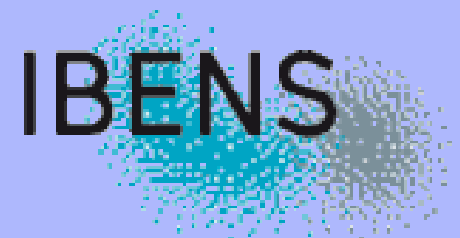


RED: Recommendations Encouraging Diversity

Project check-in #2

10.06.2025

Clemence Reda



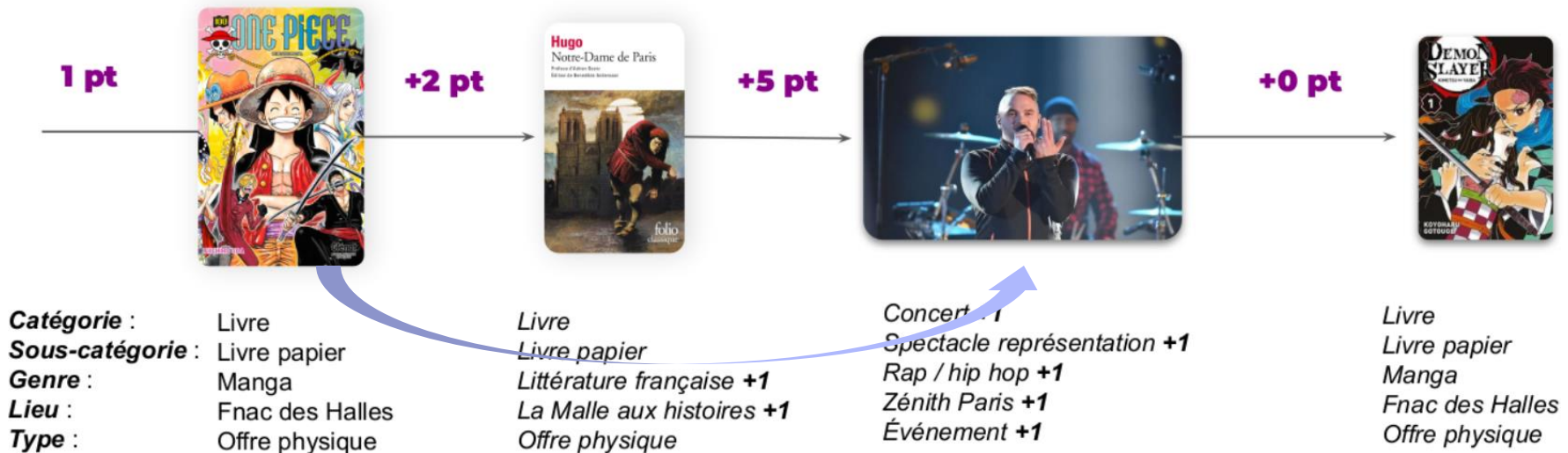
Background

1. Motivation: Diversity in recommendation (p.3-4)
2. Objectives of the RED project (p.5)
3. Data & diversity metrics (p.6-8)
4. Intro to MABs and DPPs (p.9-11)
5. Main idea (p.12)

Pass Culture a phone app for French teens (<20yr) to browse and book cultural goods nearby with credits.

Diversification points obtained for each new category / subcategory / genre / location / type (a bit like set cover; achievement score); those are not visible to the user

Comment mesurer la diversification ?



Courtesy of Jill-Jenn Vie (Inria SODA).

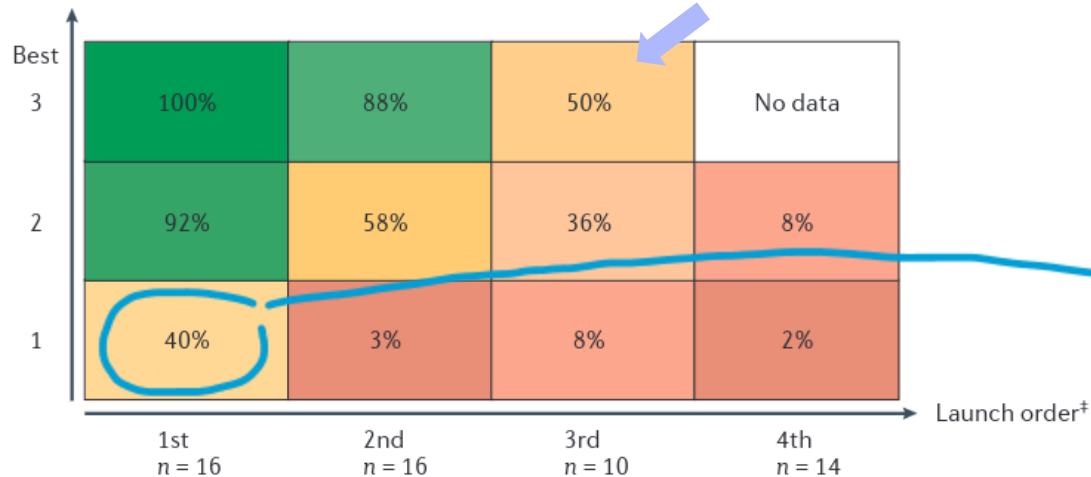
Drug repurposing First-in-class versus Best-in-class

First launched in that mechanistic class

Highest therapeutic advantage [...]

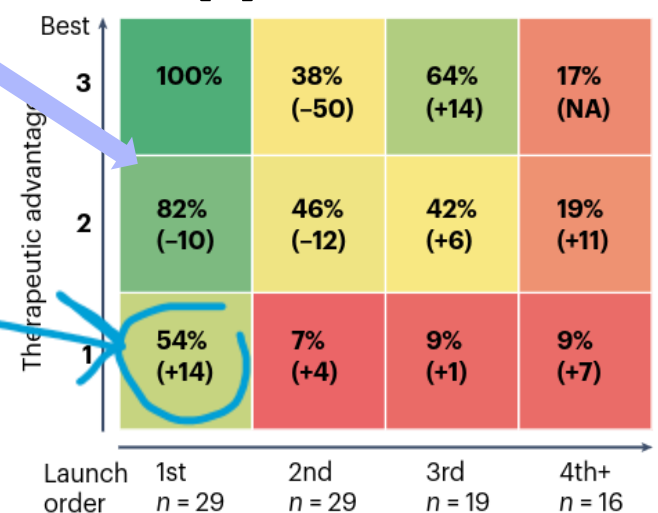
In 2013 [1]

Therapeutic advantage*



"The data indicated that it is slightly better to be first than to be best"

In 2023 [2]



"[...] products that are first-to-launch increasingly tend to perform better [...]"

[1] Schulze, U., & Ringel, M. (2013). What matters most in commercial success: first-in-class or best-in-class?. *Nature Reviews Drug Discovery*, 12(6), 419-420.

[2] Spring, L., Demuren, K., Ringel, M., & Wu, J. (2023). First-in-class versus best-in-class: an update for new market dynamics. *Nat Rev Drug Discov*, 22(7), 531-532.

Objective of RED to design recommender systems for personalized good and diverse items.

taking into account the **user's interests/item history**

suggesting items with a **high probability of positive feedback** from the user

somewhat "**controllable**" recommendations: *e.g.*, out of the user's comfort zone

Scientific challenges

- Incorporate a tradeoff between quality and diversity
- Versatile enough to be applicable in all use cases
- Not too computationally expensive (millions of items)

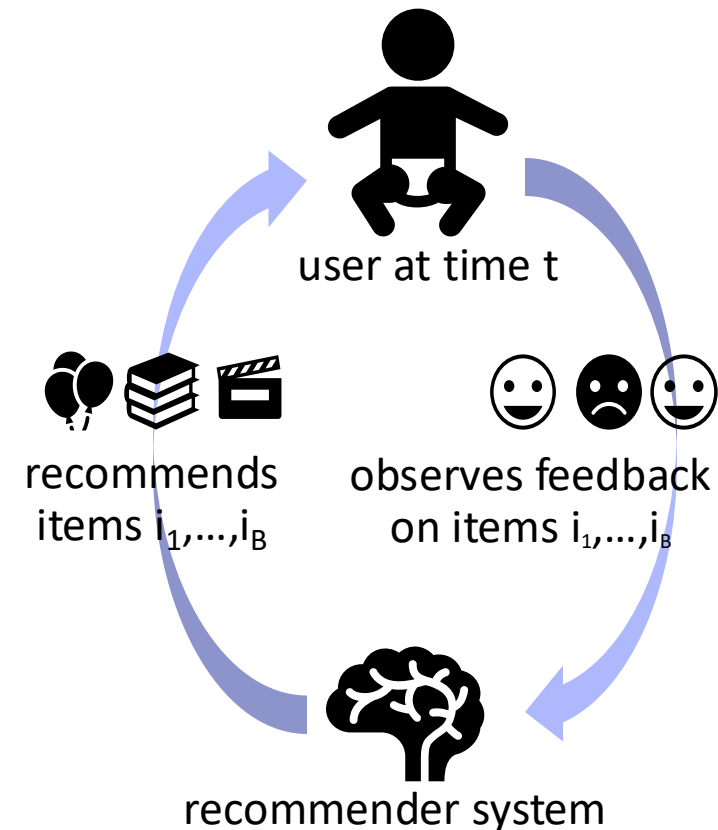
Data

Φ_i : embedding for item i with $|\Phi_i|=1$

H_u : item history for user u (at all times)

q : (un)known feedback model with positive values

B : the size of the batch of items to recommend at each time



Data sets I don't have access to Pass Culture data

Synthetic : Gaussian Generate item and user embeddings Φ_i and c_u at random, $q(i,u) \propto \Phi_i^T c_u$

Pseudo-real :  MovieLens [1] for movie recommendation : 0 (not seen), 1, 2, ..., 5 stars

Apply universal-sentence-encoder to movie title & keywords

$q(i,u) := \text{SVD}(\text{rating matrix})[i,u]$ # fill out 0's with Singular Value Decomposition

RMSE: 0.96 (on held-out data)

Pseudo-real : PREDICT [2] for drug repurposing : 0 (not tested), 1 (success), -1 (failure)

Select 10 first PCs from drug and disease feature vectors to get Φ_i and c_u

% of cumulative explained variance: 85% (drugs), 68% (diseases)

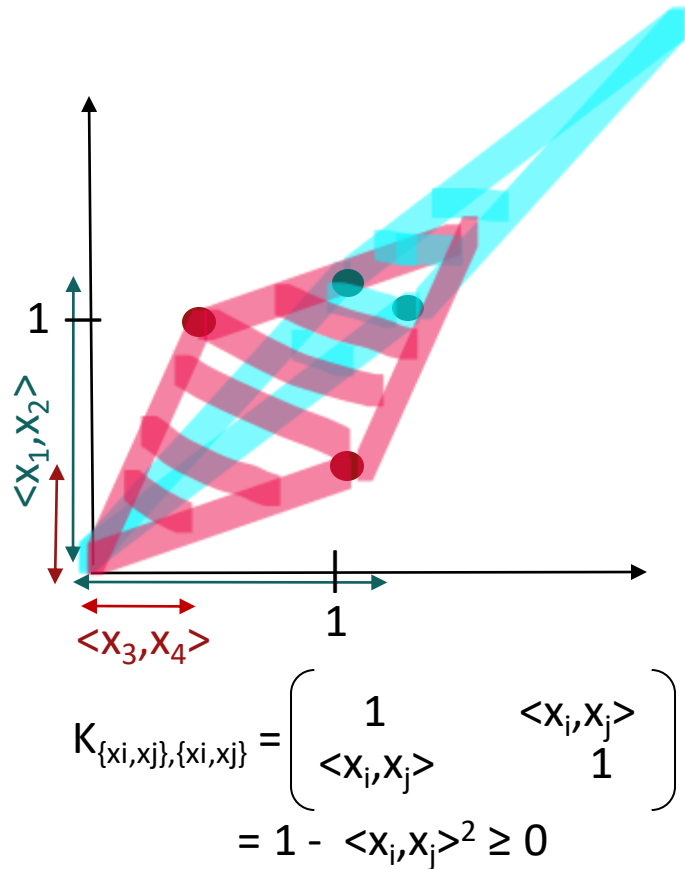
$q(i,u) := \text{RF}([\Phi_i, c_u])$ # classification with Random Forest

Average accuracy: 0.97 (on held-out data)

[1] <https://movielens.org/>

[2] Réda, C. (2023). PREDICT drug repurposing dataset (2.0.1) [Data set]. Zenodo. [doi:10.5281/zenodo.7983090](https://doi.org/10.5281/zenodo.7983090)

Diversity metrics introduction to similarity kernels



(Item embedding) kernel $K_{\{x\}, \{y\}} = \langle \Phi(x), \Phi(y) \rangle$
 e.g., linear kernel: $\Phi = \text{Id}$, \sim similarity b/w items with
 K positive definite

← (Log)-determinant of (a subset of) the kernel Volume
 in space occupied by a set S of items:

$$\text{vol}(S) = |\det(K_{S,S})|^{1/2}$$

Drawn for all $i \mid |x_i|^2 = 1$, K linear kernel in 2D
 ... but works for any # of dimension and # of
 points (and any symmetric kernel)

$|K_{\{x_1, x_2\}, \{x_1, x_2\}}| \leq |K_{\{x_3, x_4\}, \{x_3, x_4\}}|$
 where $\langle x_3, x_4 \rangle \leq \langle x_1, x_2 \rangle$ (cosine similarity)

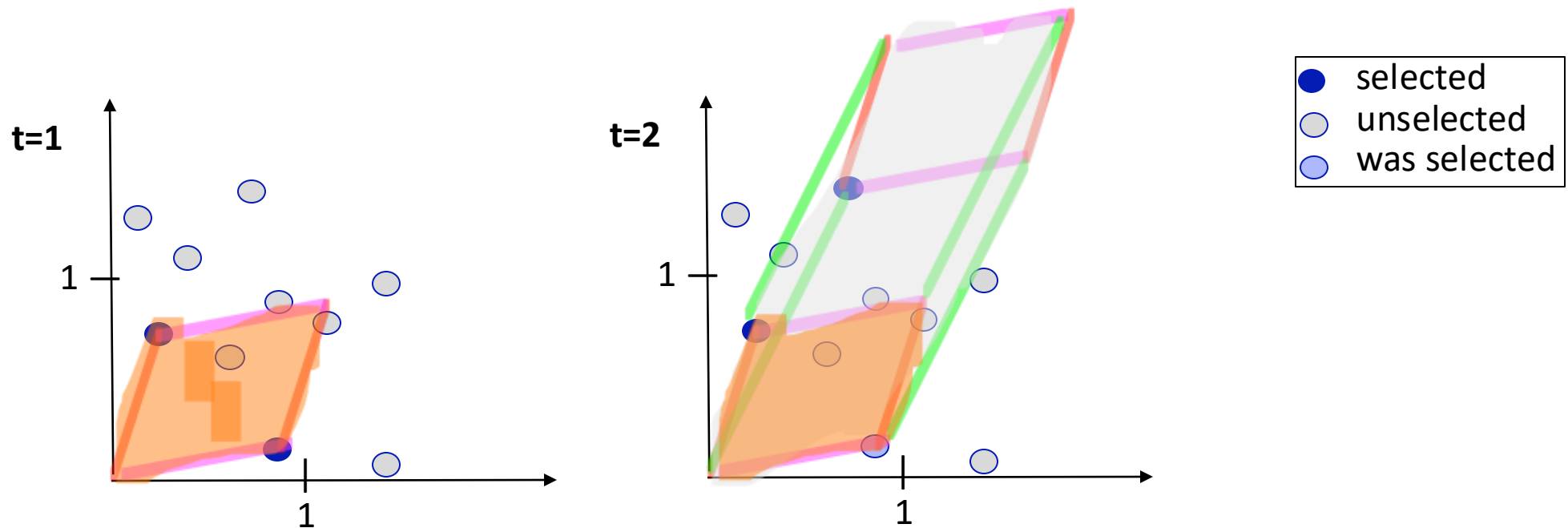
(Pointwise) relevance and diversity metrics with kernels

Relevance: "Click-through-rate" $R(S=\{i_1, \dots, i_B\}, H_t) = \sum_{i \in B} q(i, u_t) / B$

Intrabatch diversity: "inside a batch" $AD(S) = \text{vol}(K_{S,S})$

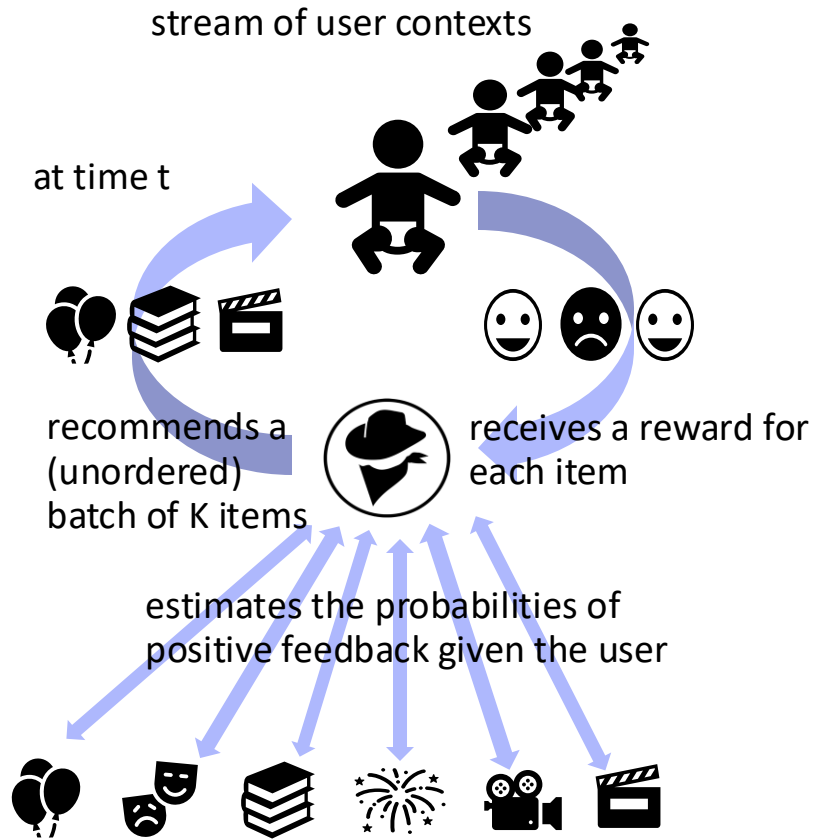
Interbatch diversity: "across history" $ED(S, H_t) = \text{vol}(K_{S \cup H_t, S \cup H_t}) - \text{vol}(K_{H_t, H_t})$

Avg # unique recommendations in history Complementary to the interbatch metric



(Stochastic, contextual) multi-armed bandits or MABs

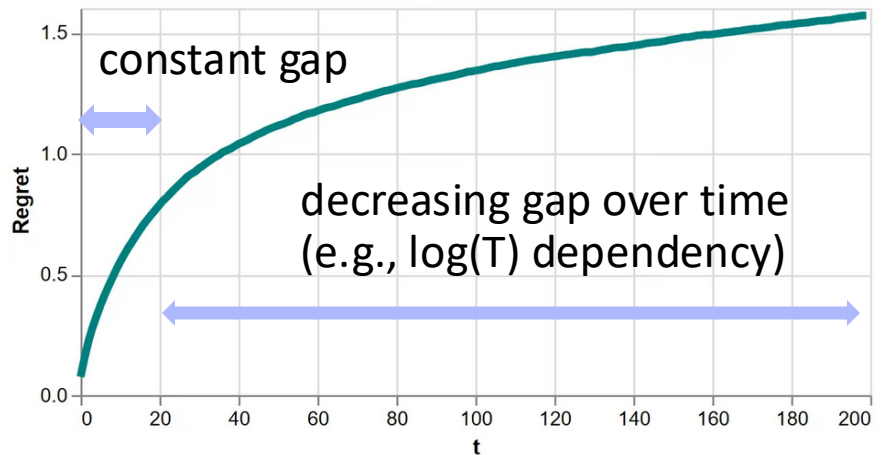
online reward maximization when q is *unknown*



Definition A MAB is defined by its **sampling rule**: arm with highest probability of positive feedback?

Reward maximization Performance is compared to a deterministic oracle with access to the true q

$$\text{Regret}(T) = \sum_t \text{Perf}(q, t) - \text{Perf}(q_t, t)$$



(Finite) determinantal Point Processes or DPPs sampling diverse set of points in a data-driven fashion

Point Process A distribution over finite subsets of a (finite) set $X=\{x_1, x_2, \dots, x_N\}$

Definition* The probability of a subset is correlated to its determinant

$\text{Prob}(S \subseteq X) \propto \det(K_{S,S})$ where $K_{S,S}$ is the kernel built on the subset $S \subseteq X$

- **Sampling** Algorithms in $O(N^3)$ [1] or $O(N \cdot \text{poly}(k))$ [2] to sample k points out of N from a DPP
- **Conditioning (SAMPLE)** [3] $\text{Prob}(S \subseteq X \mid H \subseteq X) \propto \det(K_{S,S} - K_{S,H} K_{H,H}^{-1} (K_{S,H})^T)$
- **MAP inference (MAX)** Finding a subset S of size k maximizing $\text{Prob}(S \subseteq X)$ is NP-hard

[1] Theorem 7 and Algorithm 18 in Hough, J. B., Krishnapur, M., Peres, Y., & Virág, B. (2006). Determinantal processes and independence.

[2] Calandriello, D., Derezhinski, M., & Valko, M. (2020). Sampling from a k -DPP without looking at all items. *Advances in Neural Information Processing Systems* 33, 6889-6899.

[3] Borodin, A., & Rains, E. M. (2005). Eynard–Mehta theorem, Schur process, and their Pfaffian analogs. *Journal of statistical physics*, 121, 291-317.

(Finite) determinantal Point Processes or DPPs sampling diverse set of points in a data-driven fashion

Point Process A distribution over finite subsets of a (finite) set $X=\{x_1, x_2, \dots, x_N\}$

Definition* The probability of a subset is correlated to its determinant

$\text{Prob}(S \subseteq X) \propto \det(K_{S,S})$ where $K_{S,S}$ is the kernel built on the subset $S \subseteq X$

- **Sampling** Algorithms in $O(N^3)$ [1] or $O(N \cdot \text{poly}(k))$ [2] to sample k points out of N from a DPP ✗
- **Conditioning (SAMPLE)** [3] $\text{Prob}(S \subseteq X \mid H \subseteq X) \propto \det(K_{S,S} - K_{S,H} K_{H,H}^{-1} (K_{S,H})^T)$
in $O(k^3 + |H|^3 + k^2 |H|^2 + (N - |H|)k^2)$ [4] includes history
- **MAP inference (MAX)** Finding a subset S of size k maximizing $\text{Prob}(S \subseteq X)$ is NP-hard greedy approx. [5] in $O(k^2 N)$

[1] Theorem 7 and Algorithm 18 in Hough, J. B., Krishnapur, M., Peres, Y., & Virág, B. (2006). Determinantal processes and independence.

[2] Calandriello, D., Derezhinski, M., & Valko, M. (2020). Sampling from a k -DPP without looking at all items. *Advances in Neural Information Processing Systems* 33, 6889-6899.

[3] Borodin, A., & Rains, E. M. (2005). Eynard–Mehta theorem, Schur process, and their Pfaffian analogs. *Journal of statistical physics*, 121, 291-317.

[4] Mariet, Z., Gartrell, M., & Sra, S. (2019, April). Learning determinantal point processes by corrective negative sampling. In *The 22nd International Conference on Artificial Intelligence and Statistics* (pp. 2251-2260). PMLR.

[5] Chen, L., Zhang, G., & Zhou, E. (2018). Fast greedy map inference for determinantal point process to improve recommendation diversity. *Advances in Neural Information Processing Systems*, 31.

(Finite) determinantal Point Processes or DPPs leveraging the quality-diversity decomposition

Quality-diversity decomposition (QD) [1] $K = Q\Phi^T\Phi Q$

where Φ is the item embedding matrix $|\Phi_i|=1$ and $Q = \text{diag}(q_1, q_2, \dots, q_N)$ where $q_i \geq 0$

$$\text{Prob}(S \subseteq X) \propto \det(K_{S,S}) = \prod_{i \in S} q_i^2 \det(\Phi_{S,:}^T \Phi_{S,:})$$

$$f_{\text{QD}}(S, u) = 2 \log \text{vol}(Q_{S,u}^{2\lambda} \Phi_{S,:}^T \Phi_{S,:}^{2(1-\lambda)} Q_{S,u}^{2\lambda})$$

Scientific challenges (*bis*)

- Incorporate a tradeoff between quality and diversity
→ **use (a version of) QD decomposition as objective function**
- Versatile enough to be applicable in all use cases
→ **for any kernel and any feedback model, including *unknown* ones**
- Not computationally too expensive (millions of items)
→ **a time complexity with a dependency larger than $O(N)$ will be too expensive**

[1] Kulesza, A., & Taskar, B. (2010). Structured determinantal point processes. *Advances in neural information processing systems*, 23.

My progress on the project since last time

1. Generalized QD objective function (p.14-15)
2. Case where the reward model is known (p.16-19)
3. Case where the reward model is unknown (p.20-21)
4. Adaptive QD tradeoff (p.22-24)
5. Perspectives (p.25)

Generalized quality-diversity objective function

Data

Φ_i : embedding for item i with $|\Phi_i|=1$

H_u : item history for user u

q : feedback model with positive values

B : the size of the batch of items to recommend at each time

K : kernel similarity

Hyperparameters

λ : weight of the relevance task

η : regularization factor

Objective function f for user u

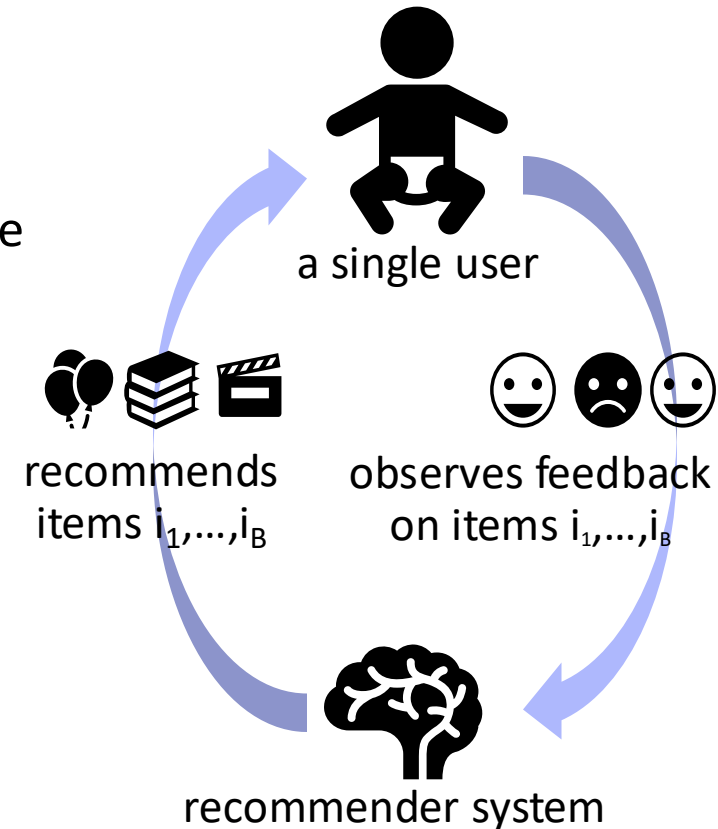
For all S , $f_{\text{obj}}(S, u) = 2 \log \text{vol}(Q_{S,u}^{2\lambda} \mathbf{M}_{S|Hu,S|Hu}^{2(1-\lambda)} Q_{S,u}^{2\lambda})$

- $Q_{S,u} = \text{diag}(\{ q(i, u) \text{ for } i \text{ in } S \})$

- $\mathbf{M}_{S|Hu,S|Hu} = K_{S,S} - K_{S,Hu}(K_{Hu,Hu} + \eta \text{Id})^{-1}(K_{S,Hu})^T$

Conditioning:
dependency in $\Omega(|H|^2 \log |H|)$

a trajectory of length T



Generalized quality-diversity objective function with a fuzzy denuding/masking approach

Practical function f for user u

$$f_{\text{Prac}}(S, u) = 2 \log \text{vol}(Q_{S,u}^{2\lambda} M_{S \cap H_u, S \cap H_u}^{2(1-\lambda)} Q_{S,u}^{2\lambda}) \quad \text{where } M_{S \cap H_u, S \cap H_u} = K_{S,S} - K'_{S,S}$$

and for all x, y in S , $K'_{\{x\}, \{y\}} = \langle \Phi'(x), \Phi'(y) \rangle$ and $\Phi'(x) = 0$ if $\min_{\psi \in H_u} |\psi - x| > \alpha$, else $\Phi(x)$



dependency in $\Omega(k \log |H|)$
using a k-d tree for all of S
can be faster using GPUs

Synthetic Gaussian data, $B=5$, $T=10$, $\alpha=0$, $\eta=0.01$, K linear, average across 10 iters (MAX)

	QD $\lambda=0.5$	Obj $\lambda=0.5$	Prac $\lambda=0.5$	QD $\lambda=0.8$	Obj $\lambda=0.8$	Prac $\lambda=0.8$	QD $\lambda=0.1$	Obj $\lambda=0.1$	Prac $\lambda=0.1$	
Relevance	0.586	0.516	<u>0.559</u>	0.571	0.524	<u>0.558</u>	0.586	0.496	<u>0.560</u>	rewards in [0,1]
Intrabatch diversity	0.956	<u>0.971</u>	0.976	0.990	0.965	<u>0.971</u>	0.956	0.961	<u>0.960</u>	
Interbatch diversity	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	Volume collapses
%unique recomm	<u>0.100</u>	1.000	1.000	<u>0.100</u>	1.000	1.000	<u>0.100</u>	1.000	1.000	



Case where the reward model is known

Synthetic Gaussian data: reward $q(i,u) \propto \Phi_i^T c_u$

$B=5, T=10, \alpha=0, \lambda=0.5, \eta=0.01, K$ linear, average across 10 iters (MAX)

	ϵ -greedy ($\epsilon=\lambda$)	QD	Markov DPP [1]	Obj	Prac
Relevance	0.536	0.592	0.524	<u>0.526</u>	0.572
Intrabatch diversity	0.582	0.956	<u>0.953</u>	0.972	0.972
Interbatch diversity	0.000	0.000	0.000	0.000	0.000
%unique recomm.	0.200	0.100	<u>0.600</u>	1.000	1.000



Baseline for
diversity



Baseline for
relevance



Baseline for integrating
user history

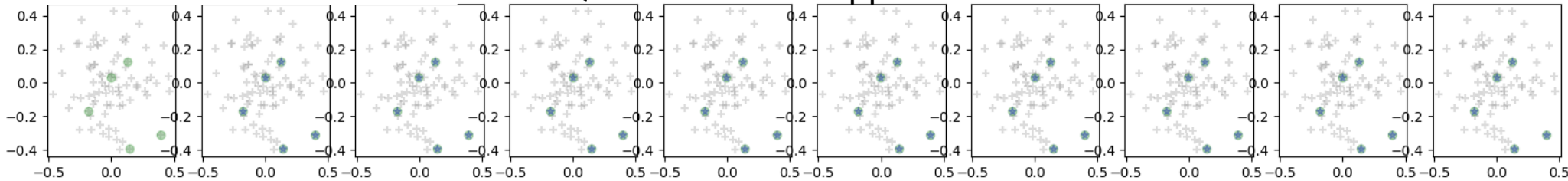
[1] Affandi, R. H., Kulesza, A., & Fox, E. B. (2012). Markov determinantal point processes. *arXiv preprint arXiv:1210.4850*.

Case where the reward model is known (SAMPLE)

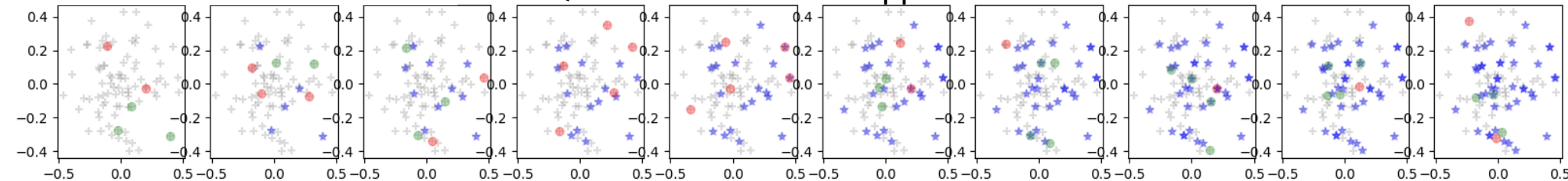
Synthetic Gaussian data: reward $q(i,u) \propto \Phi_i^T c_u$

2D PCA plot of selected points (positive ●, negative ●) at each point of the trajectory of length $T=10$
Previously selected ●

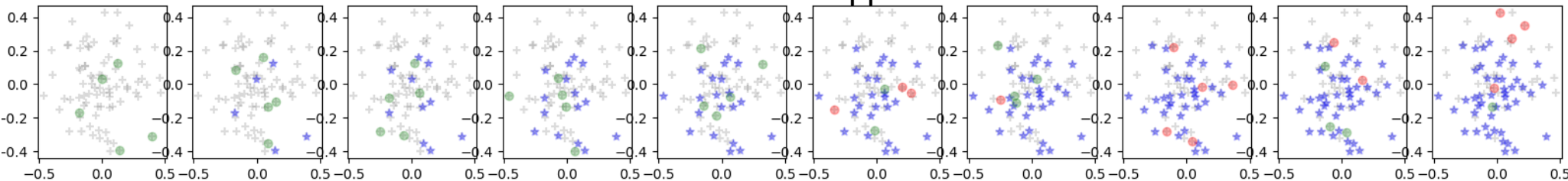
QD with the MAX approach



QD with the SAMPLE approach



Prac with the MAX approach



Case where the reward model is known

MovieLens data: reward $q(i,u) := \text{SVD}(\text{rating matrix})[i,u]$

$B=5, T=10, \alpha=0, \lambda=0.5, \eta=0.01, K$ linear, average across 10 iters (MAX)

	ϵ -greedy ($\epsilon=\lambda$)	QD	Markov DPP [1]	Obj	Prac
Relevance	4.581	4.577	<u>4.430</u>	4.421	4.481
Intrabatch diversity	0.261	0.761	0.543	<u>0.692</u>	0.711
Interbatch diversity	0.000	0.000	0.000	0.000	0.000
%unique recomb.	0.200	0.100	0.860	1.000	1.000
Runtime (sec.)	353 +- 35	510 +- 6	1,328 +- 12	508 +- 3	<u>515 +- 4</u>

rewards in [1,5]

Baseline for diversity

Baseline for relevance

Baseline for integrating user history

[1] Affandi, R. H., Kulesza, A., & Fox, E. B. (2012). Markov determinantal point processes. *arXiv preprint arXiv:1210.4850*.

Case where the reward model is known

PREDICT data: reward $q(i,u) := \text{RandomForest}([\Phi_i, c_u])$

$B=5, T=10, \alpha=0, \lambda=0.5, \eta=0.01, K$ linear, average across 10 iters (MAX)

	ϵ -greedy ($\epsilon=\lambda$)	QD	Markov DPP [1]	Obj	Prac	rewards in {0,1} (overfitting)
Relevance	1.000	1.000	1.000	1.000	1.000	
Intrabatch diversity	0.637	1.000	0.881	1.000	<u>0.961</u>	
Interbatch diversity	0.000	0.000	0.000	0.000	0.000	
%unique recomb.	0.100	0.100	<u>0.967</u>	1.000	0.933	
Runtime (sec.)	51 +- 1	52 +- 1	56 +- 1	51 +- 0	<u>52 +- 0</u>	

Baseline for diversity

Baseline for relevance

Baseline for integrating user history

[1] Affandi, R. H., Kulesza, A., & Fox, E. B. (2012). Markov determinantal point processes. *arXiv preprint arXiv:1210.4850*.

Case where the reward model is unknown

Assumption: linear reward $q^*(i,u) := ([\Phi_i, c_u])^T \theta^*$

Cumulative regret at time T $\text{Regret}(T) = \sum_{t \leq T} f_{\text{obj}}(S_t, u_t; q^*) - f_{\text{obj}}(S_t, u_t; q_t)$

At each time, we update the model $q_t(i,u) = [\Phi_i, H_u]^T \theta_t$ based on the received feedback and then use MAX on the Upper Confidence Bound of the empirical reward model q_t

Synthetic data : $B=5, T=16, \alpha=0, \lambda=0.5, \eta=0.01, K$ linear (MAX)

(δ -uniform) Upper Confidence Bound

There is β such that for all t, i, u ,

$$q^*(i,u) \leq q_t(i,u) + \beta(t,i)$$

with probability $1-\delta$

The baselines (almost)
never learn because we
forced the MAX approach

What happens if we allow
SAMPLE?

	LinUCB bandit [1]	LinOASM bandit [2]	Prac+UCB
Relevance	0.94	<u>0.95</u>	0.99
Intrabatch diversity	0.94	1.00	1.00
%unique recomm.	0.10	<u>0.10</u>	1.00
Regret(T)	16.15	<u>14.84</u>	5.28

Baseline for learning

Baseline for learning \times diversity

[1] Li et al. "A contextual-bandit approach to personalized news article recommendation." *Proceedings of the 19th international conference on World wide web*. 2010.

[2] Gabillon et al. "Large-scale optimistic adaptive submodularity." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 28. No. 1. 2014.

Case where the reward model is unknown

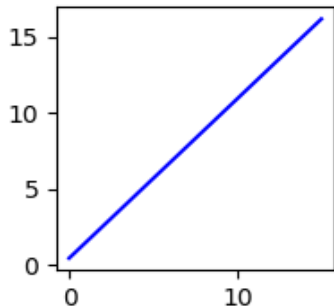
Assumption: linear reward $q^*(i, u) := ([\Phi_i, c_u])^T \theta^*$

Cumulative regret at time T $\text{Regret}(T) = \sum_{t \leq T} f_{\text{obj}}(S_t, u_t; q^*) - f_{\text{obj}}(S_t, u_t; q_t)$

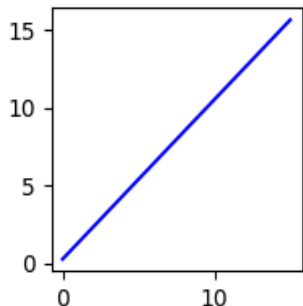
Comparing the regret curves for each bandit algorithm at $T=16$

LinUCB

MAX: cumulative regret 16.15

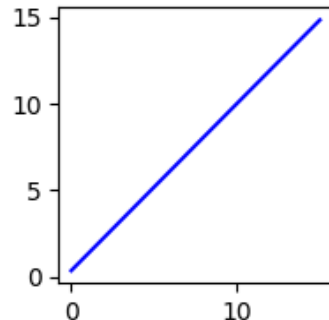


SAMPLE: cumulative regret 15.62

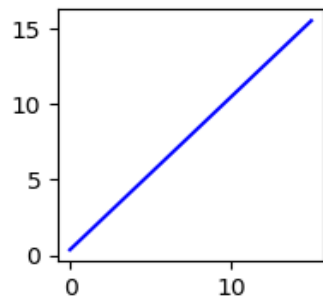


LinOASM

MAX: cumulative regret 14.84

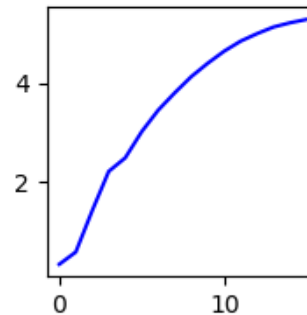


SAMPLE: cumulative regret 15.52



Prac+UCB

MAX: cumulative regret 5.28



The regret may decrease when we allow SAMPLE

But the baselines are still not learning (linear curve)...

... contrary to Prac+UCB (using MAX)

Adaptive quality-diversity tradeoff automatically tune λ in $[0,1]$ to the user feedback

Goal: Finetune on the fly λ to a value which maximizes $\sum_{t < T} f(S_t, u_t ; \lambda)$

We can use any "good" learner (e.g., AdaHedge, EXP3, etc.) to learn the value of λ

Pseudo-code (for any bandit)

- Start with a value λ_0

- Initialize the online learner with the starting value

- For each time t up to T

 - Apply the bandit strategy for sampling S_t

 - Update the learner with the "gain" obtained with λ_t on S_t : $-\nabla_{\lambda} f(S_t, u_t ; \lambda_t)$

 - Obtain the new value λ_{t+1}

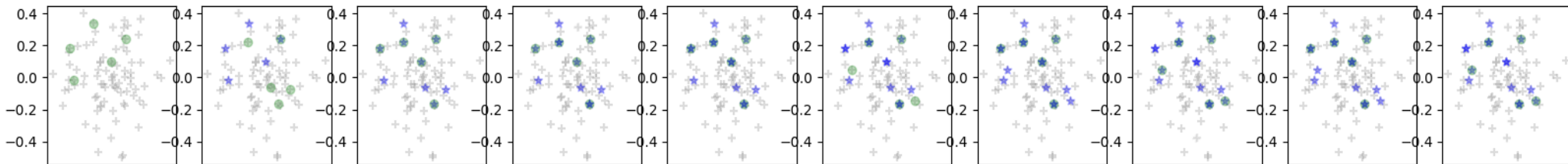
- Return λ_T and the best a posteriori $\lambda^* = \operatorname{argmax}_{\lambda} \sum_{t < T} f(S_t, u_t ; \lambda)$

Adaptive quality-diversity tradeoff

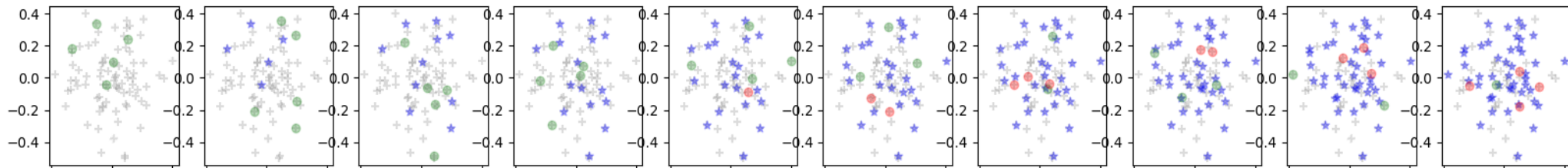
Synthetic data: reward $q(i,u) \propto \Phi_i^\top c_u$

2D PCA plot of selected points (positive ●, negative ●) at each point of the trajectory of length $T=10$
Previously selected ● Starting from $\lambda_0 = 0.3$

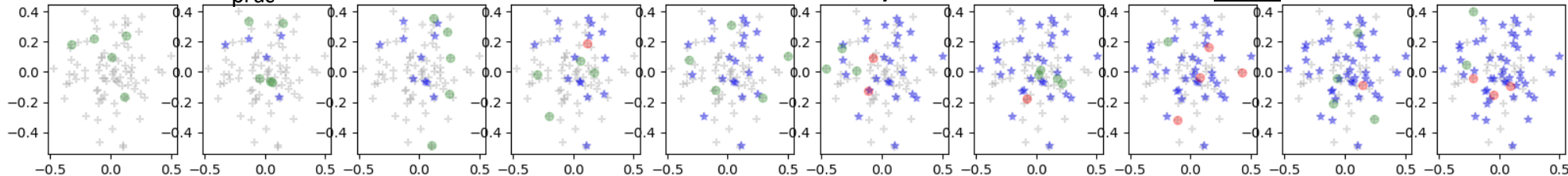
f_{QD} relevance=**0.57** intrabatch diversity=**0.99** $\lambda^*=0.49$ $\lambda=0.59$



f_{obj} relevance=0.52 intrabatch diversity=0.98 $\lambda^*=0.50$ $\lambda=0.51$



f_{prac} relevance=0.53 intrabatch diversity=**0.99** $\lambda^*=0.50$ $\lambda=0.56$

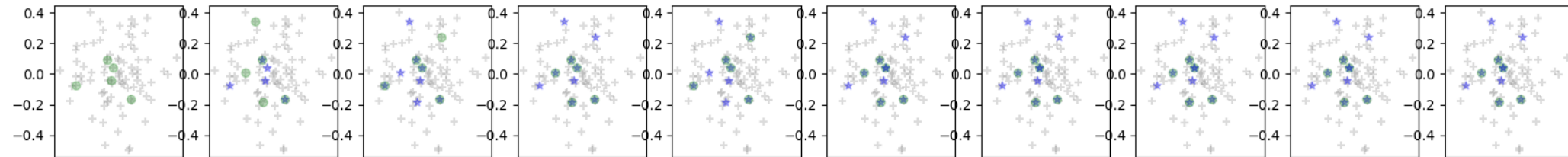


Adaptive quality-diversity tradeoff

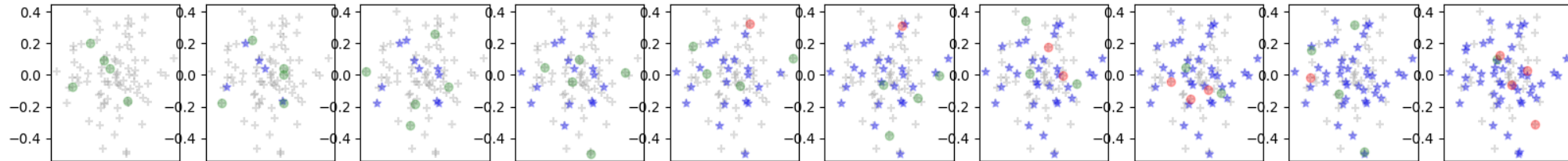
Synthetic data: reward $q(i,u) \propto \Phi_i^T c_u$

2D PCA plot of selected points (positive ●, negative ●) at each point of the trajectory of length $T=10$
Previously selected ● Starting from $\lambda_0 = 0.1$

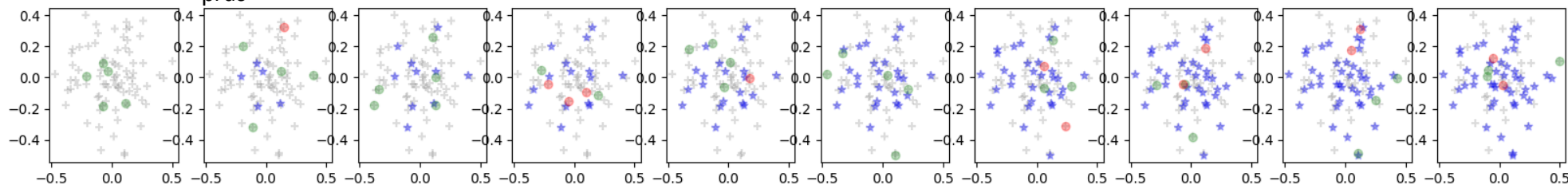
f_{QD} relevance=**0.58** intrabatch diversity=**0.99** $\lambda^*=0.47$ $\lambda=0.61$



f_{obj} relevance=0.54 intrabatch diversity=0.98 $\lambda^*=0.46$ $\lambda=0.51$



f_{prac} relevance=0.53 intrabatch diversity=**0.99** $\lambda^*=0.46$ $\lambda=0.57$



What happens next

1. Find a better measure of interbatch diversity
2. Solve the overfitting problem in the PREDICT data set
3. Optimize the implementation of Prac(+UCB) for millions of items
4. Test in "real conditions": alternating users
5. Derive theoretical guarantees on relevance/diversity