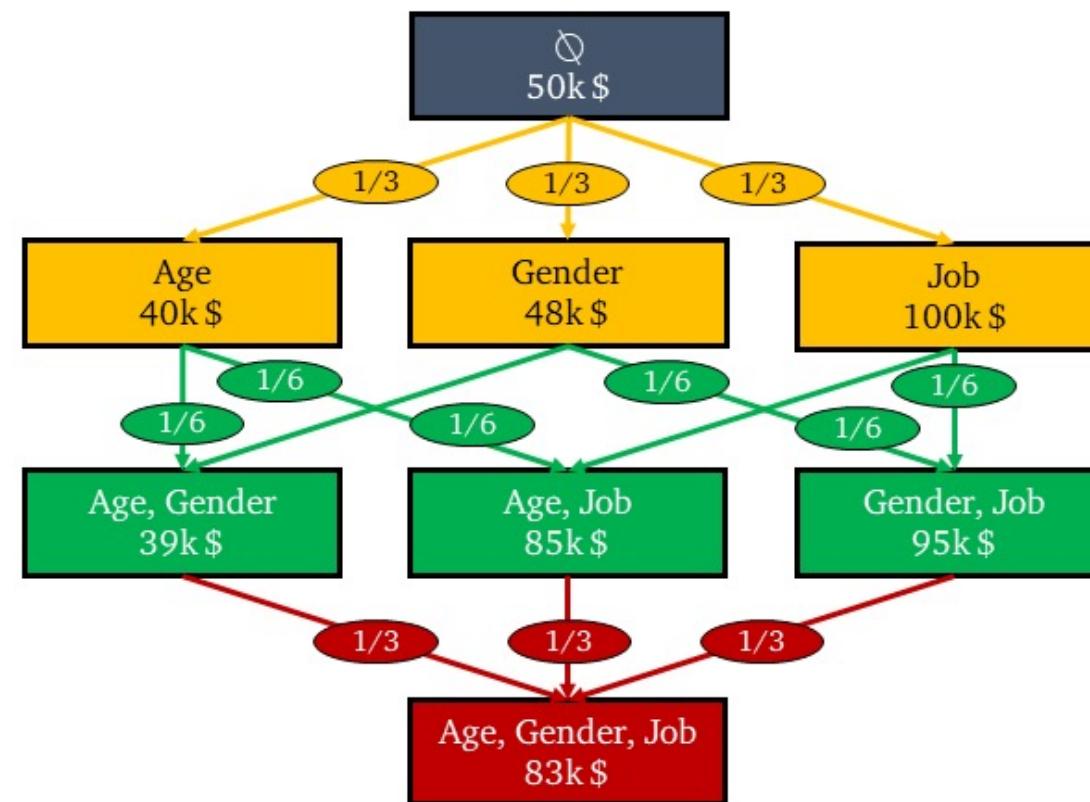
$$\Rightarrow \omega_1 = 1/3$$

$$\omega_2 = \omega_3 = 1/6$$

$$\omega_4 = 1/3$$

F_{Cf}

	N. of Nodes $\binom{F}{f}$	N. of Edges $f \times \binom{F}{f}$
$f = 0$	1	
$f = 1$		3
$f = 2$	3	6
$f = 3$		3
Sum	$2^F = 8$	$F \times 2^{F-1} = 12$



$$\begin{aligned}
 SHAP_{Age}(x_0) &= [(1 \times \binom{3}{1})^{-1} \times MC_{Age, \{Age\}}(x_0) + \\
 &\quad [(2 \times \binom{3}{2})^{-1} \times MC_{Age, \{Age, Gender\}}(x_0) + \\
 &\quad [(2 \times \binom{3}{2})^{-1} \times MC_{Age, \{Age, Job\}}(x_0) + \\
 &\quad [(3 \times \binom{3}{3})^{-1} \times MC_{Age, \{Age, Gender, Job\}}(x_0) + \\
 &= \frac{1}{3} \times (-10k\$) + \frac{1}{6} \times (-9k\$) + \frac{1}{6} \times (-15k\$) + \frac{1}{3} \times (-12k\$) \\
 &= -11.33k\$
 \end{aligned}$$

$$SHAP_{feature}(x) = \sum_{set: feature \in set} [|set| \times \binom{F}{|set|}]^{-1} [Predict_{set}(x) - Predict_{set \setminus feature}(x)]$$

<https://towardsdatascience.com/shap-explained-the-way-i-wish-someone-explained-it-to-me-ab81cc69ef30>

SHapely Additive exPlanations → SHAP

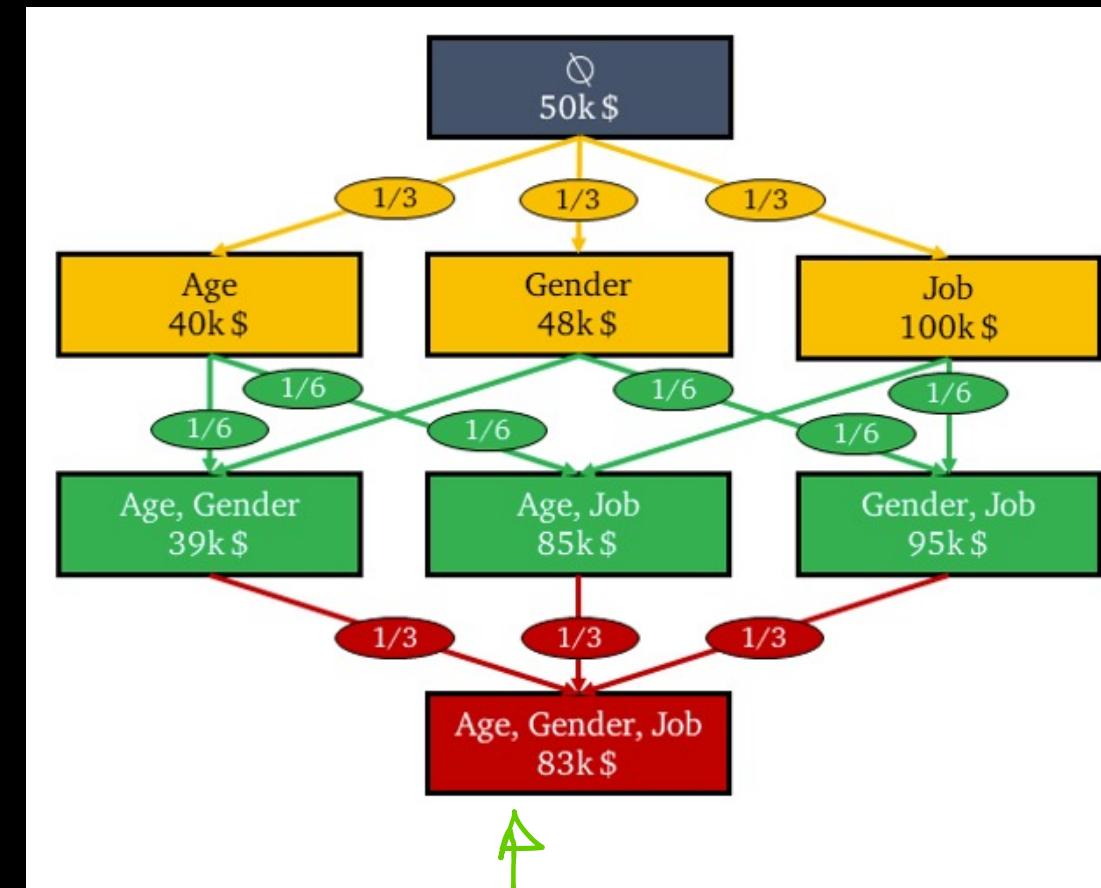
$$SHAP_{Age}(x_0) = -11.33k$$

$$SHAP_{Gender}(x_0) = -2.33k$$

$$SHAP_{Job}(x_0) = +46.66k$$

$$SUM = \frac{33k}{}$$

83-50



sum of shapely values

$$= \hat{y}_F - \hat{y}_{\phi} \rightarrow \text{no-features}$$

↓
all features

Setup of explainability \rightarrow same as LIME

$x \rightarrow x' \in \{0,1\}^M$: interpretable representation

$$x = h_x(x')$$

$g(x')$: Surrogate model

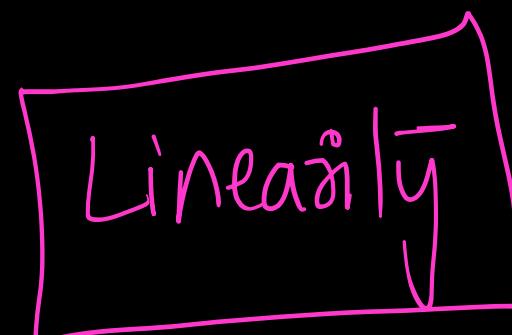
$$g(x') \approx f(h_x(x'))$$

Additive Feature Attribution:

#dim in interpretable representation

$$g(z') = \phi_0 + \sum_{i=1}^M \phi_i z'^i$$

shapely values ← effect of feature i



Expectations in Probability

$$E(x) = \sum_x x \cdot P(x=x) \quad (\text{or}) \quad \int x \cdot P(x) dx$$

\hookrightarrow PMF PDF

\hookrightarrow expected or average value of X

$$E(f(x)) = \sum_x f(x) \cdot P(x=x) \quad (\text{or}) \int_x f(x) P(x) dx$$



Say: X^2

$$E(X|Y=y_1) = \sum_x x \cdot P(X|Y=y_1)$$

Example:

https://en.wikipedia.org/wiki/Conditional_expectation#Example_1:_Dice_rolling

Expectations & Shapely values

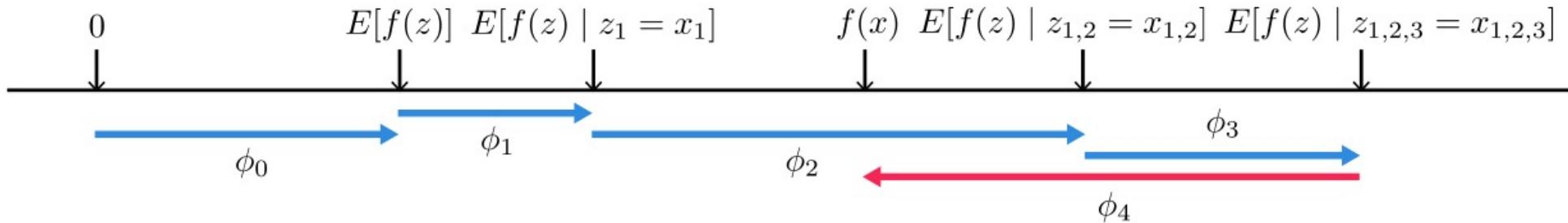


Figure 1: SHAP (SHapley Additive exPlanation) values attribute to each feature the change in the expected model prediction when conditioning on that feature. They explain how to get from the base value $E[f(z)]$ that would be predicted if we did not know any features to the current output $f(x)$. This diagram shows a single ordering. When the model is non-linear or the input features are not independent, however, the order in which features are added to the expectation matters, and the SHAP values arise from averaging the ϕ_i values across all possible orderings.

<https://arxiv.org/pdf/1705.07874.pdf>

Computing shapely values

= computing conditional expectations

Very time consuming exp-time-complexity
[NP-Hard]

Lots of approximate Algorithms

Kernel SHAP \leftarrow LIME + shapley values.

(
as LIME uses a $\pi_{x'}(\cdot)$: local kernel

\rightarrow Model agnostic

Kernel SHAP

$$\pi_{\alpha}(z') = \frac{M-1}{\binom{M}{|z'|} |z'| \binom{M-|z'|}{M-1}}$$

Theorem(2)

features
present in z'



VS LIME

$$\pi_\alpha(z') = \frac{M-1}{\binom{M}{|z'|} |z'| (M-|z'|)}$$

↗ constant

intuition { $|z'|$: v.small or v.large $\Rightarrow \pi_\alpha(z')$: larger
 $|z'|$: in between $\Rightarrow \pi_\alpha(z')$: smaller

Model specific approximate Shapley computation

↳ TreeSHAP (2018)

↳ fast & model specific (RF, GBDT..)

Pros

(+) Solid Theory

(+) LIME + SHAP

(+) Tree-SHAP is fast

Cons

(-) Kernel SHAP is slow

(-) Tree SHAP's unintuitive results

(-) Approximations
(ignores feature dependence)

Code : <https://shap.readthedocs.io/en/latest/index.html>

Q&A