

CS480/680: Introduction to Machine Learning

Lec 23: Model Interpretability

Yaoliang Yu

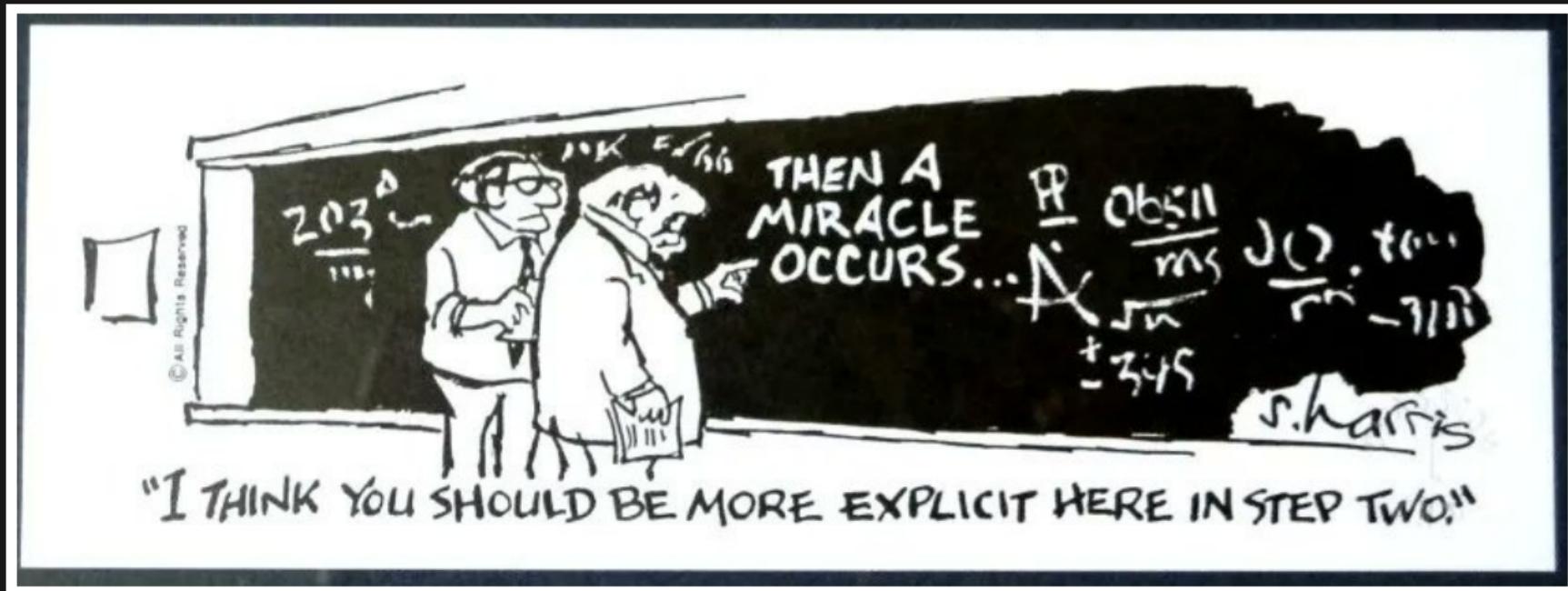


UNIVERSITY OF
WATERLOO

| FACULTY OF MATHEMATICS
DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE

April 3, 2025

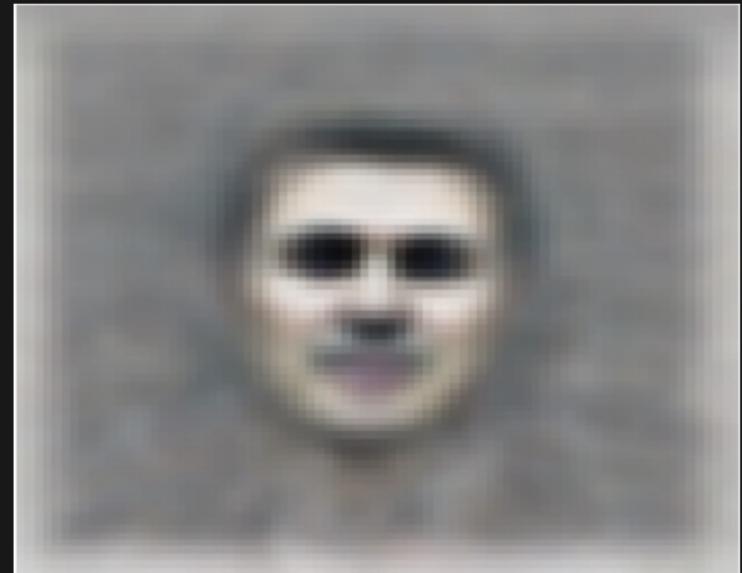
What is Attribution?



- Feature attribution: which features are responsible for this prediction?
- Data valuation: which data point is more valuable for training?

Activation Maximization

- To understand a neuron activation, fix the network weights
- Enumerate test set or run (projected) grad ascent on input



Q. V. Le et al. "Building high-level features using large scale unsupervised learning". In: *Proceedings of the 29th International Conference on Machine Learning*. 2012.

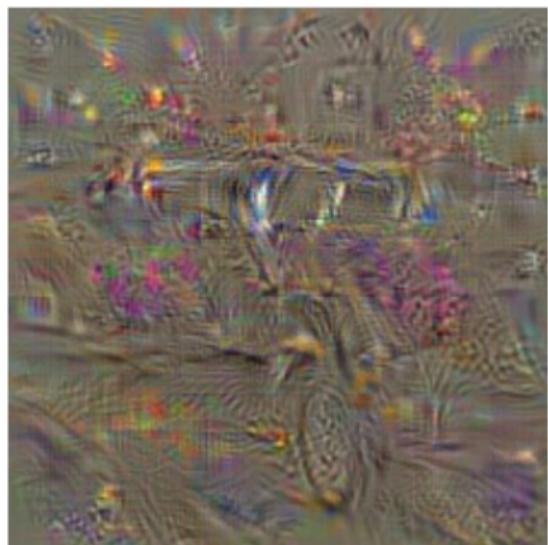
Gradient Saliency



goose



ostrich



limousine

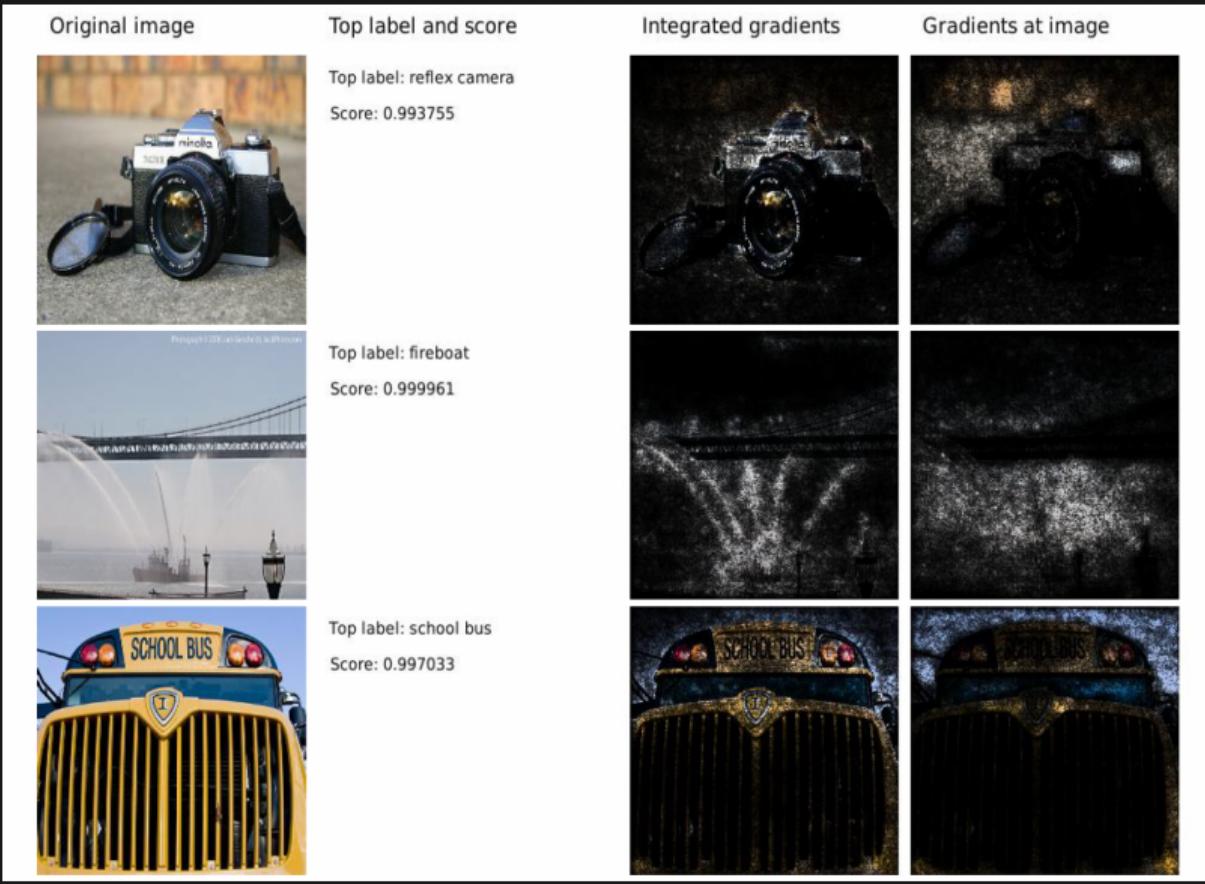
K. Simonyan, A. Vedaldi, and A. Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps". In: *ICLR workshop*. 2017, R. R. Selvaraju et al. "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization". In: *IEEE International Conference on Computer Vision*. 2017, pp. 618–626.

Integrated Gradient

- Fix a baseline \mathbf{z}
- Pick a path $\gamma : [0, 1] \rightarrow \mathbb{R}^d$ from baseline $\mathbf{z} = \gamma(0)$ to data $\mathbf{x} = \gamma(1)$
 - e.g., $\gamma(t) := (1 - t)\mathbf{z} + t\mathbf{x}$
- Integrated gradient:

$$IG(\mathbf{x}; \mathbf{z}) := \int_0^1 \nabla u(\gamma(t)) \odot \gamma'(t) dt$$

- Efficiency: $\text{sum}[IG(\mathbf{x}; \mathbf{z})] = u(\mathbf{x}) - u(\mathbf{z})$



M. Sundararajan, A. Taly, and Q. Yan. "Axiomatic Attribution for Deep Networks". In: *Proceedings of the 34th International Conference on Machine Learning*. 2017, pp. 3319–3328.

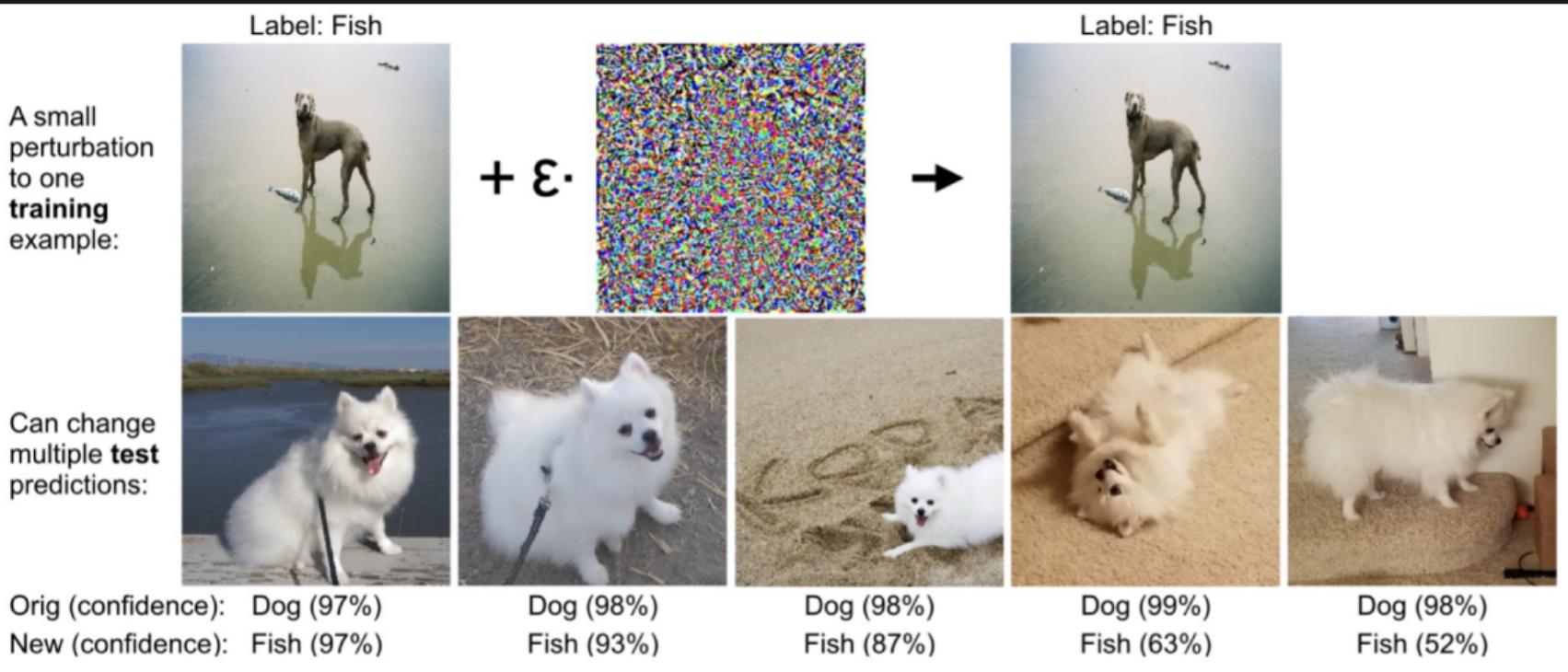
Influence Function

$$\mathbf{w}_\epsilon = \operatorname{argmin}_{\mathbf{w}} \ell(\mathbf{w}) + \epsilon \cdot r(\mathbf{w})$$

- Optimality condition: $\nabla \ell(\mathbf{w}_\epsilon) + \epsilon \cdot \nabla r(\mathbf{w}_\epsilon) = 0$
- Sensitivity: $\frac{d\mathbf{w}_\epsilon}{d\epsilon} \mid_{\epsilon=0} = -[\nabla^2 \ell(\mathbf{w}_0)]^{-1} \cdot \nabla r(\mathbf{w}_0)$
- Chain rule: $\frac{df(\mathbf{w}_\epsilon)}{d\epsilon} \mid_{\epsilon=0} = \nabla f(\mathbf{w}_0) \cdot \frac{d\mathbf{w}_\epsilon}{d\epsilon} \mid_{\epsilon=0} = -\nabla f(\mathbf{w}_0) \cdot [\nabla^2 \ell(\mathbf{w}_0)]^{-1} \cdot \nabla r(\mathbf{w}_0)$
- Good approximation of leave-one-out test:

$$-\ell(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}_i), \quad \epsilon = -\frac{1}{n}, \quad r(\mathbf{w}) = \ell(\mathbf{w}; \mathbf{x}_i)$$

P. W. Koh and P. Liang. “[Understanding black-box predictions via influence functions](#)”. In: *Proceedings of the 34th International Conference on Machine Learning (ICML)*. 2017, pp. 1885–1894.



P. W. Koh and P. Liang. “[Understanding black-box predictions via influence functions](#)”. In: *Proceedings of the 34th International Conference on Machine Learning (ICML)*. 2017, pp. 1885–1894.

OR Example

$$y = x_1 \text{ or } x_2$$

- $n = 2$; consider $x_1 = x_2 = 1$
- Gradient methods do not work
 - fix x_1 , conclude that x_2 does not matter
 - fix x_2 , conclude that x_1 does not matter
 - conclude neither x_1 or x_2 matters ...
- $u(1) = u(2) = u(1, 2) = 1$ and 0 else
- **Banzhaf value:** $p_s \equiv \frac{1}{2^{n-1}} \implies \phi_1 = \phi_2 = \frac{1}{2}$
- **Shapley value:** $p_s = \frac{s!(n-s-1)!}{n!} \implies \phi_1 = \phi_2 = \frac{1}{2}$



Coalition Game

- n is the number of “players”
- $u : 2^{[n]} \rightarrow \mathbb{R}$ the “payoff” function
 - w.l.o.g. $u(\emptyset) = 0$, where \emptyset is the baseline
- Examples:
 - each feature is a player (**feature valuation**)
 - each training example is a player (**data valuation**)
 - each neuron is a player
 - performance metric (e.g. accuracy) is payoff
- **Valuation:** what is the value of each player i ?

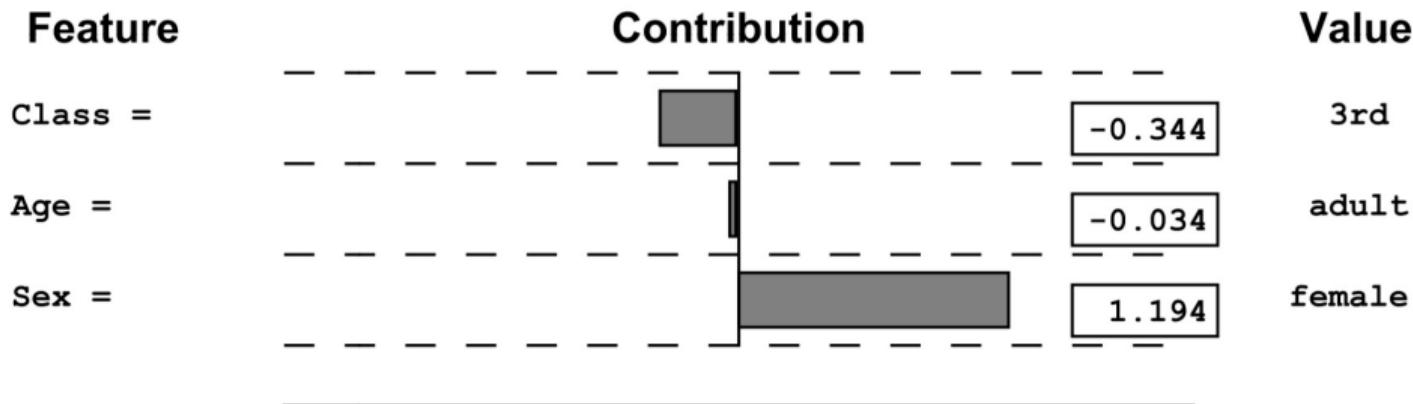
Data: titanic

naive Bayes Explanation

Model: NB

Prediction: $p(\text{survived} = \text{yes} | x) = 0.671$

Actual class label for this instance: yes

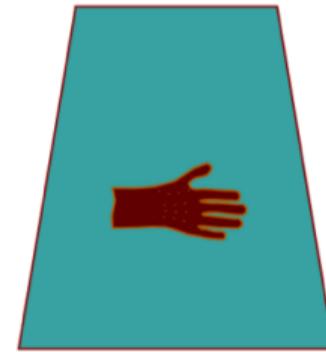
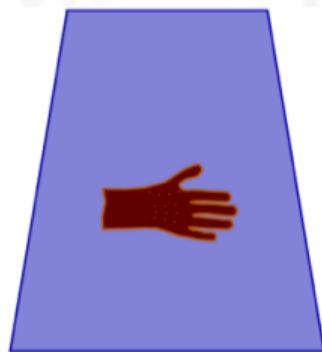
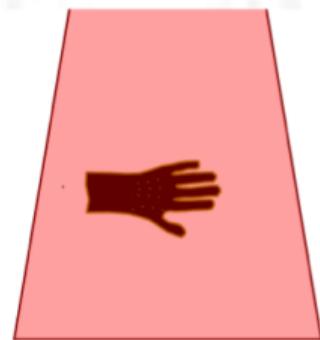


E. Štrumbelj and I. Kononenko. "An Efficient Explanation of Individual Classifications using Game Theory". *Journal of Machine Learning Research*, vol. 11 (2010), pp. 1–18.

Probabilistic Value

- Given a set function $u : 2^{[n]} \rightarrow \mathbb{R}$, e.g.,
 - prediction based on a subset of features
 - accuracy obtained through training on a subset of data
- Find an **additive** approximation $\phi : 2^{[n]} \rightarrow \mathbb{R}$, where $\phi(S) = \sum_{i \in S} \phi(\{i\})$
- **Marginal contribution** of i : $u(S \cup \{i\}) - u(S \setminus \{i\})$
- **Leave-one out**: $u([n]) - u([n] \setminus \{i\})$, i.e., set $S = [n]$
- **(Symmetric) probabilistic value**: $\phi_i^p = \phi^p(\{i\}) = \sum_{S \not\ni i} p_s \cdot [u(S \cup \{i\}) - u(S)]$

R. J. Weber. "Probabilistic values for games". In: *The Shapley Value: Essays in Honor of Lloyd S. Shapley*. Ed. by A. E. Roth. Cambridge University Press, 1988, pp. 101–120.



$$\phi_i = \sum_{S \not\ni i} p_s \cdot [u(S \cup \{i\}) - u(S)]$$

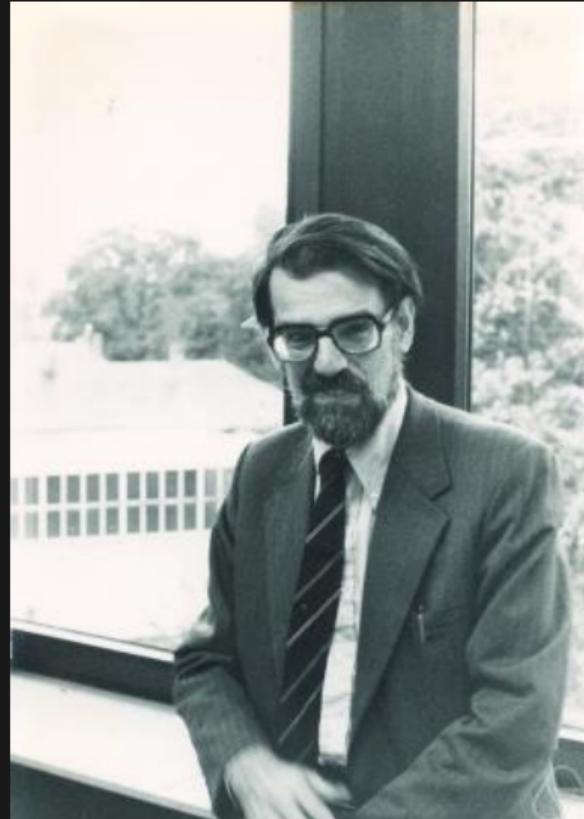
- $n = 3$
- $u(1, 2) = u(1, 3) = u(1, 2, 3) = 1$ and 0 else
- **Banzhaf value:** $p_s \equiv \frac{1}{2^{n-1}} \implies \phi_1 = \frac{1}{4} + \frac{1}{4} + \frac{1}{4}, \quad \phi_2 = \phi_3 = \frac{1}{4}$
- **Shapley value:** $p_s = \frac{s!(n-s-1)!}{n!} \implies \phi_1 = \frac{1}{6} + \frac{1}{6} + \frac{1}{3}, \quad \phi_2 = \phi_3 = \frac{1}{6}$

J. F. Banzhaf III. "Weighted Voting Doesn't Work: A Mathematical Analysis". *Rutgers Law Review*, vol. 19 (1965), pp. 317–343.

L. S. Shapley. "A Value for n-person Games". In: *Contributions to the Theory of Games*. Vol. 2. 1953, pp. 307–318.



John Francis Banzhaf III



Lloyd Stowell Shapley, Nobel Prize (2012)

One Man, One Vote?

<i>State Name (1)</i>	<i>Population 1960 Census</i>	<i>Electoral Vote 1964</i>	<i>Relative Voting Power (2)</i>	<i>Percent Voting Power (3)</i>	<i>Percent Deviation From Average Voting Power (4)</i>	<i>Mississippi</i>	<i>2178141.</i>	<i>7</i>	<i>1.392</i>	<i>39.2</i>	<i>-17.3</i>
Alabama	3266740.	10	1.632	63.2	-3.0	Missouri	4319813.	12	1.710	71.0	1.6
Alaska	226167.	3	1.838	83.8	9.2	Montana	674767.	4	1.421	42.1	-15.5
Arizona	1302161.	5	1.281	28.1	-23.9	Nebraska	1411330.	5	1.231	23.1	-26.9
Arkansas	1786272.	6	1.315	31.5	-21.9	Nevada	285278.	3	1.636	63.6	-2.8
California	15717204.	40	3.162	216.2	87.9	New Hampshire	606921.	4	1.499	49.9	-10.9
Colorado	1753947.	6	1.327	32.7	-21.1	New Jersey	6066782.	17	2.063	106.3	22.6
Connecticut	2535234.	8	1.477	47.7	-12.2	New Mexico	951023.	4	1.197	19.7	-28.9
Delaware	446292.	3	1.308	30.8	-22.3	New York	16782304.	43	3.312	231.2	96.8
Dist. of Columbia	763956.	3	1.000	.0	-40.6	North Carolina	4556155.	13	1.807	80.7	7.4
Florida	4951560.	14	1.870	87.0	11.1	North Dakota	632446.	4	1.468	46.8	-12.8
Georgia	3943116.	12	1.789	78.9	6.3	Ohio	9706397.	26	2.539	153.9	50.9
Hawaii	632772.	4	1.468	46.8	-12.8	Oklahoma	2328284.	8	1.541	54.1	-8.4
Idaho	667191.	4	1.429	42.9	-15.1	Oregon	1768687.	6	1.321	32.1	-21.5
Illinois	10081158.	26	2.491	149.1	48.0	Pennsylvania	11319366.	29	2.638	163.8	56.8
Indiana	4662498.	13	1.786	78.6	6.1	Rhode Island	859488.	4	1.259	25.9	-25.2
Iowa	2757537.	9	1.596	59.6	-5.2	South Carolina	2382594.	8	1.524	52.4	-9.5
Kansas	2178611.	7	1.392	39.2	-17.3	South Dakota	680514.	4	1.415	41.5	-15.9
Kentucky	3038156.	9	1.521	52.1	-9.6	Tennessee	3567089.	11	1.721	72.1	2.3
Louisiana	3257022.	10	1.635	63.5	-2.9	Texas	9579677.	25	2.452	145.2	45.7
Maine	969265.	4	1.186	18.6	-29.5	Utah	890627.	4	1.237	23.7	-26.5
Maryland	3100689.	10	1.675	67.5	-.4	Vermont	389881.	3	1.400	40.0	-16.8
Massachusetts	5148578.	14	1.834	83.4	9.0	Virginia	3966949.	12	1.784	78.4	6.0
Michigan	7823194.	21	2.262	126.2	34.4	Washington	2853214.	9	1.569	56.9	-6.8
Minnesota	3413864.	10	1.597	59.7	-5.1	West Virginia	1860421.	7	1.506	50.6	-10.5
						Wisconsin	3951777.	12	1.788	78.8	6.2
						Wyoming	330066.	3	1.521	52.1	-9.6

J. F. Banzhaf III. "One Man, 3.312 Votes: A Mathematical Analysis of the Electoral College". *Villanova Law Review*, vol. 13, no. 2 (1968), pp. 304–332.

Shapley's Axioms

- **Linear:** $\phi_i(u + v) = \phi_i(u) + \phi_i(v)$
- **Symmetry:** if $u(S \cup i) = u(S \cup j)$ for all S with $i, j \notin S$, then $\phi_i = \phi_j$
- **Null:** if $u(S \cup i) = u(S)$ for all S with $i \notin S$, then $\phi_i = 0$
- **Efficient:** $\sum_i \phi_i = u([n])$

$$\phi_i(u) = \sum_{S \not\ni i} \frac{s!(n-s-1)!}{n!} \cdot [u(S \cup \{i\}) - u(S)]$$

The Power of Linearity

$$\mathbb{R}^n \ni \phi = \begin{matrix} \text{green bar} \\ | \\ \mathbb{R}^{n \times 2^n} \\ | \\ \text{red dotted grid} \\ | \\ u \in \mathbb{R}^{2^n} \end{matrix} \times \begin{matrix} \text{blue bar} \end{matrix}$$

M. Grabisch, J.-L. Marichal, and M. Roubens. "Equivalent Representations of Set Functions". *Mathematics of Operations Research*, vol. 25, no. 2 (2000), pp. 157–178.

Random Order Value

- Let π be a permutation of $[n] := \{1, 2, \dots, n\}$
- Suppose $i = \pi(k)$ and define

$$\psi_i(u, \pi) = u[\underbrace{\pi(1), \dots, \pi(k)}_{\text{when } i \text{ joins}}] - u[\underbrace{\pi(1), \dots, \pi(k-1)}_{\text{before } i \text{ joins}}]$$

- Randomize over permutations: $\boxed{\phi_i(u) = \mathbb{E}_\pi \psi_i(u, \pi)}$
- What happens if we sum all values?

$$\sum_i \phi_i(u) = ?$$

R. J. Weber. "Probabilistic values for games". In: *The Shapley Value: Essays in Honor of Lloyd S. Shapley*. Ed. by A. E. Roth. Cambridge University Press, 1988, pp. 101–120.

How to Estimate Probabilistic Value?

$$\phi_i^p = \phi^p(\{i\}) = \sum_{S \not\ni i} p_s \cdot [u(S \cup i) - u(S)]$$

Algorithm 1: Monte Carlo estimation of probabilistic value

Input: utility u , probability p

```
1 for  $i = 1, \dots, n$  do
2    $\varphi_i \leftarrow 0$ 
3   for  $k = 1, \dots, m$  do
4     sample a random subset  $S \not\ni i$  with probability  $\propto \binom{n-1}{s} p_s$ 
5      $\varphi_i \leftarrow \varphi_i + [u(S \cup \{i\}) - u(S)]$  // 2 evals of utility
6    $\hat{\phi}_i \leftarrow \varphi_i/m$ 
```

X. Deng and C. H. Papadimitriou. "On the Complexity of Cooperative Solution Concepts". *Mathematics of Operations Research*, vol. 19, no. 2 (1994), pp. 257–266.

- Suppose w.l.o.g. utility $u \in [0, 1]$
- $\hat{\phi}_i$ is an average over m i.i.d. samples
- From Hoeffding's inequality: $\Pr[|\hat{\phi}_i - \phi_i| \geq \epsilon] \leq 2 \exp(-m\epsilon^2/2)$
- To achieve $\|\hat{\phi} - \phi\|_\infty \leq \epsilon$ with probability $1 - \delta$, need $O(\frac{n}{\epsilon^2} \log \frac{n}{\delta})$ samples
- Maximum sample reuse for the Banzhaf value: $O(\frac{n}{\epsilon^2} \log \frac{n}{\delta})$ for ℓ_2 norm
 - $\|\hat{\phi} - \phi\|_2 \leq \epsilon$ vs. $\|\hat{\phi} - \phi\|_\infty \leq \epsilon/\sqrt{n}$

Least-square Value

$$\min_{\phi \in \mathbb{R}^n} \sum_{S \subseteq [n]} q_s \cdot [u(S) - \phi(S)]^2 \quad \text{s.t.} \quad u([n]) = \sum_i \phi_i$$

- Take $q_s = p_s + p_{s-1}$ recovers (efficient normalization of) probabilistic value
 - $\sum_{S \not\ni i} p_s [u(S \cup \{i\}) - u(S)] = \sum_{S \ni i} p_{s-1} u(S) - \sum_{S \not\ni i} p_s u(S) = \sum_{S \ni i} [p_{s-1} + p_s] u(S) - \sum_S p_s u(S)$
 - Shapley value corresponds to $q_s = \frac{(s-1)!(n-1-s)!}{(n-1)!} \equiv \frac{1}{\binom{n-2}{s-1}}$
- Shapley value \subseteq Random order value \subseteq Probabilistic value \subseteq Least-square value

A. Charnes, B. Golany, M. S. Keane, and J. J. Rousseau. "Extremal Principle Solutions of Games in Characteristic Function Form: Core, Chebychev and Shapley Value Generalizations". In: *Econometrics of Planning and Efficiency*. 1988, pp. 123–133, L. M. Ruiz, F. Valenciano, and J. M. Zarzuelo. "The Family of Least Square Values for Transferable Utility Games". *Games and Economic Behavior*, vol. 24, no. 1-2 (1998), pp. 109–130.

W. Li and Y. Yu. "Faster Approximation of Probabilistic and Distributional Values via Least Squares". In: *International Conference on Learning Representations (ICLR)*. 2024.



A Simple Observation

$$\begin{aligned}\phi_i &= \sum_{S \ni i} p_s [u(S) - u(S \setminus i)] = \sum_{S \ni i} p_s u(S) - \sum_{S \not\ni i} p_{s+1} u(S) \\ &= \sum_{S \ni i} [p_{s+1} + p_s] u(S) - \sum_S p_{s+1} u(S) \\ &\equiv \sum_{S \ni i} [p_{s+1} + p_s] u(S) \\ &\equiv \mathbb{E}[u(S) \cdot \mathbf{1}_{i \in S}]\end{aligned}$$

Algorithm 2: Maximum sample reuse for Rankings

Input: utility u , probability p and total number of samples T

Output: An estimate $\hat{\mathbf{r}}$ up to some scalar

```
1 define  $\mathbf{q} \in \Delta_n$  by  $q_s \propto \binom{n}{s}(p_s + p_{s+1})$            // distribution over coalition size
2  $\hat{\mathbf{r}} \leftarrow \mathbf{0}_n, \mathbf{t} \leftarrow \mathbf{0}_n$                                 // initialization
3 for  $k = 1, 2, \dots, T$  do
4   sample  $s_k \in [n]$  using  $\mathbf{q}$                                 // sample coalition size
5   uniformly sample  $S_k$  from  $\{S \subseteq [n] \mid |S| = s_k\}$       // sample coalition
6   for  $i \in S_k$  do
7      $t_i \leftarrow t_i + 1$                                          // number of updates for  $i$ 
8    $\hat{r}_i \leftarrow (1 - \frac{1}{t_i})\hat{r}_i + \frac{1}{t_i}u(S_k)$     // averaging
```

- Each utility eval $u(S_k)$ is used for updating $|S_k| = s_k$ values
- Total number of updates t_i for each i is now random

Theorem: Maximum sample reuse for ranking

Achieve (ϵ, δ) guarantee under ℓ_2 norm with $O(\frac{\kappa^2(n)}{\epsilon^2} \log \frac{n}{\delta})$ samples

$$\kappa(n) := \frac{\sum_{s=1}^{n-1} \binom{n}{s} (p_s + p_{s+1})}{\sum_{s=1}^{n-1} \binom{n-1}{s-1} (p_s + p_{s+1})} = \frac{n}{\mathbb{E}[s]} \leq n$$

- The larger the average coalition size $\mathbb{E}[s]$, the smaller $\kappa(n)$: greater recycling
- Worst-case: $\kappa(n) = n$, matching the bound for Monte Carlo
- For common probabilistic values, either $\kappa(n) = \Theta(1)$ or $(\log n)^2$
- Dependence on ϵ is likely optimal; some room for improving $\kappa(n)$?

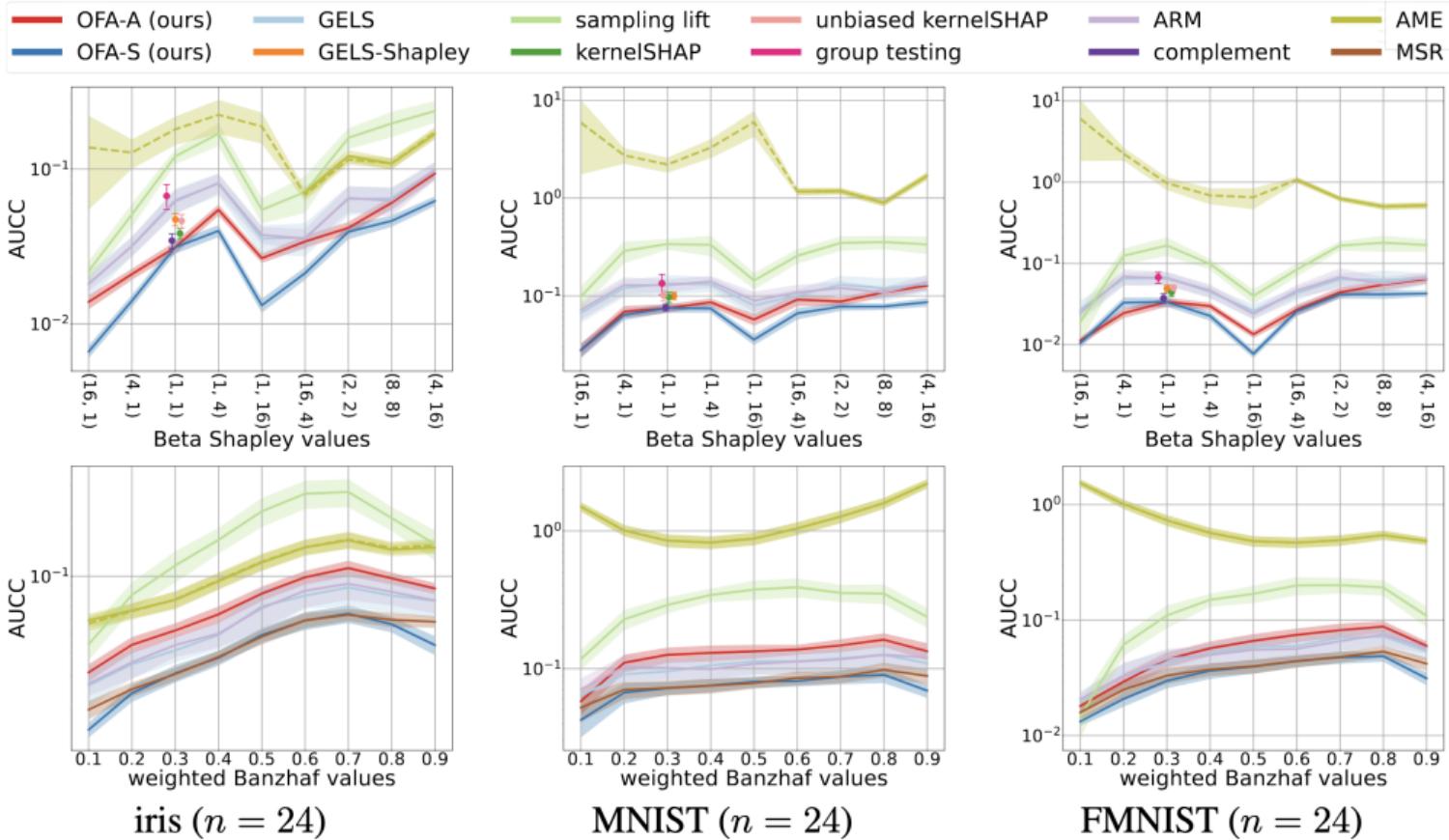
Dummy to the Rescue

- To recover the precise value of ϕ from rankings, we introduce a dummy 😺:

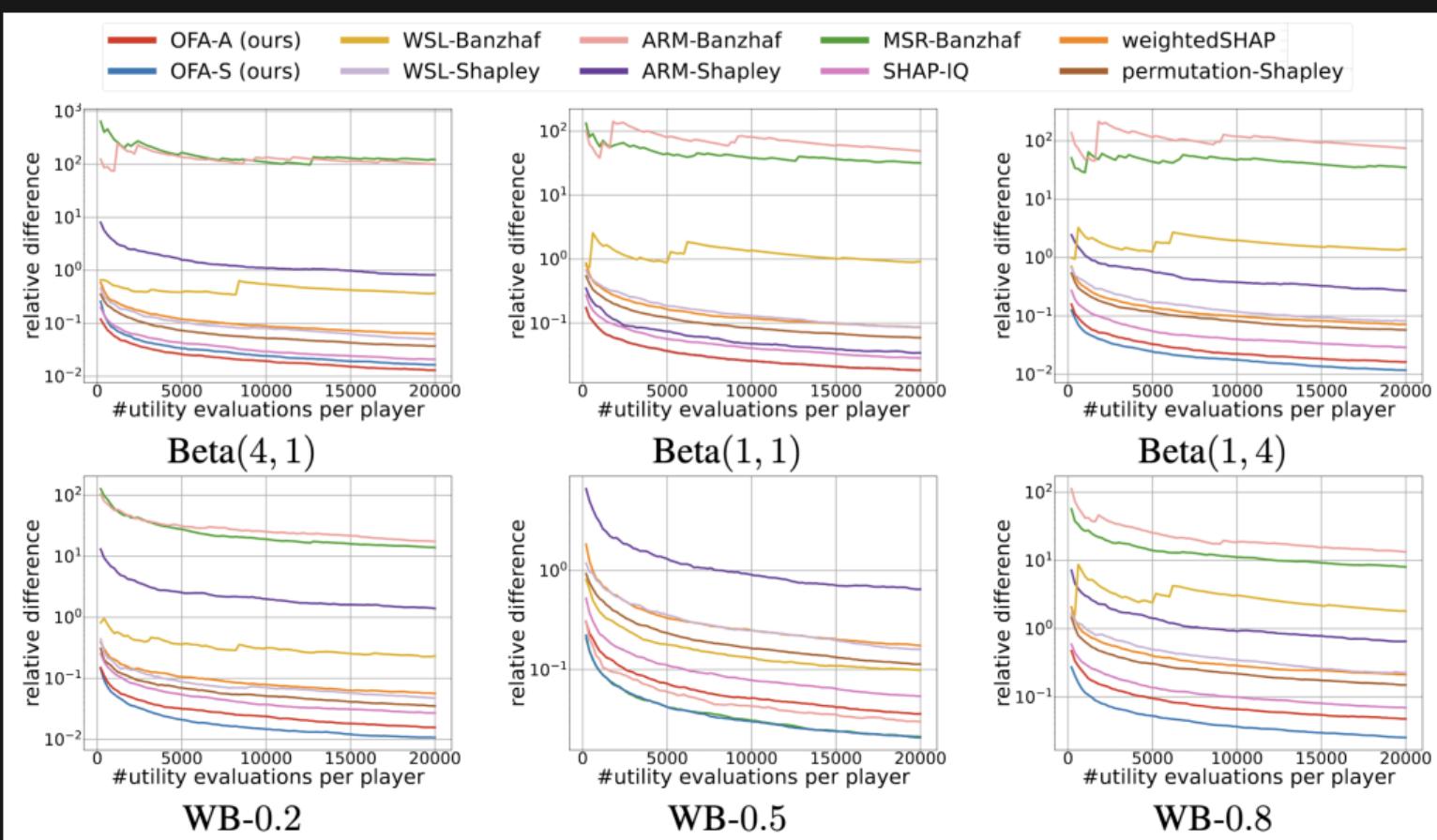
$$v(S) := u(S \setminus \text{😺})$$

- By definition, $\phi_{\text{😺}} = 0$ while ϕ_i remains the same for others
- Conclude: $\phi_i = r_i - r_{\text{😺}}$
- Similar complexity result

- For Shapley: $O(\frac{n(\log n)^2}{\epsilon^2} \log \frac{n}{\delta})$
- vs. $O(\frac{n \log n}{\epsilon^2} \log \frac{n}{\delta})$ vs. $\frac{n}{\epsilon^2} \log \frac{n}{\delta}$ (but $O(n^2)$ memory)



W. Li and Y. Yu. "One Sample Fits All: Approximating All Probabilistic Values Simultaneously and Efficiently". In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2024.



W. Li and Y. Yu. "One Sample Fits All: Approximating All Probabilistic Values Simultaneously and Efficiently". In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2024.