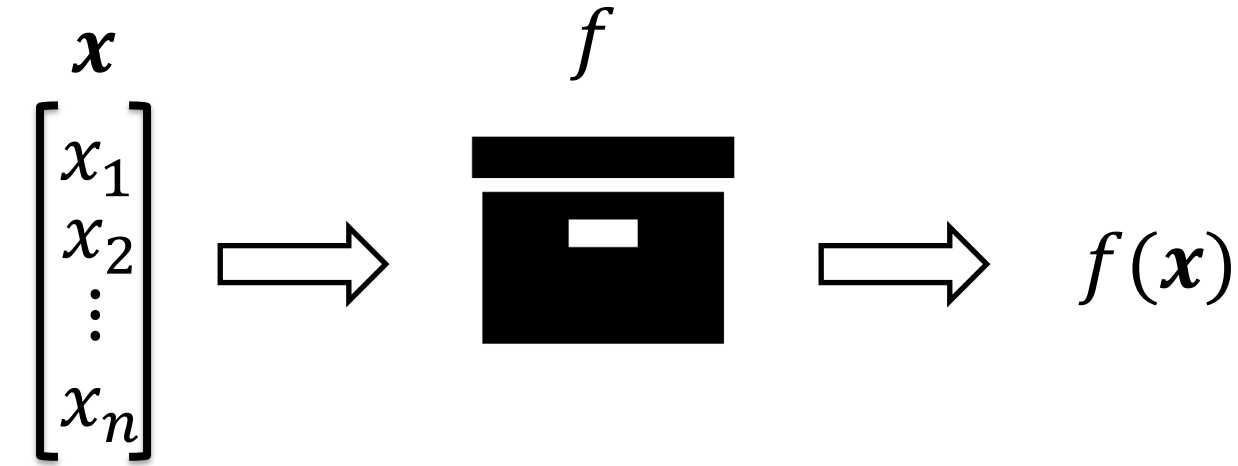


Explainable AI

Goal: Explain the output of an opaque function.



Example: How does changing traffic (an input feature) change the enjoyment of a bike ride (the output)?

	Explicand	Baseline	S	Temperature	Wind speed	Chance of rain	Helicopters	Traffic	f(x ^S)
	73°	89°	{2,3}	89°	11mph	30%	5	3	5/10
	11mph	1mph	{2,3,5}	89°	11mph	30%	5	8	4/10
	30%	0%	{1}	73°	1mph	0%	5	3	6/10
	2	5	{1,5}	73°	1mph	0%	5	8	8/10
	8	3							
	⋮	⋮							

Our answer will attribute the model output to each input using...

Shapley Values

For a set function $v: 2^{[n]} \rightarrow \mathbb{R}$, the i th Shapley value is

$$\phi_i = \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i\}} \frac{v(S \cup \{i\}) - v(S)}{\binom{n-1}{|S|}} = \frac{1}{n} \sum_{k \in [n-1]} \underbrace{\frac{1}{\binom{n-1}{k}} \sum_{S \subseteq [n] \setminus \{i\}: |S|=k} v(S \cup \{i\}) - v(S)}_{\text{Average over sets of size } k} \underbrace{\frac{1}{\binom{n-1}{k}}}_{\text{Average over all sizes } k}$$

We will set $v(S) = f(x^S)$.

Regression Formulation

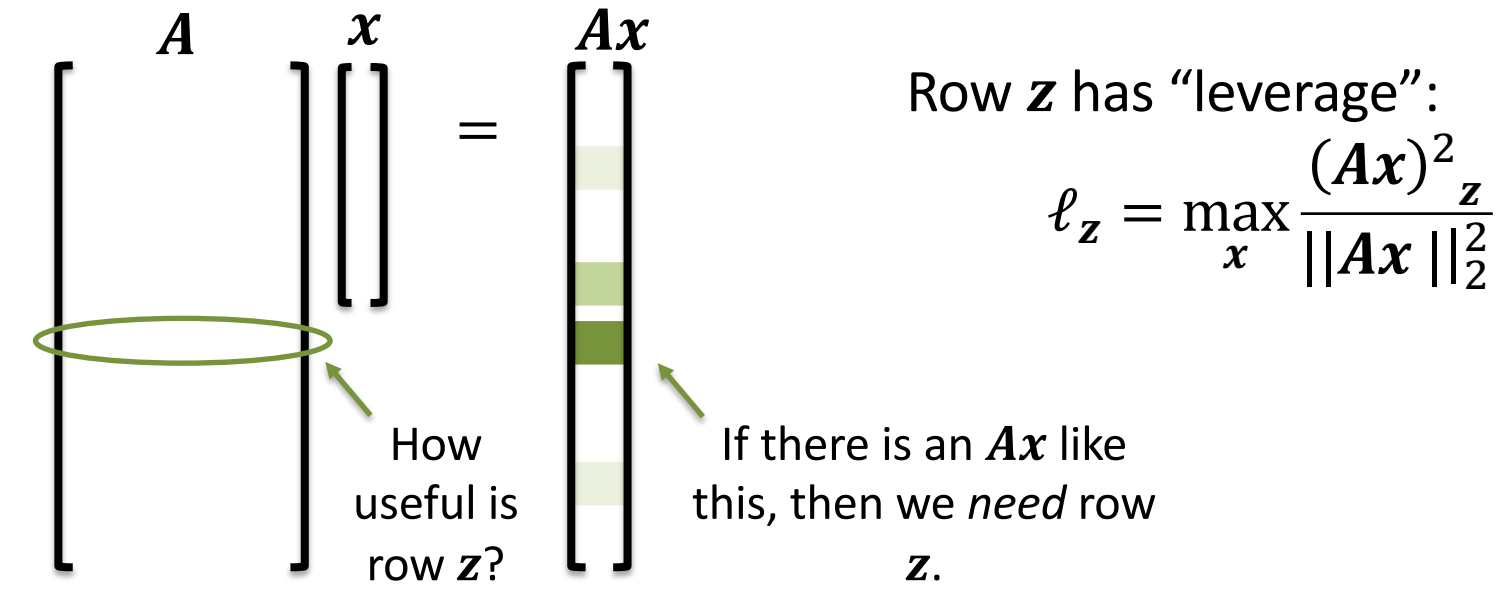
For every v , there is a matrix $A \in \mathbb{R}^{2^n \times n}$ and $b \in \mathbb{R}^{2^n}$ so that

$$\begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_n \end{bmatrix} = \phi = \arg \min_{\beta} \|A\beta - b\|_2 + 1 \frac{v([n]) - v(\emptyset)}{n}$$

We can compute Shapley values by solving a regression problem!

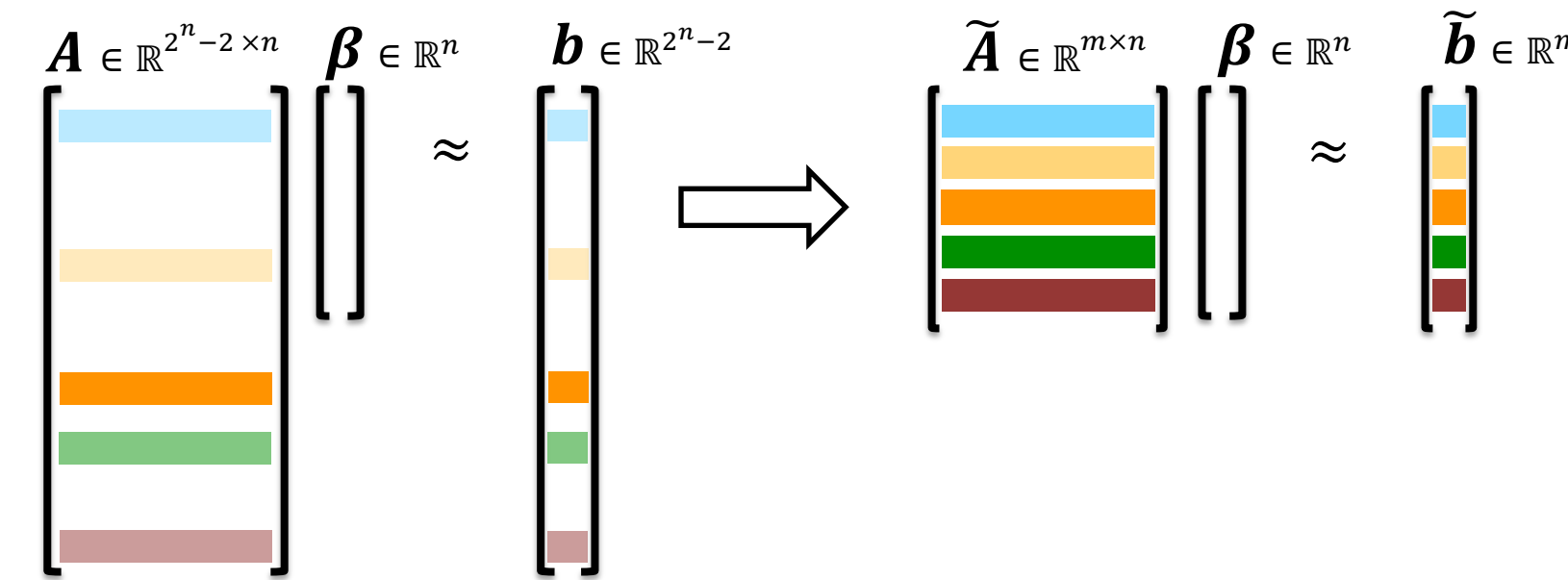
Leverage Score Sampling

If the regression problem is too large, how should we sample?



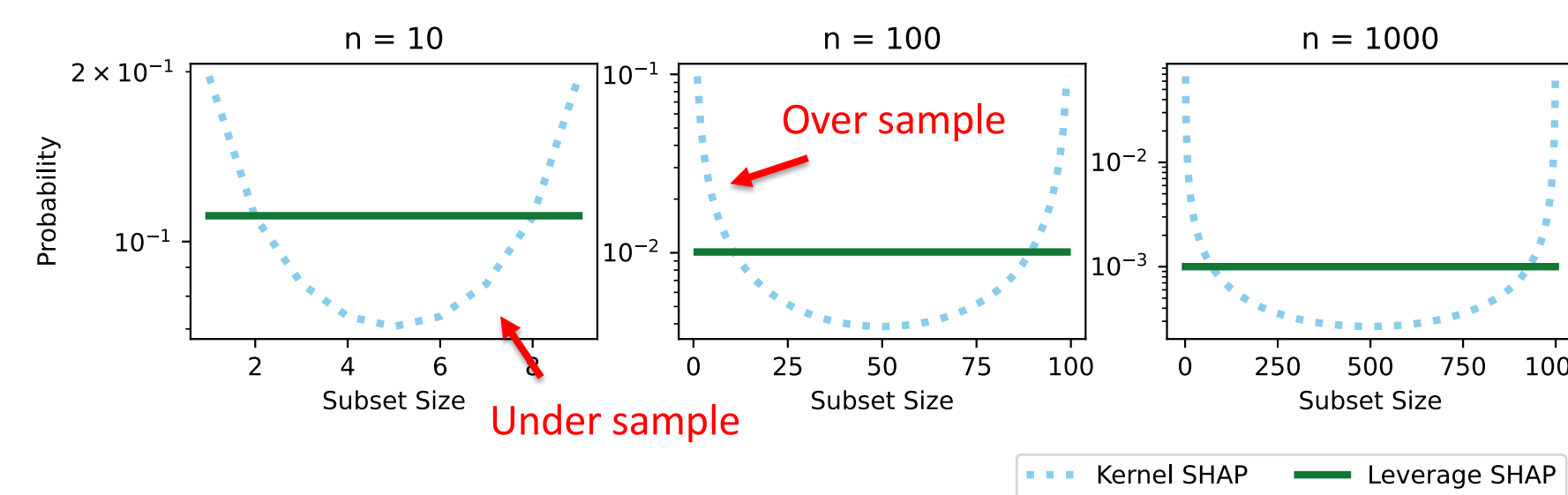
Leverage SHAP

Algorithm: *Leverage SHAP* estimates Shapley values by solving a regression problem sampled via leverage scores $\ell_z = \binom{n}{|z|}^{-1}$.



Leverage SHAP vs Kernel SHAP

Distribution Over Sampled Subset Size



Error Over 100 Runs

	IRIS	California	Diabetes	Adult	Correlated	Independent	NHANES	Communities
Kernel SHAP								
Mean	0.00261	0.0208	15.4	0.000139	0.00298	0.00324	0.0358	130.0
1st Quartile	5.69e-07	0.0031	3.71	1.48e-05	0.00166	0.00163	0.0106	33.5
2nd Quartile	9.52e-06	0.0103	8.19	3.86e-05	0.00249	0.00254	0.0221	53.6
3rd Quartile	0.00181	0.029	20.1	0.000145	0.00354	0.00436	0.0418	132.0
Optimized Kernel SHAP								
Mean	3.28e-14	0.00248	2.33	1.81e-05	0.000739	0.000649	0.00551	21.8
1st Quartile	2.12e-14	0.000279	0.549	2.16e-06	0.00027	0.000187	0.000707	5.85
2nd Quartile	3.55e-14	0.00138	1.26	5.43e-06	0.000546	0.000385	0.0024	13.0
3rd Quartile	4.22e-14	0.0036	3.03	1.63e-05	0.00101	0.000964	0.00665	25.1
Leverage SHAP								
Mean	3.28e-14	0.000186	0.63	5.21e-06	0.000458	0.000359	0.00385	14.7
1st Quartile	2.12e-14	1.91e-05	0.0631	6.3e-07	0.000139	9.51e-05	0.000333	3.6
2nd Quartile	3.55e-14	8.31e-05	0.328	2.33e-06	0.000376	0.000235	0.00149	8.9
3rd Quartile	4.22e-14	0.000231	0.769	7.09e-06	0.000617	0.000556	0.00401	15.3

Theoretical Guarantees

With $O\left(n \log n + \frac{n}{\epsilon}\right)$ samples, the Leverage SHAP solution $\tilde{\phi}$, with probability 9/10, satisfies

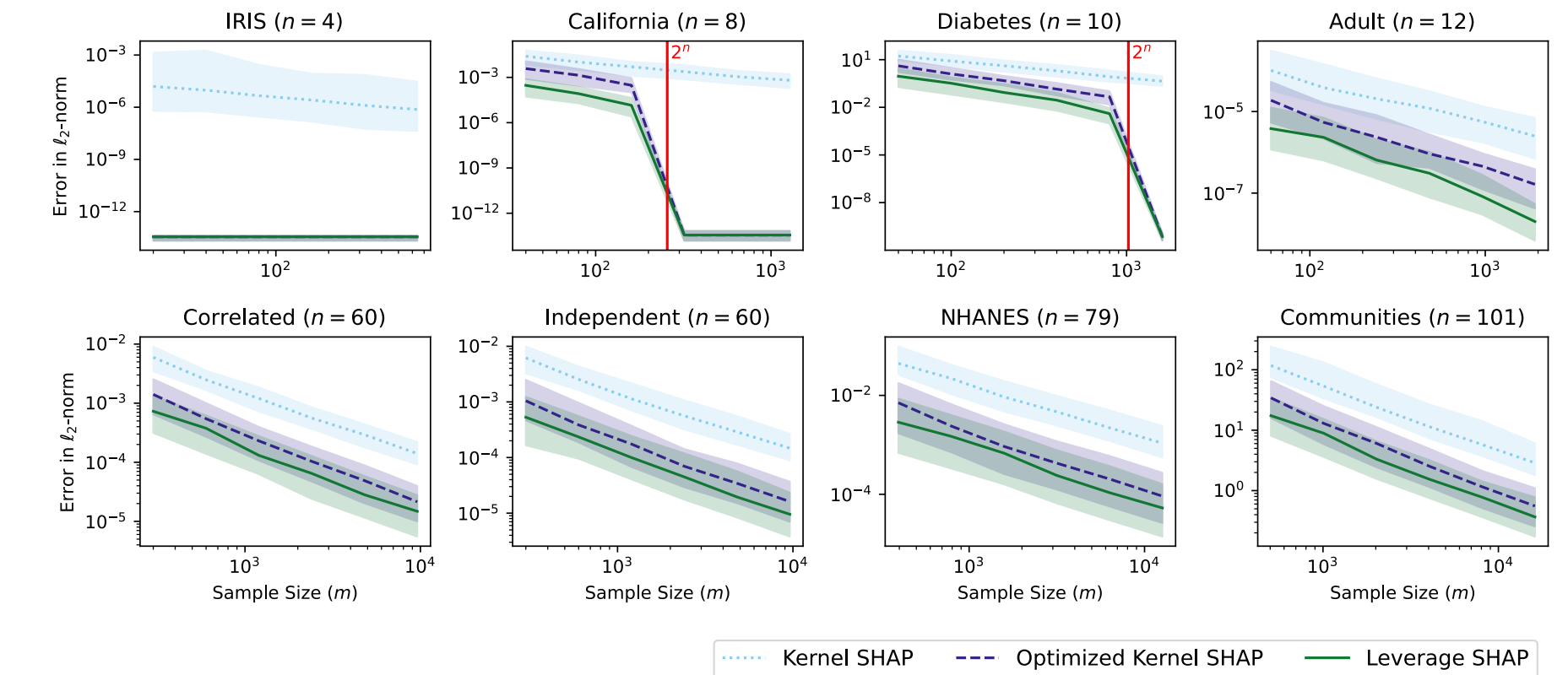
$$\|A\tilde{\phi} - b\|_2^2 \leq (1 + \epsilon) \|A\phi - b\|_2^2$$

and, for $\gamma = \frac{\|A\phi - b\|_2^2}{\|A\phi\|_2^2}$,

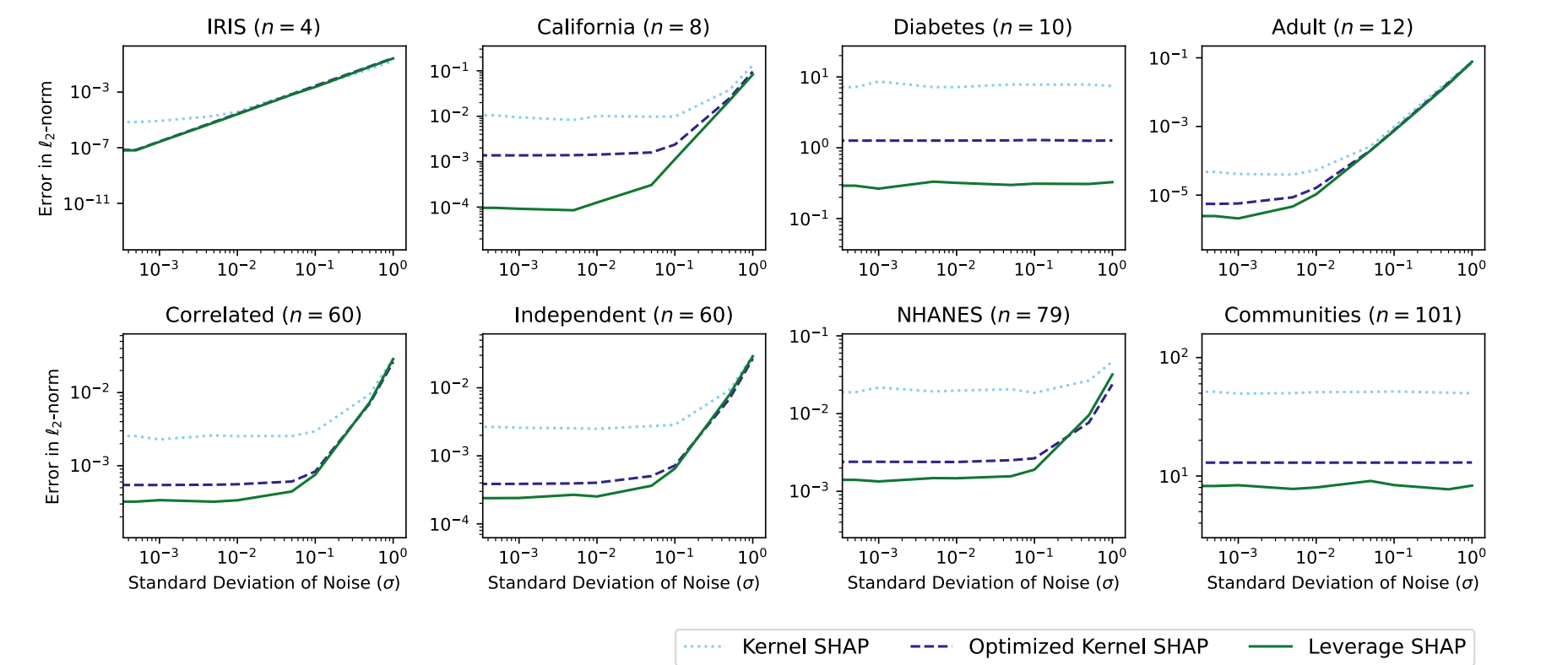
$$\|\tilde{\phi} - \phi\|_2^2 \leq \epsilon \gamma \|\phi\|_2^2.$$

Performance by...

Sample Size



Noise



Noise is relevant because set functions can often only be approximated e.g., $v(S) = \mathbb{E}_{x^S}[f(x^S)]$

Notes

¹Our work is available on arXiv at <https://arxiv.org/abs/2410.01917>

²This work was supported by the National Science Foundation under Grant No. 2045590 and Graduate Research Fellowship Grant No. DGE-2234660