

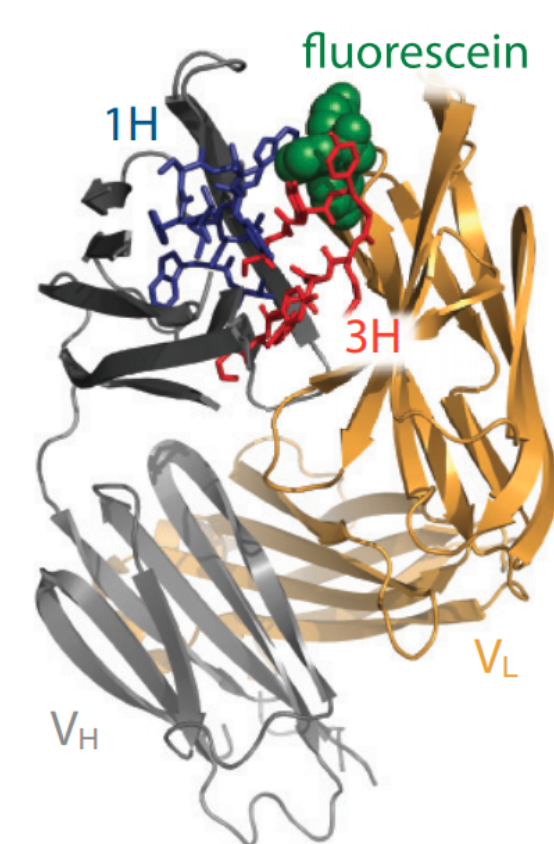
Learning to Make Decisions via Submodular Regularization

Ayya Alieva³, Aiden Aceves¹, Jialin Song¹, Stephen Mayo¹, Yisong Yue¹, Yuxin Chen²

¹Caltech, ²University of Chicago, ³Stanford University

Motivation

- Active learning, experimental design and protein engineering are instances of **combinatorial optimization**.
- These problems are usually NP-hard.
- Greedy could be intractable and/or highly suboptimal.



Our Approach

- ML framework that learns a greedy heuristic based on the **submodular-norm loss**.
- The learned objective is permutation-invariant and approximately submodular with **strong approximation guarantees**.
- The model is easily integrated with ML pipelines.

Submodular Surrogates

Given: ground set V ; value function $u: 2^V \rightarrow \mathbb{R}$; budget k .

Goal: Policy $\pi: 2^V \rightarrow V$ that maximizes $u(S_{\pi,k})$ by iteratively adding to the set:

$$S_{\pi,i} = S_{\pi,i-1} \cup \{\pi(S_{\pi,i-1})\}.$$

Approach: Design a surrogate objective function $f: 2^V \rightarrow \mathbb{R}$ which is

- Monotone:** $A \subseteq B \subseteq V$ and $e \in V \setminus A$, $f(A \cup \{e\}) \geq f(A)$.
- Submodular:** $A, B \subseteq V$ and $e \in V \setminus (A \cup B)$, $f(A \cup \{e\}) + f(B \cup \{e\}) - f(A \cup B \cup \{e\}) \leq f(A \cup \{e\}) + f(B) - f(A \cup B)$.
- Approximately equal to oracle u .**

If we had a "good" **submodular surrogate** f :

- Let $g(A, e) = f(A \cup \{e\}) - f(A)$ for $A \subseteq V$ and $e \in V$.
- $g(A, e) \approx g^{\text{exp}}(A, e) = u(A \cup \{e\}) - u(A)$
- $(S) = \arg \max_{e \in V} g(S, e)$ is **near-optimal**.

Problem: Hand-engineering f is hard; evaluating g^{exp} is expensive.

Learning with Submodular Regularization (LeaSuRe)

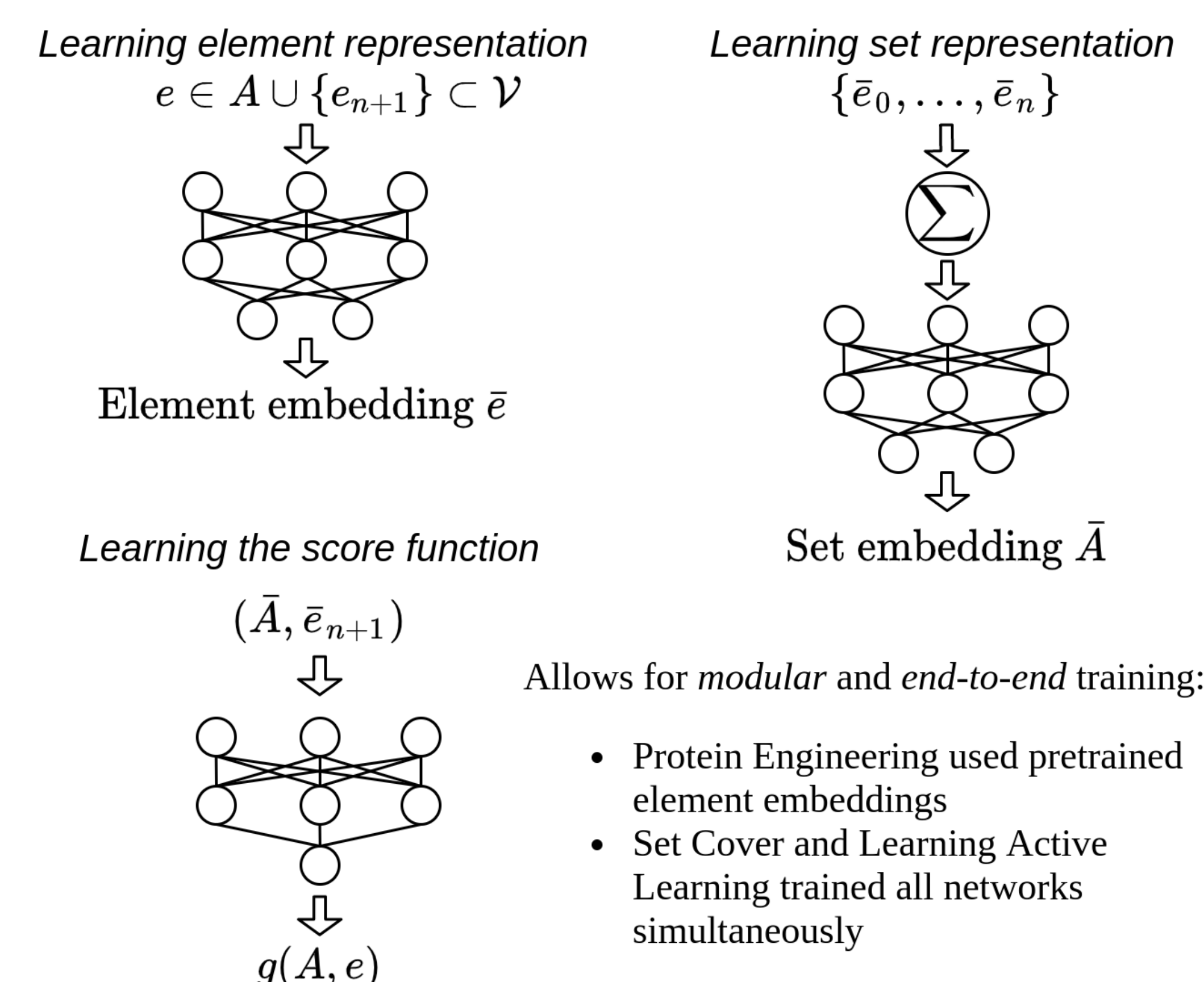
We train NN $g: 2^V \times V \rightarrow \mathbb{R}$ to approximate g^{exp} :

- DAGger [5] collects **real labelled** dataset $D_{\text{real}} = \{(A, x, g^{\text{exp}}(A, x))\}$.
- For $A, x \in D_{\text{real}}$ we generate random **superset** A and form an **unlabelled synthetic** dataset $D_{\text{synth}} = \{(A, x, A, x) | A \supseteq A, A, x \in D_{\text{real}}\}$
- We update g using **submodular-norm** loss:

$$\begin{aligned} \text{Loss}(g, g^{\text{exp}}) = & \sum_{(A, x) \in D_{\text{real}}} (g^{\text{exp}}(A, x) - g(A, x))^2 \\ & + \sum_{(A, x, A, x) \in D_{\text{synth}}} (|g(A, x) - g(A, x)|) \\ & + \sum_{(A, x) \in D_{\text{synth}}} \text{ReLU}(-g(A, x)), \end{aligned}$$

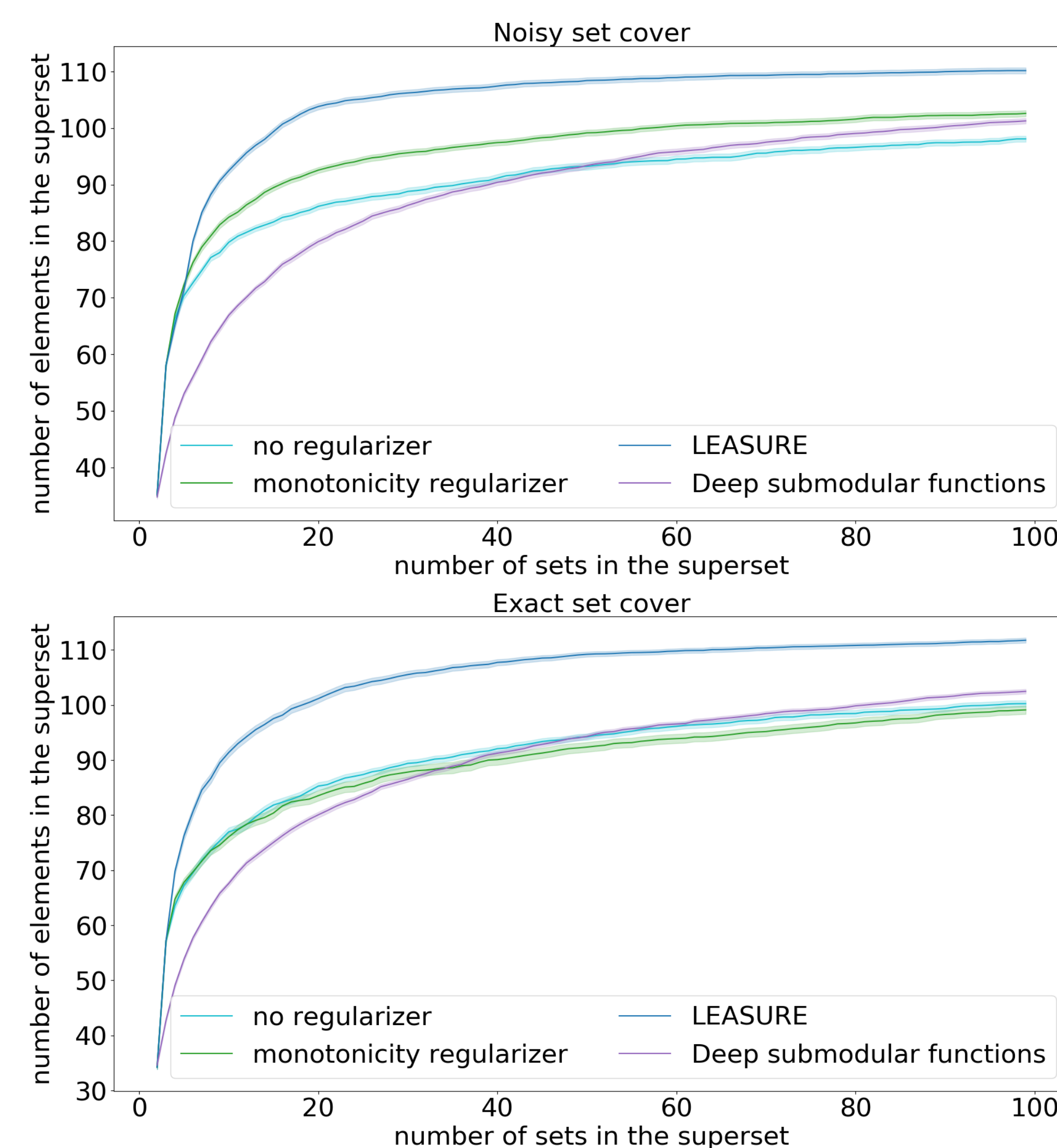
LeaSuRe encourages g to model the gain of a **monotone, submodular** function approximating the oracle u .

Architecture



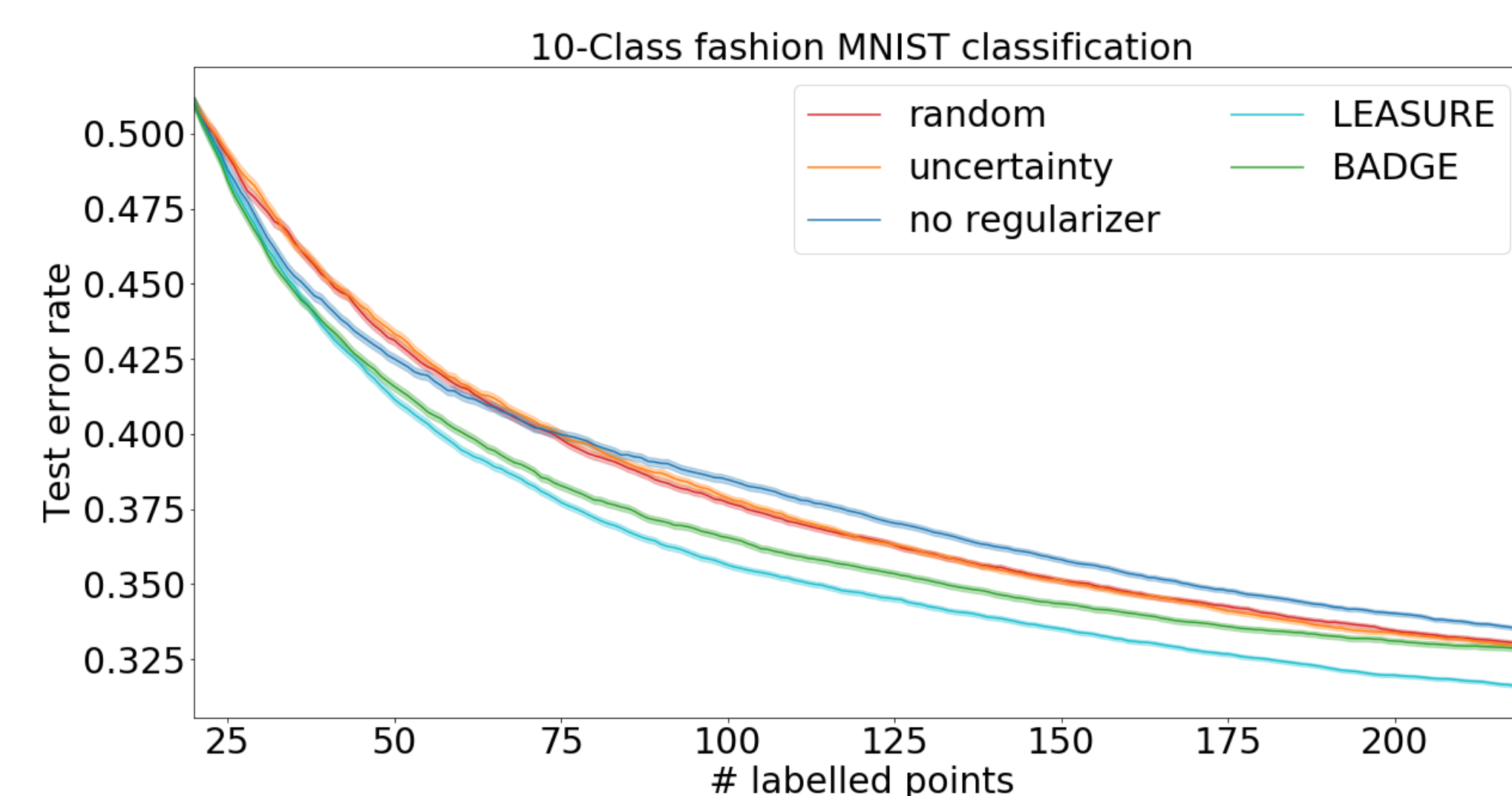
Protein Engineering used UniRep [2, 1]; Set Cover and Learning Active Learning used a two layer DNN.

Set Cover



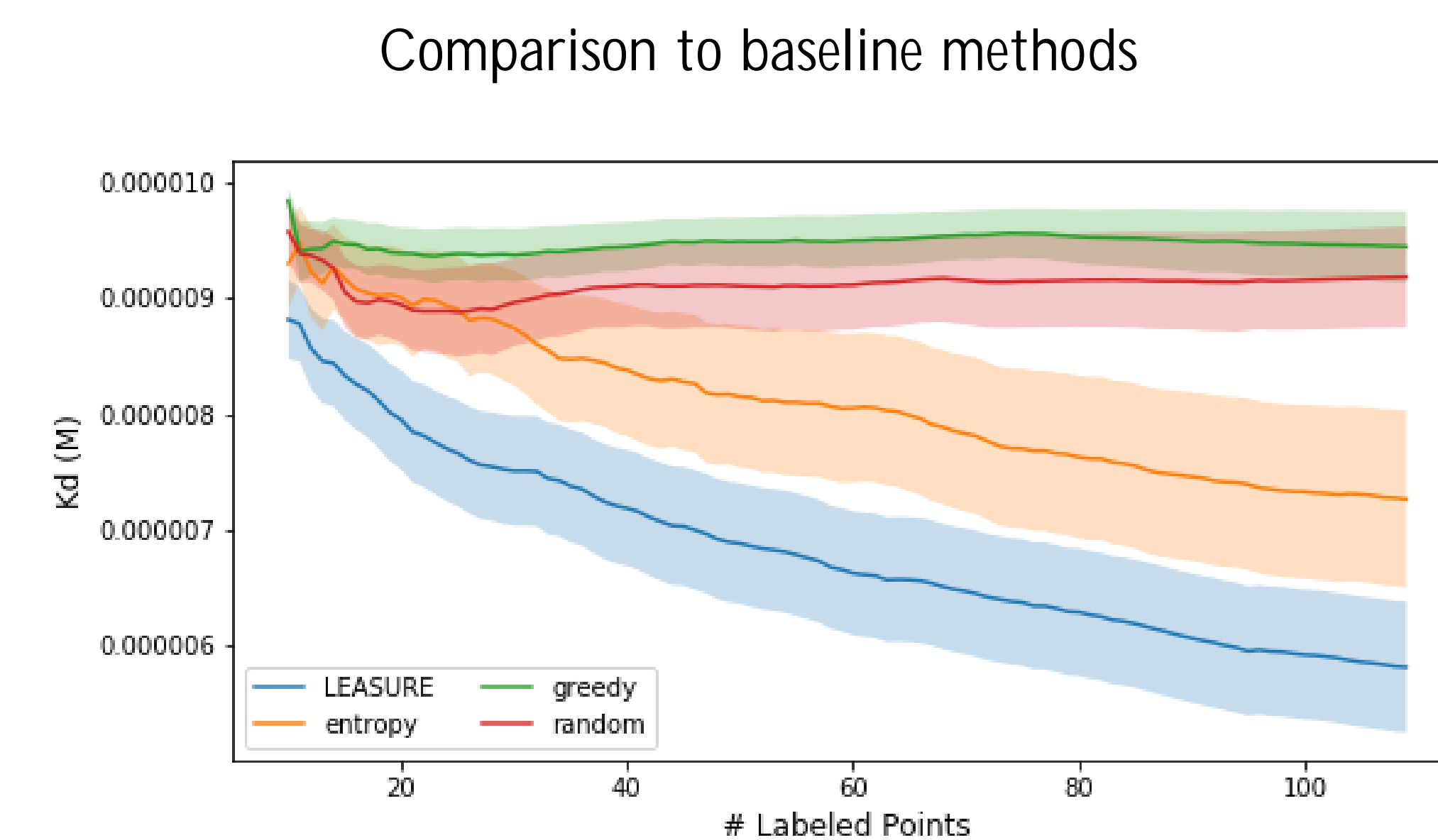
Baselines: "no regularizer" - MSE loss; "monotonicity regularizer" - MSE + ReLU loss; Deep Submodular Functions [4].

Learning active learning on Fashion MNIST

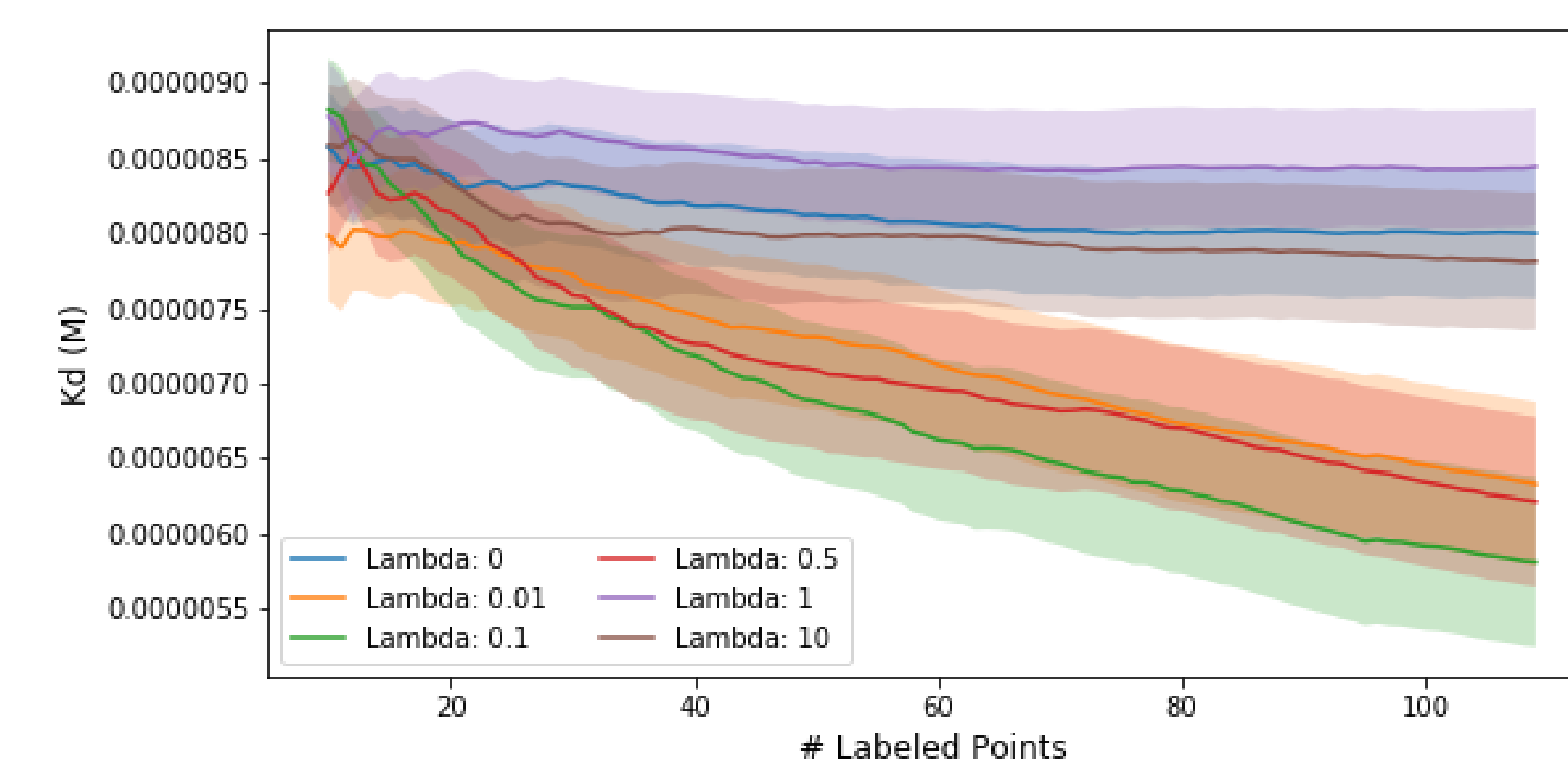


Baselines: "random" - random sampling; "uncertainty" - uncertainty sampling; "no regularizer" - DAGger + MSE loss; BADGE [3].

Protein Engineering



Effect of scaling parameter λ for LeaSuRe



Given: A large set of unlabelled proteins.
Goal: Determine a subset of high-potency (low K_D) proteins.

Problem: The protein set is too large to test the potency of all elements.

Baselines: "random" - random sampling; "greedy" - hand-engineered surrogate g ; "entropy" - a version of uncertainty sampling.

References

- [1] Aceves, A. (2021). Our code. <https://github.com/ajaceves/alilpe>.
- [2] Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M., and Church, G. M. (2019).
- [3] Ash, J. T., Zhang, C., Krishnamurthy, A., Langford, J., and Agarwal, A. (2020).
- [4] Dolhansky, B. W. and Bilmes, J. A. (2016).
- [5] Ross, S., Gordon, G., and Bagnell, D. (2011).