

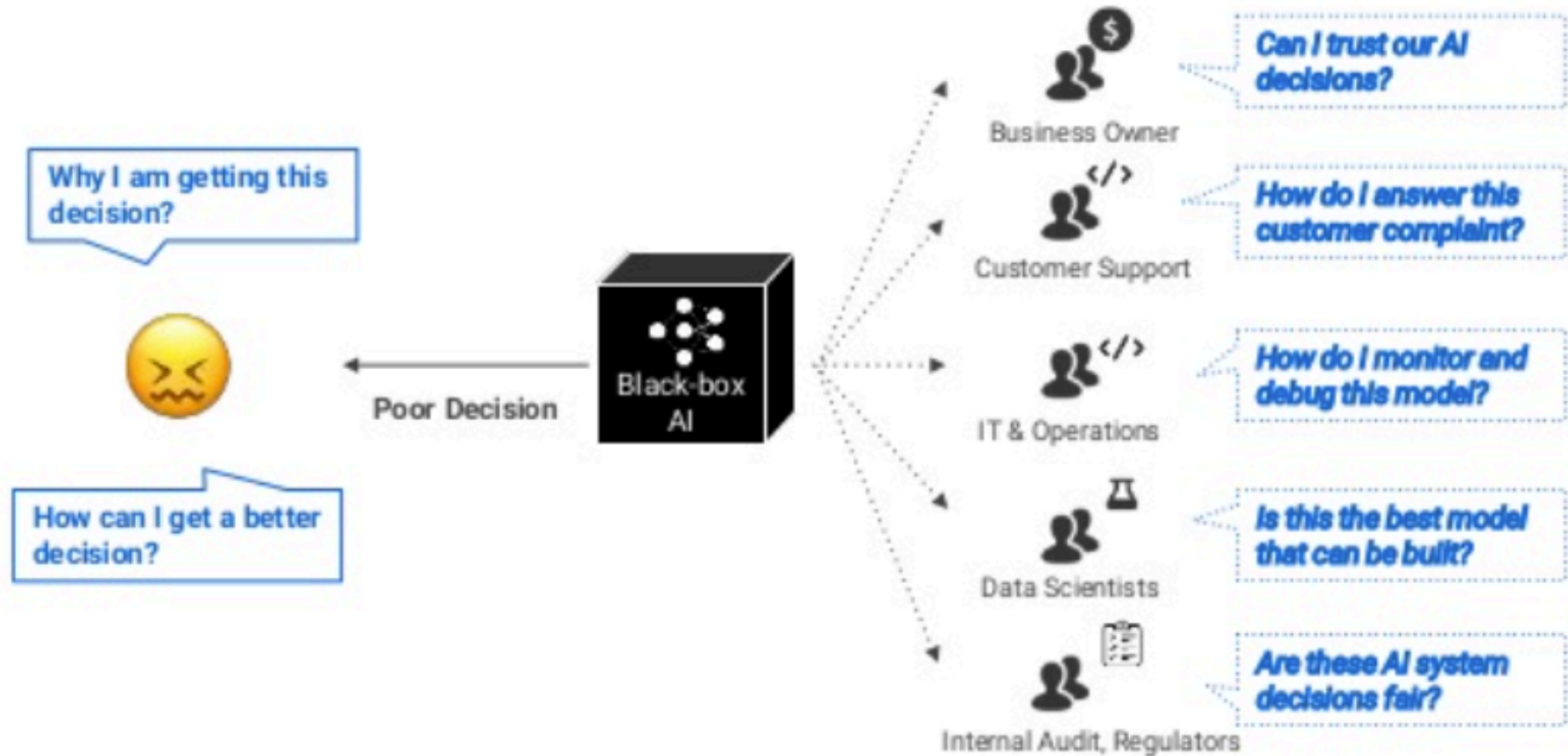


# Interpretabilidade em Modelos de ML

Luciano Barbosa



# Modelos Black Box



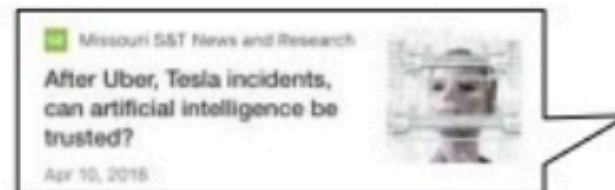


# Viés em Modelos

	
<b>DYLAN FUGETT</b>	<b>BERNARD PARKER</b>
<b>Prior Offense</b> 1 attempted burglary	<b>Prior Offense</b> 1 resisting arrest without violence
<b>Subsequent Offenses</b> 3 drug possessions	<b>Subsequent Offenses</b> None
<b>LOW RISK</b> <b>3</b>	<b>HIGH RISK</b> <b>10</b>
<i>Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.</i>	



# Viés em Modelos





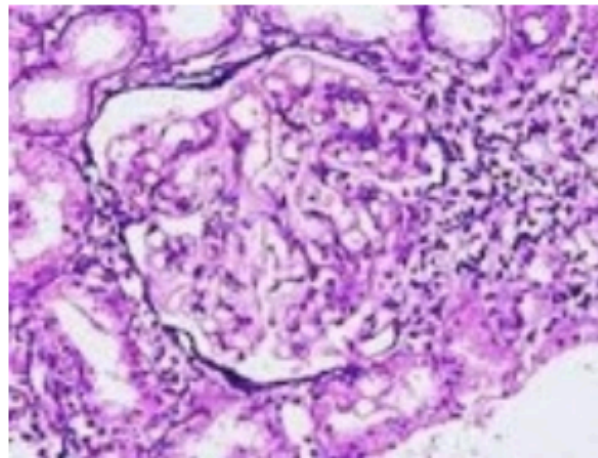
# Para que Interpretabilidade?

- Verificar se o modelo funciona como o esperado

*“Autonomous car crashes, because it wrongly recognizes ...”*



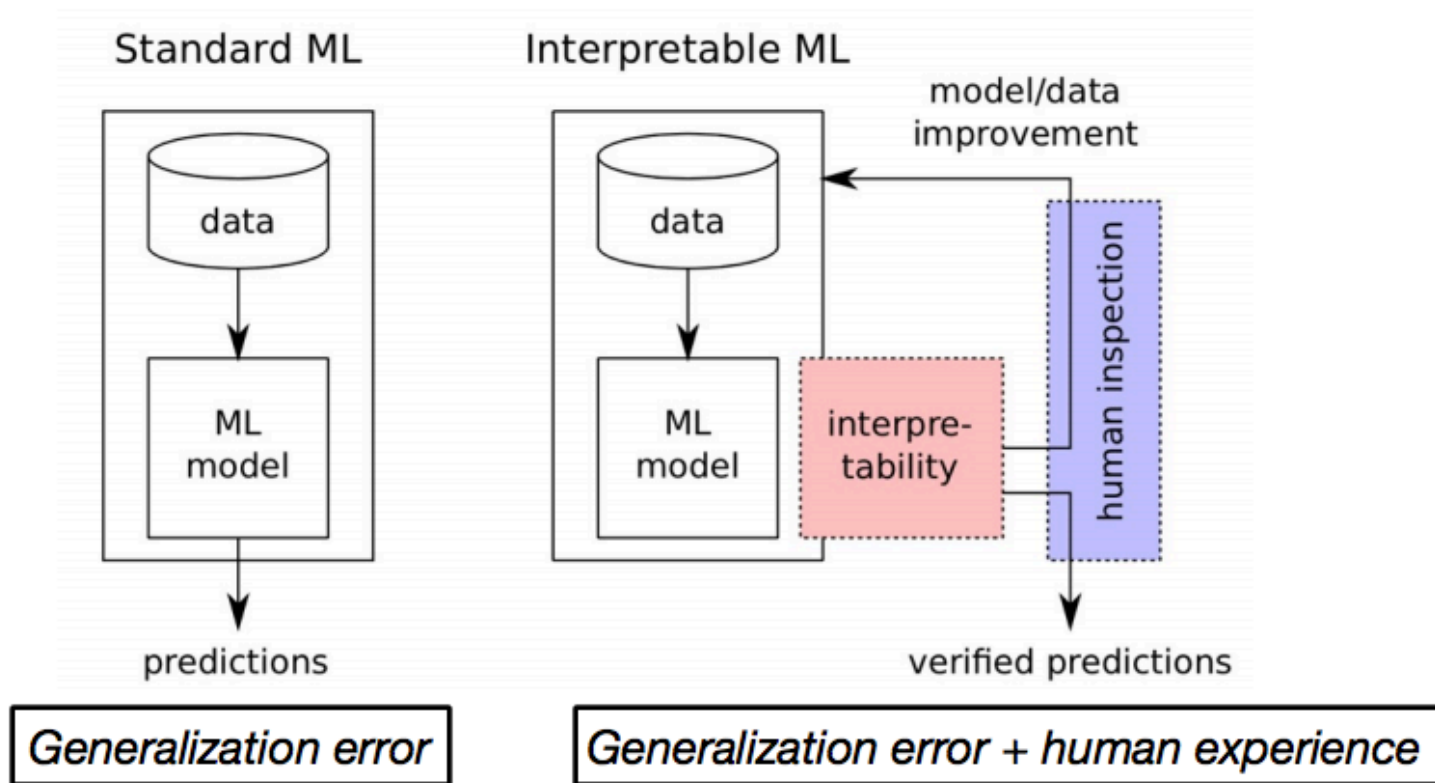
*“AI medical diagnosis system misclassifies patient’s disease ...”*





# Para que Interpretabilidade?

- Verificar problemas e melhorar modelos





# Para que Interpretabilidade?

- Verificar problemas e melhorar modelos



Top label: **"clog"**

Why did the network label this image as **"clog"**?





# Para que Interpretabilidade?

- Aprender novos insights

*“It's not a human move. I've never seen a human play this move.” (Fan Hui)*







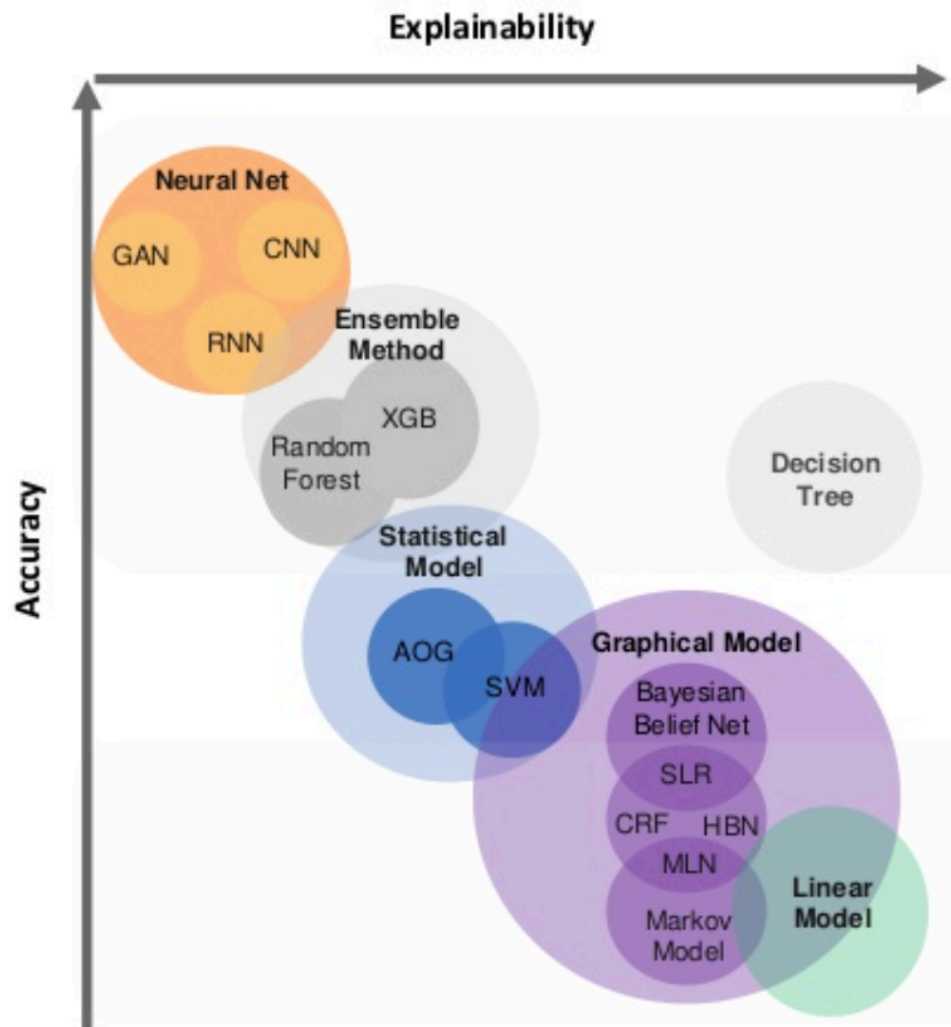
# Para que Interpretabilidade?

- Confirmed with legislation





# Acurácia X Explicabilidade





# Estratégias

- Explicação dado um modelo de ML
  - Predição individual
  - Predição global
- Construção de um modelo interpretável
  - Decision trees, regressão linear etc



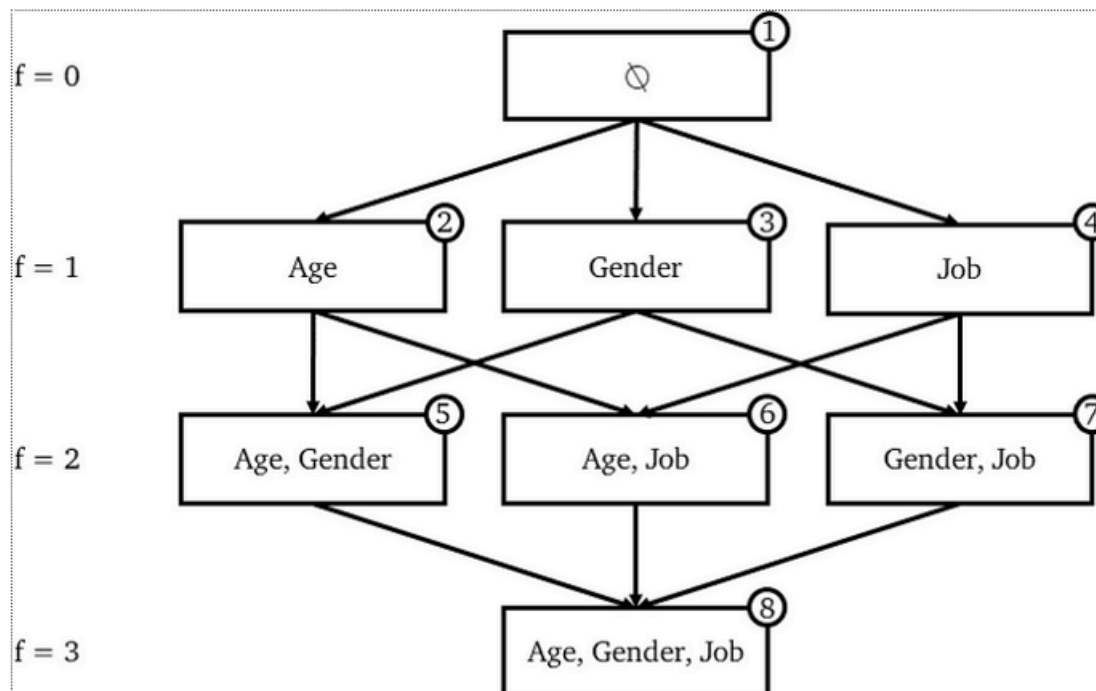
# SHAP: SHapley Additive exPlanations

- Predição individual
- Utiliza teoria dos jogos
- Jogo: saída do modelo
- Jogadores: features do modelo
- Contribuição de cada feature (jogador) para o jogo (predição)



# SHAP: Exemplo

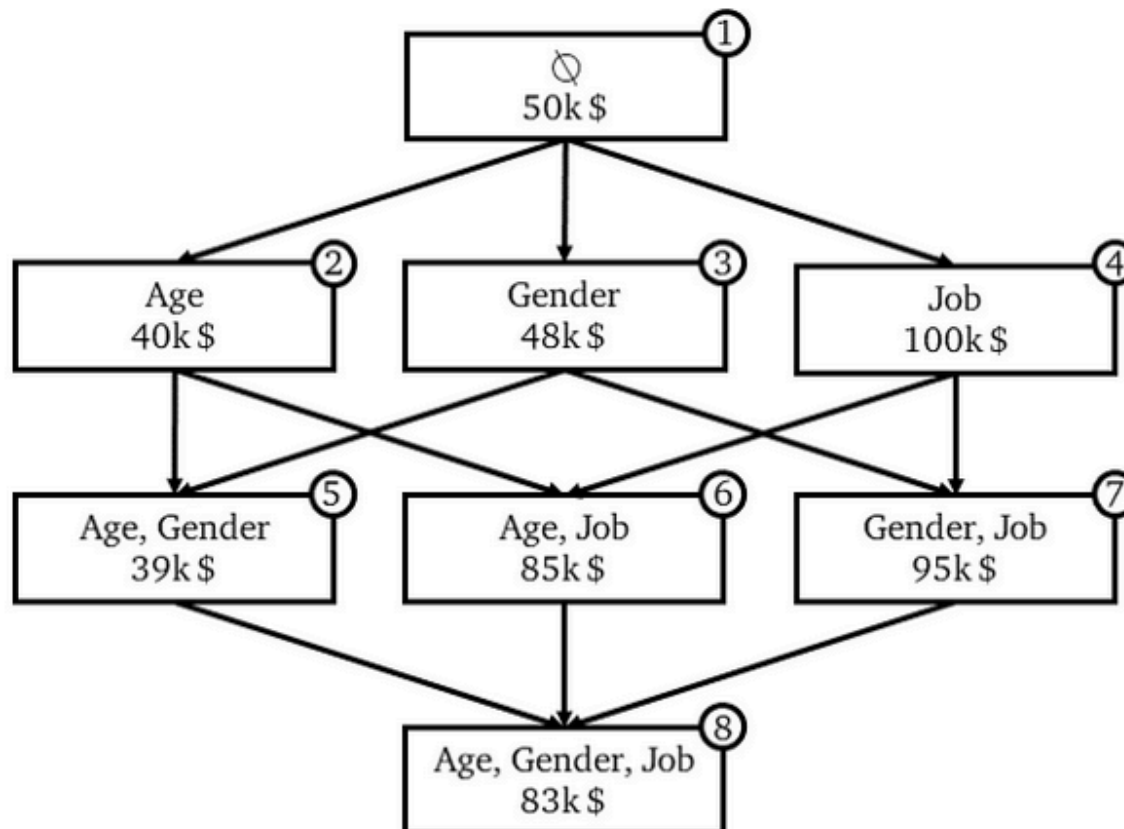
- Cada possível coalizão de jogadores determina a importância de um jogador
- Predizer renda a partir de idade, sexo e emprego





# SHAP: Exemplo

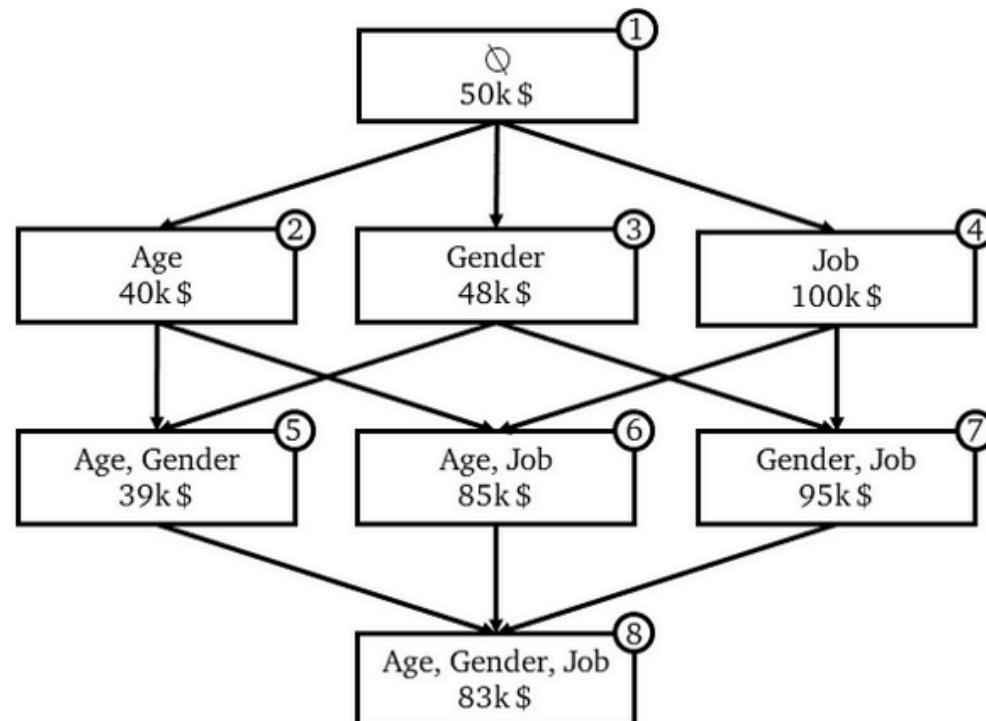
- Treina um modelo para cada coalizão (mesmo conjunto de treinamento e hiper-parâmetros)





# SHAP: Exemplo

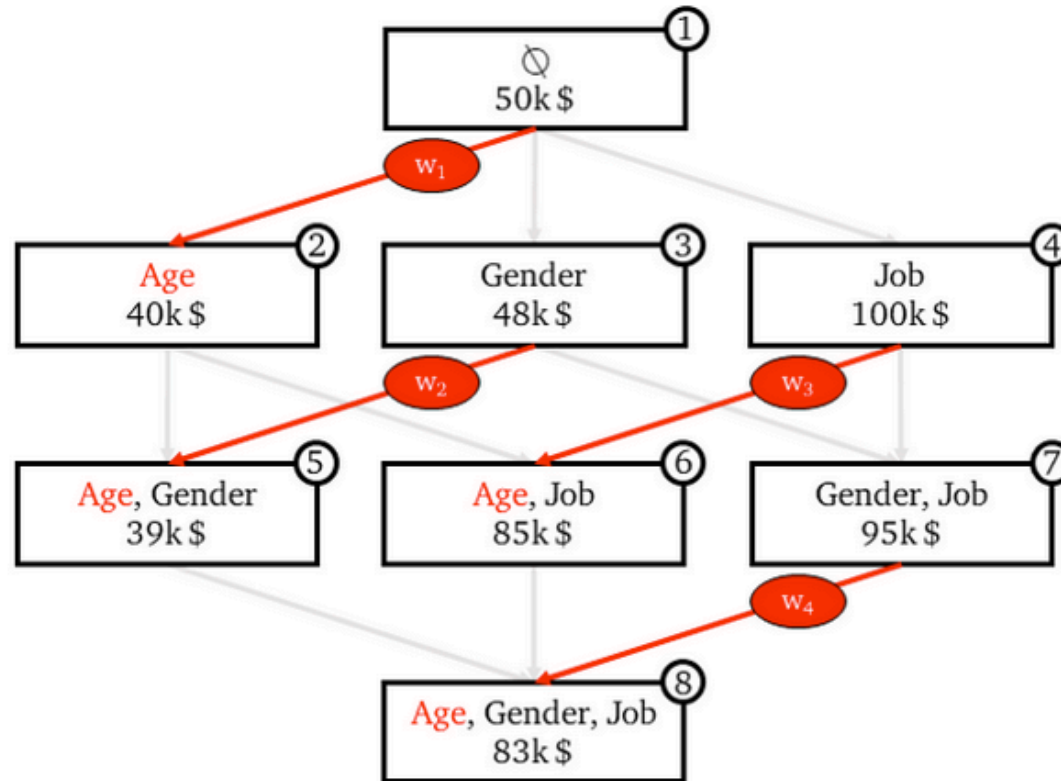
- Contribuição marginal (aresta): diferença entre previsões de dois nós conectados é o efeito da feature adicional







# Contribuição da idade

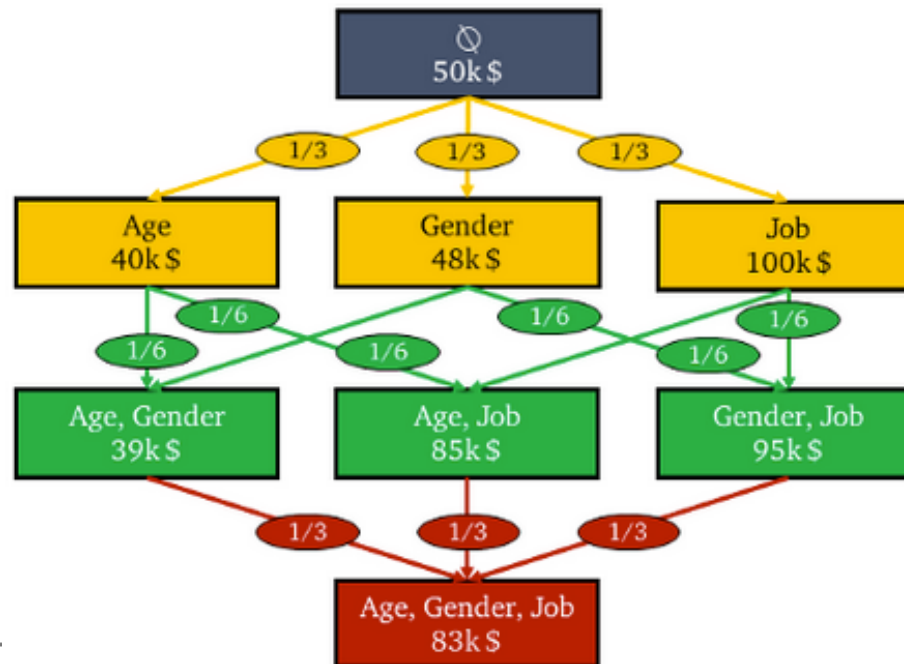


$$\begin{aligned} SHAP_{Age}(x_0) = & w_1 \times MC_{Age, \{Age\}}(x_0) + \\ & w_2 \times MC_{Age, \{Age, Gender\}}(x_0) + \\ & w_3 \times MC_{Age, \{Age, Job\}}(x_0) + \\ & w_4 \times MC_{Age, \{Age, Gender, Job\}}(x_0) \end{aligned}$$



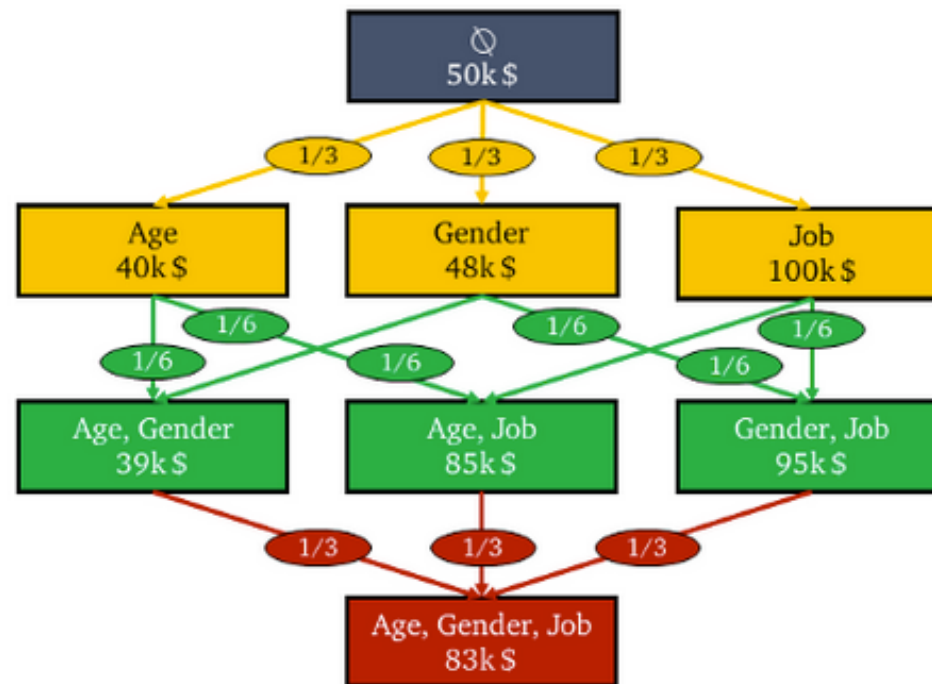
# Pesos

- Soma dos pesos para cada linha tem que ser igual:  $w_1 = w_2 + w_3 = w_4$
- Os pesos na mesma linha são iguais:  $w_1 = 1/3$ ;  $w_2 = 1/6$ ;





# Resultado para Idade



$$\begin{aligned} SHAP_{Age}(x_0) &= \frac{1}{3} \times (-10k\$) + \frac{1}{6} \times (-9k\$) + \frac{1}{6} \times (-15k\$) + \frac{1}{3} \times (-12k\$) \\ &= -11.33k\$ \end{aligned}$$



# Resultado Final

$$\text{SHAP\_Age}(x_0) = -11.33\text{k \$}$$

$$\text{SHAP\_Gender}(x_0) = -2.33\text{k \$}$$

$$\text{SHAP\_Job}(x_0) = +46.66\text{k \$}$$

