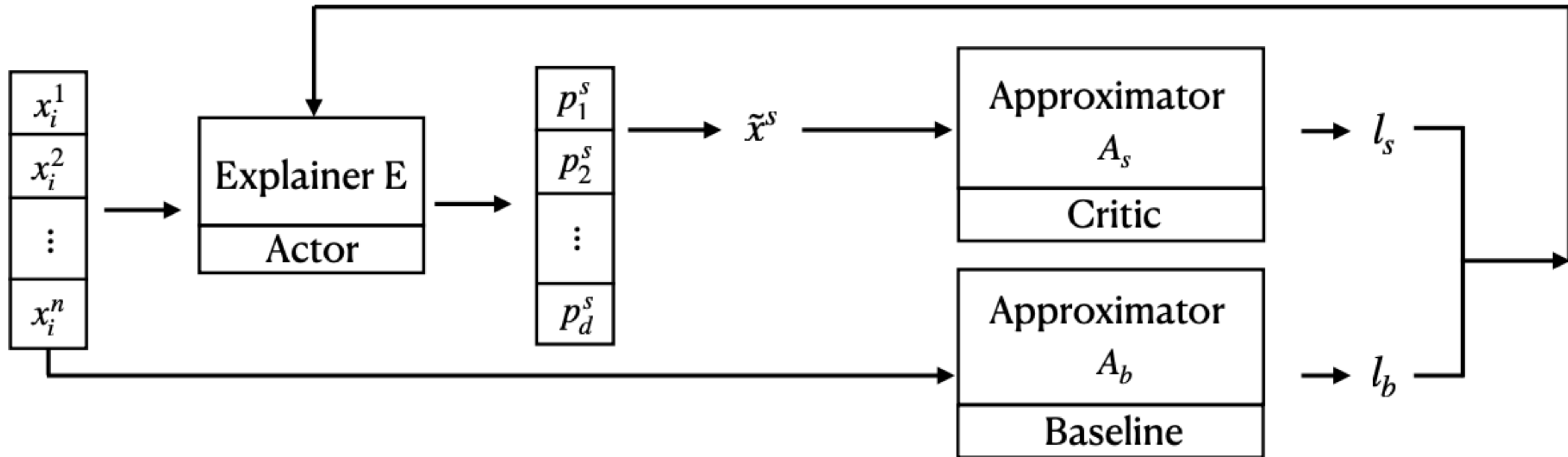




Related Work

**INVASIVE (Actor-Critic Based)**

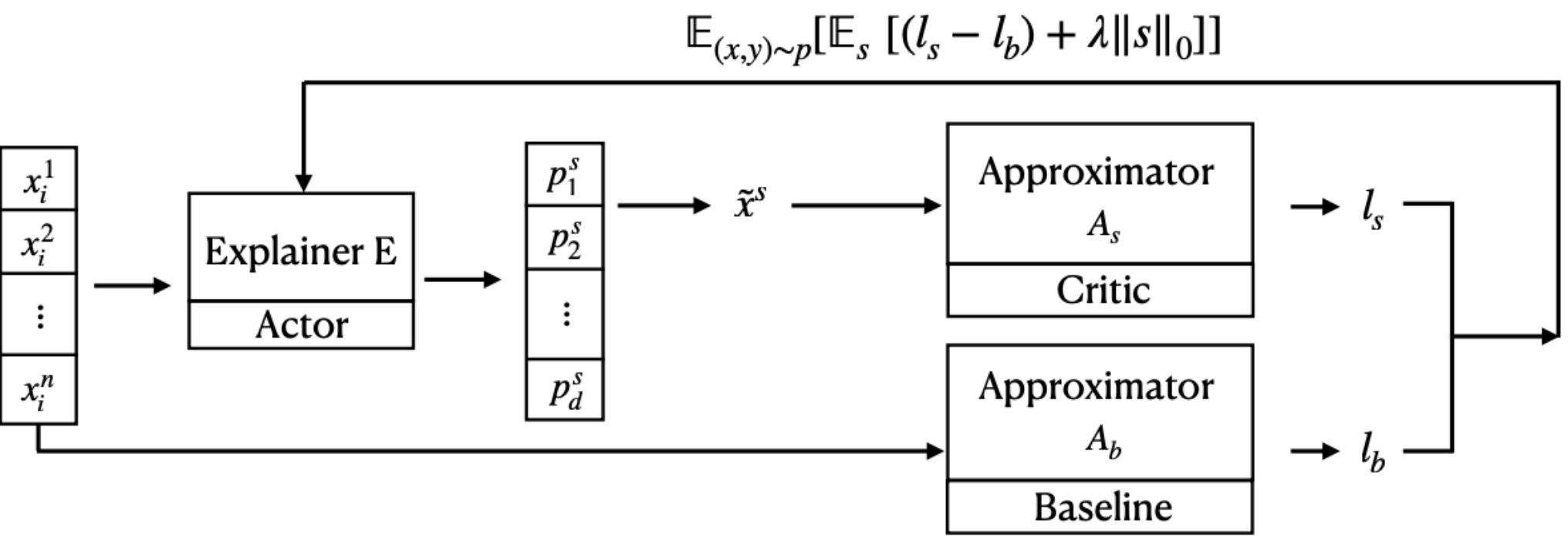
$$\mathbb{E}_{(x,y) \sim p} [\mathbb{E}_s [(l_s - l_b) + \lambda \|s\|_0]]$$



- Actor (Explainer): Neural network Model
- Critic: Neural network Model that use partial feature values to approximate original black-box model
- ✓ Achieve nice performance on synthetic datasets in terms of TPR & FDR

# Related Work

## INVASE (Actor-Critic Based)



# Related Work

## INVASE (Actor-Critic Based)

