

Bandits for Algorithmic Trading with Signals

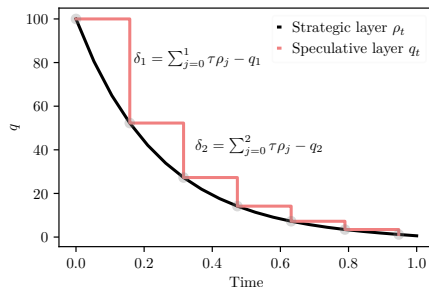
Álvaro Cartea, Fayçal Drissi, Pierre Osselin

Fields-CFI Conference on Recent Advances in Mathematical Finance and Insurance
Oxford-Man Institute

25 September 2023

In practice: slices and child orders

- The trading schedule is sliced into child orders.
- Discrete observation times $t \in \{0, \dots, T\}$. At time t :
 - The trader compares the target inventory $\sum_{j=0}^t \tau \rho_j$ with her current inventory q_t .
 - Child order with quantity: $\delta_t = \sum_{j=0}^t \tau \rho_j - q_t$.



Extensions of optimal execution models

- Realistic assumptions: complex representations of the state space of **market features** and **trading decisions**.

- Market features (contextual information):
 - Predictive price signals.
 - Volatility and liquidity estimators.
 - Fill probability of limit orders.

- Trading decisions (actions):
 - Order type (limit order, market order, ...).
 - Depth of a limit order.
 - Trading venue (dark pool, lit LOB, OTC).

Extensions of optimal execution models in the literature

- Dynamic programming

- + Interpretability and robustness.

- - Need assumptions, dynamics, parameter estimation.
Curse of dimensionality.

- Solutions: NN PDE solvers.

Approximation techniques.

- Machine learning (RL, NN)

- + Scalability to complex environments.

- - Poor interpretability and robustness.
Noisy data.

Training requires data (simulations).

- Solutions: focus on interpretable problems and architectures.

Learning in algorithmic trading

Strategic and Speculative layers

Our approach: two layers for execution algorithms

■ Strategic layer

- An agent liquidates Q shares throughout $[0, T]$.
- Encodes urgency, risk aversion, execution costs, **market impact**.
- Result: trading schedule $(\rho_t)_{t \in \mathbb{R}}$.
- Split the optimal schedule into **child orders**.

■ Speculative layer

- **Targets** ρ .
- Optimises trading performance of child orders using **unsupervised learning**.
- Actions: placement, timing, routing, ..
- Market features: estimators/predictors.

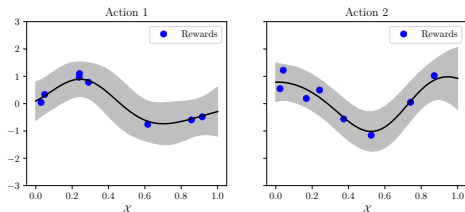
Optimal behaviour of the agent in the speculative layer

- New contextual non-stationary **bandit** (MTGP-LR).
- **Actions**: order type, placement, timing, routing of child orders.
- **Context**: market features (estimators and predictors).
- **Reward** for each arm k : Trading performance.
- **Reward functions**: multi-task Gaussian Process.

Multi-task Gaussian Process for the reward functions

- The agent learns the optimal mapping from contextual market features to optimal actions.
- UCB: Trade-off **exploration** (large posterior variance) / **exploitation** (large posterior mean).
- Two-layer approach for optimal execution
 - Guarantees of the strategic layer.
 - Interpretability of MTGPs.
 - No need for pre-training.
 - Adapts to non-stationary markets.

Multi-task Gaussian Process for the reward functions



- The agent learns the optimal mapping from contextual market features to optimal actions.
- UCB: Trade-off **exploration** (large posterior variance) / **exploitation** (large posterior mean).
- Two-layer approach for optimal execution
 - Guarantees of the strategic layer.
 - Interpretability of MTGPs.
 - No need for pre-training.
 - Adapts to non-stationary markets.

Speculative layer with MTGP-LR

Algorithm 1: MTGP-LR

Input:

Action space \mathcal{A} , Contextual features \mathcal{X} , time horizon $T > 1$.

while $t \leq T$ do

Observe vector of contexts $\mathbf{x}_t \in \mathcal{X}$;

Compute the vector $\mathbf{UCB} \leftarrow \left\{ \tilde{\mu}(a, \mathbf{x}_t) + \tilde{\beta}_t^{1/2} \tilde{\sigma}(a, \mathbf{x}_t) \right\}_{a \in \mathcal{A}}$;

Select action (trading decision) $a_t \leftarrow \arg \max_{a \in \mathcal{A}} \mathbf{UCB}_a$;

Observe reward (trading performance) y_t ;

if **Change-point detected** then

└ Reset the bandit;

⇒ Sub-linear regret in stationary environments !

Non-stationary environments: Likelihood ratio test

- Let \mathcal{W} be the time window that contains the last P rewards. Let $\mathcal{W} = \underline{\mathcal{W}} + \overline{\mathcal{W}}$.

- Hypothesis test:

- H_0 : f_t in $\overline{\mathcal{W}}$ is the same as in $\underline{\mathcal{W}}$.
- H_1 : f_t in $\overline{\mathcal{W}}$ is from a new MTGP.

- Likelihood ratio statistic

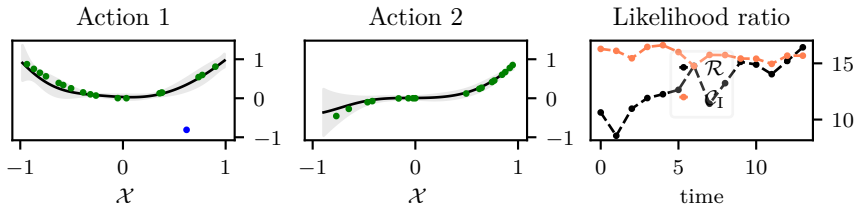
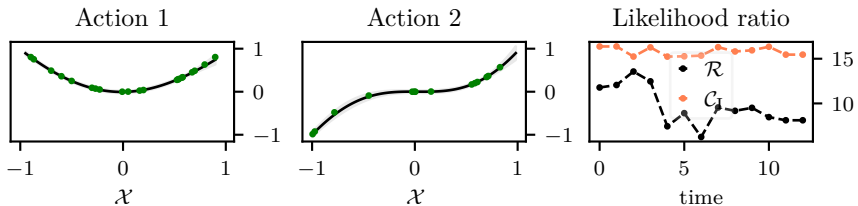
$$\begin{aligned} \mathcal{R} = 2 \log \frac{\rho(\bar{\mathbf{y}}|\mathbf{H}_1)}{\rho(\bar{\mathbf{y}}|\mathbf{H}_0)} &= -\bar{\mathbf{y}}^T (\bar{K} + \sigma_1^2 I)^{-1} \bar{\mathbf{y}} - \log |\bar{K} + \sigma_1^2 I| \\ &\quad + (\bar{\mathbf{y}} - \tilde{\boldsymbol{\mu}})^T (\tilde{K} + \sigma_0^2 I)^{-1} (\bar{\mathbf{y}} - \tilde{\boldsymbol{\mu}}) + \log |\tilde{K} + \sigma_0^2 I|. \end{aligned}$$

- Statistical test: $\mathcal{R} \geq \mathcal{T} \implies$ reject null hypothesis.

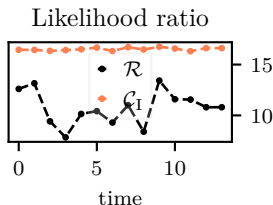
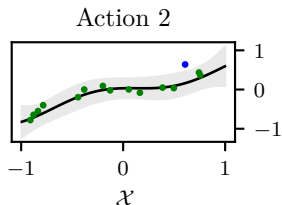
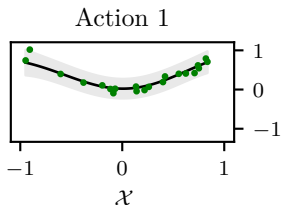
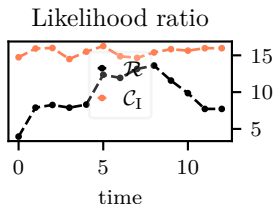
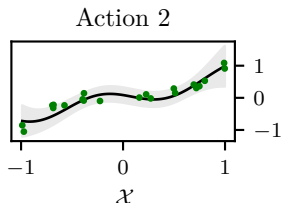
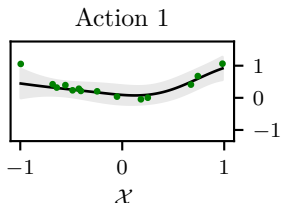
- Inference error:

- **Type I**: wrong detection of regime change.
- **Type II**: missed detection of regime change.

- Results: we choose \mathcal{T} to target a specific Type I or Type II inference error.

How to reset in non-noisy environments:MTGP bandit with $\sigma = 0$, $P = 20$, $p = 10$, $\delta_I = 0.4$, $\beta = 0.6$ $f^{a_1} : x \mapsto x^2$, $f^{a_2} : x \mapsto x^3$.

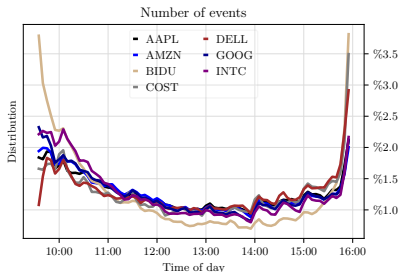
Learning in algorithmic trading

How to reset in noisy environmentsMTGP bandit with $\sigma = 0.2$, $P = 20$, $p = 10$, $\delta_l = 0.4$, $\beta = 0.6$ $f^{a_1} : x \mapsto x^2$, $f^{a_2} : x \mapsto x^3$.

Application: optimal timing with short term predictive signals

LOB data: 1 October 2022 → 31 December 2022

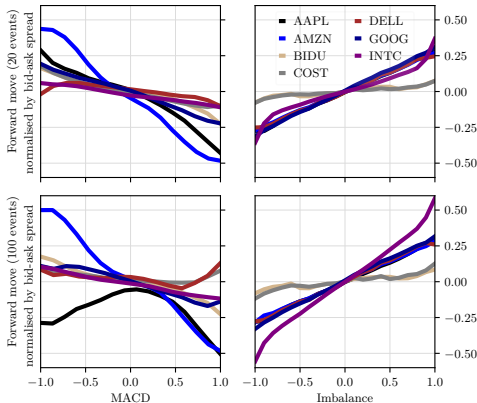
Ticker	Avg. spread	Avg. queue size	Avg. queue size	Events per minute
	(in ticks)	best bid	best ask	
AAPL	1.48	532.72	543.38	4472.36
AMZN	1.44	512.35	517.36	3939.42
BIDU	11.37	141.48	152.81	187.09
COST	24.62	73.39	72.33	177.15
DELL	1.60	396.65	390.34	361.98
GOOG	1.43	496.99	514.38	2766.57
INTC	1.16	5247.56	5197.65	1458.23



Short term predictive signals

■ Volume imbalance: $I_t = \frac{Q_t^B - Q_t^A}{Q_t^B + Q_t^A}$

■ MACD:
$$\begin{cases} \tilde{S}_t & = \text{EMA}^{\varepsilon_1}(S_t) - \text{EMA}^{\varepsilon_2}(S_t) \\ I_t^2 & = 10^5 \left(\tilde{S}_t - \text{EMA}^{\varepsilon_3}(\tilde{S}_t) \right) / S_{t-\varepsilon_2-\varepsilon_3}, \\ \text{EMA}^\varepsilon(x_t) & = \varepsilon x_t + (1 - \varepsilon) \text{EMA}(x_{t-\Delta t}) \end{cases}$$



Signals: predictive power depends on signal values

Ticker	Accuracy MACD	Accuracy IMB	Accuracy MACD extreme values	Accuracy IMB extreme values
AAPL	50.34%	62.84 %	60.24 %	77.25 %
AMZN	49.61 %	64.01 %	60.48 %	77.9 %
BIDU	49.23 %	52.33 %	50.38 %	58.16 %
COST	51.63 %	55.48 %	53.71 %	63.76 %
DELL	48.29 %	62.64 %	57.62 %	77.54 %
GOOG	48.77 %	64.82 %	60.97 %	79.27 %
INTC	52.58 %	81.93 %	72.80 %	98.03 %

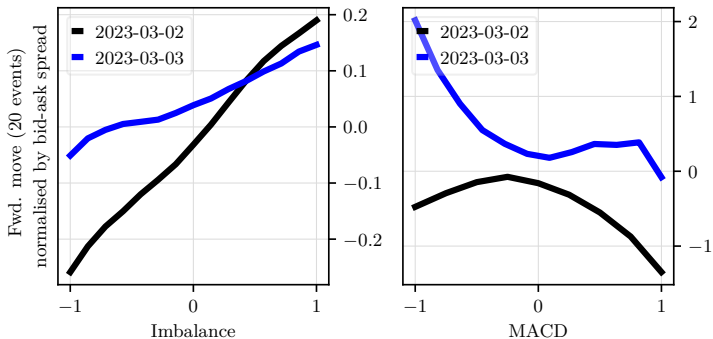
Table: Accuracy of MACD and imbalance computed as the hit ratio between signal predictions 20 events ahead and the realised price move.

Signals: noise

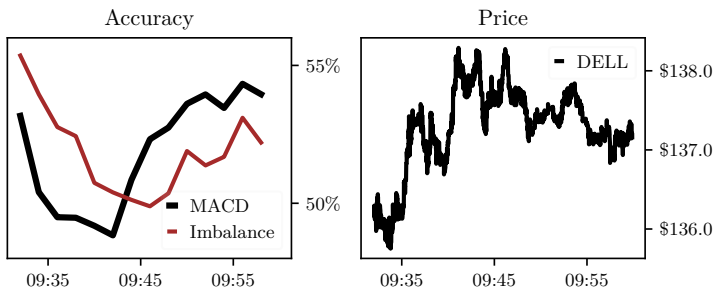
Mixing (very) noisy signals reduces statistical inference.

Ticker	Accuracy MACD	Accuracy IMB	Accuracy MACD+IMB
AAPL	50.34%	62.84 %	35.44 %
AMZN	49.61 %	64.01 %	34.21 %
BIDU	49.23 %	52.33 %	29.36 %
COST	51.63 %	55.48 %	31.58 %
DELL	48.29 %	62.64 %	33.01 %
GOOG	48.77 %	64.82 %	33.81 %
INTC	52.58 %	81.93 %	44.19 %

Signals: market regimes



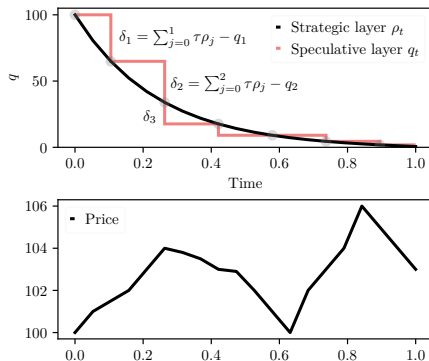
Signals: market regimes



Optimal timing with MTGP-LR

- **Context:** signal values.
- **Action k :** trade with an MO of size δ_t immediately, or at $t + 1$.
 - If $\delta_t \leq 0$, and $\mathbb{E}[S_{t+1}|I_t^k] \geq 0$, then sell at $t + 1$.
 - If $\delta_t \leq 0$, and $\mathbb{E}[S_{t+1}|I_t^k] \leq 0$, then sell at t .
- **Reward:** improvement in execution price over threshold.

$$\left(S_t - \tilde{S}_{t+1}^k \right) \times \text{sign}(\delta_t),$$



Short-term predictive signals in the literature

- Dynamics of the price and the signals

$$dS_t = \pi_t dt + \sigma dW_t + \kappa \nu_t$$

$$d\pi_t = \theta (\bar{\pi} - \pi_t) dt + \gamma dB_t .$$

- Short term (noisy) signals + long term objective.
- Non-stationarity: Continuous-time Markov Chain.
- A model for several signals needs prior knowledge :
 - Multivariate joint dynamics for the signals.
 - Exact form for how each signal drives the price.
 - Number of possible market regimes.
 - Transition probabilities between the regimes.

Performance study

Setup

- 7 securities from Nasdaq.
- 1000 execution programmes.
- Execution programme: 100 shares to liquidate in 10 minutes: buy/sell with probability 1/2).

Ablation study

- **GP-LR**: bandit without transfer learning (multi-output GP).
- **GP**: bandit without change-point detection and transfer learning (multi-output GP).
- **UCB**: bandit without contextual features, change-point detection, and transfer learning.

- **Imbalance**: only volume imbalance.
- **MACD**: only MACD.

Average performance in USD per 10^6 USD traded.

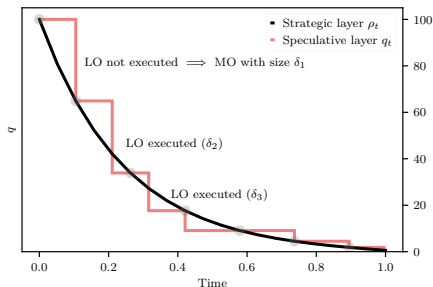
Ticker	Oracle	GP-LR	MTGP-LR	GP	UCB	Imbalance	MACD
BIDU	38.10	8.63	8.39	6.02	5.22	6.09	3.72
	±24.47	±18.21	±18.61	±18.37	±19.49	±19.24	±18.77
COST	16.44	4.56	4.96	4.66	3.15	3.76	2.15
	±12.03	±8.87	± 9.07	±9.42	±9.57	±10.22	±8.67
AAPL	18.03	4.74	5.18	4.32	4.02	4.09	1.90
	±6.49	±5.51	± 6.61	±5.01	±5.25	±5.04	±7.17
AMZN	11.72	4.70	4.74	4.04	2.80	4.70	0.70
	±2.66	±3.46	± 3.52	±4.61	±4.40	±3.81	±3.63
DELL	25.97	3.97	4.09	3.45	3.65	3.85	2.85
	±16.36	±10.96	± 10.84	±11.04	±11.64	±11.75	±11.16
GOOG	18.46	3.66	4.08	2.39	2.20	2.75	0.57
	±9.59	±7.11	± 7.13	±7.35	±7.95	±7.54	±7.70
INTC	4.47	2.53	2.78	2.69	1.24	3.56	-0.54
	±4.89	±3.89	±4.02	±4.23	±3.85	± 4.14	±3.77
Average	19.02	4.68	4.89	3.94	3.18	4.11	1.62

Application: optimal limit order placement

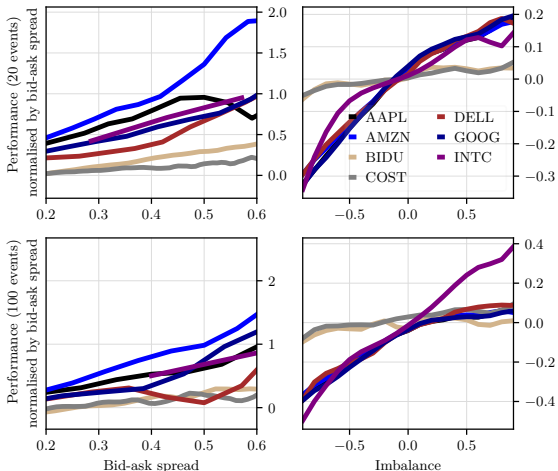
Optimal placement with MTGP-LR

- **Context:** bid-ask spread and imbalance.
- **Action k :** post limit order at the k^{th} limit at t , and send MO at $t + 1$ if the order is not executed.
- **Reward:** improvement in execution price over threshold (MO at t).

$$\left(S_t - \tilde{S}_{t+1}^k \right) \times \text{sign}(\delta_t),$$



First limit sell LOs



Learning with GPs

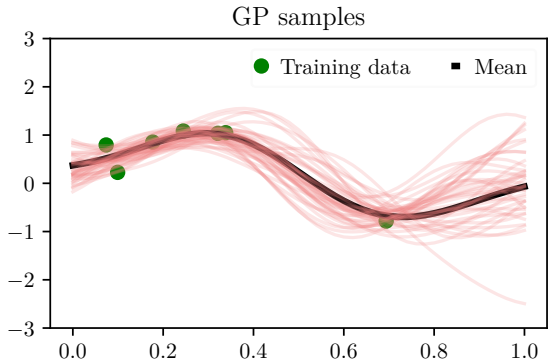
■ Inference:

$$\begin{cases} \mu_{post}(\mathbf{x}_*) & = \mathbf{k}(\mathbf{x}_*, \mathbf{X})(K + \sigma^2 I)^{-1} \mathbf{y}, \\ k_{post}(\mathbf{x}_*, \mathbf{x}'_*) & = k(\mathbf{x}_*, \mathbf{x}'_*) - \mathbf{k}(\mathbf{x}_*, \mathbf{X})(K + \sigma^2 I)^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}'_*) \end{cases},$$

Sampling with GPs

- Posterior sampling:

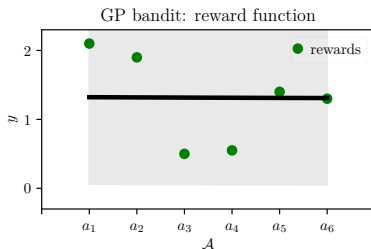
$$p(f_* | X_*, X, y, \theta) \sim \mathcal{N}(\mu_*, K_{*,*})$$



GP bandits (UCB)

Our case

- Space of actions is discrete (trading decisions for child orders).
- Similarity between all actions can be harmful (spurious transfer learning).
- Clustering effects between outcomes.



- Transfer learning only between actions that share common causal mechanisms.
- **Non stationarity**: the reward function f changes.

MTGP-LR

The setup

- \mathcal{A} : finite set of actions. \mathcal{X} : set of contextual features.
- $\mathcal{T} = \{1, \dots, T\}$: set of discrete observation times.
- The **reward function** $f_t : \mathcal{X} \times \mathcal{A} \mapsto \mathbb{R}$ is a sample from a MTGP.
- Non-stationarity: f_t is piecewise-stationarity
 - Finite number of change-points: $M = 1 + \sum_{t=1}^T \mathbb{1}_{\{f_t \neq f_{t-1}\}}$.
 - Observation times when the reward function changes: $\{\nu_j\}_{j \in \{0, \dots, M\}} \subset \mathcal{T}$.
 - Segments of stationarity: $[\nu_{j-1}, \nu_j]$, for $j \in \{1, \dots, M\}$.
- Bandit: $\{\mathcal{A}, \mathcal{X}, \mathcal{T}, (f_t)_{t \in \mathcal{T}}\}$.
- Objective: minimise the regret

$$R(T) = \sum_{t=1}^T r_t = \sum_{t=1}^T f_t(a_t^*, \mathbf{x}_t) - f_t(a_t, \mathbf{x}_t),$$

where $a_t^* = \arg \max_{a \in \mathcal{A}} f_t(a, \mathbf{x}_t)$.

The algorithm in stationary environments

Algorithm 4: MTGP-LR (stationary)

Data: Number of actions N , number of contextual features m , action space \mathcal{A} , context space \mathcal{X} , $\beta > 0$, kernel $k_{\theta}^{\mathcal{X}}$.

$\mathcal{D} \leftarrow \{\}$;

$f \sim \text{MTGP}(0, K_{\theta})$;

while $t \leq T$ **do**

- Observe vector of contexts $\mathbf{x}_t \in \mathcal{X}$;
- Compute vector $\mathbf{UCB} \leftarrow \left[\tilde{\mu}_{t-1}(a, \mathbf{x}_t) + \beta_t^{1/2} \tilde{\sigma}_{t-1}(a, \mathbf{x}_t) \right]_{a \in \mathcal{A}}$;
- $\mathbf{a}_t \leftarrow \arg \max_{a \in \mathcal{A}} \mathbf{UCB}_a$;
- Select action \mathbf{a}_t ; Sample reward $y_{\mathbf{a}_t, t}$;
- $\mathcal{D} += \{(y_t, \mathbf{x}_t, \mathbf{a}_t)\}$;
- Retrain the kernel hyper-parameters θ with data set \mathcal{D} ;

Non-stationary environments: Likelihood ratio test

- Let \mathcal{W} be the time window that contains the last P rewards.
- Let $\overline{\mathcal{W}} \subset \mathcal{W}$ be the sub-window containing the most recent $p \leq P$ points of \mathcal{W} , and let $\underline{\mathcal{W}} = \mathcal{W} \setminus \overline{\mathcal{W}}$.
- Hypothesis test:
 - Null hypothesis \mathbf{H}_0 : f_t in $\overline{\mathcal{W}}$ and the noise σ_0 are the same as in $\underline{\mathcal{W}}$.
 - Alternative hypothesis \mathbf{H}_1 : f_t in $\overline{\mathcal{W}}$ and the noise σ_1 are from a new MTGP.

- Likelihood ratio statistic

$$\mathcal{R} = 2 \log \frac{\rho(\overline{\mathbf{y}} | \mathbf{H}_1)}{\rho(\overline{\mathbf{y}} | \mathbf{H}_0)} = -\overline{\mathbf{y}}^T (\overline{\mathbf{K}} + \sigma_1^2 \mathbf{I})^{-1} \overline{\mathbf{y}} - \log |\overline{\mathbf{K}} + \sigma_1^2 \mathbf{I}|$$

$$+ (\overline{\mathbf{y}} - \tilde{\boldsymbol{\mu}})^T (\tilde{\mathbf{K}} + \sigma_0^2 \mathbf{I})^{-1} (\overline{\mathbf{y}} - \tilde{\boldsymbol{\mu}}) + \log |\tilde{\mathbf{K}} + \sigma_0^2 \mathbf{I}|.$$

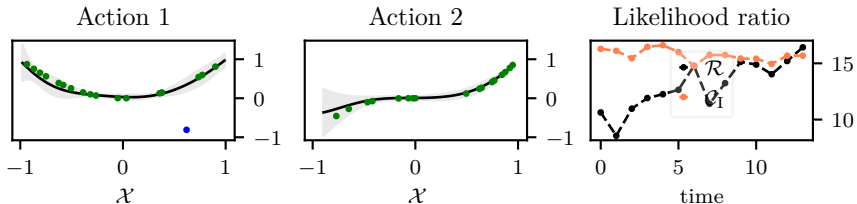
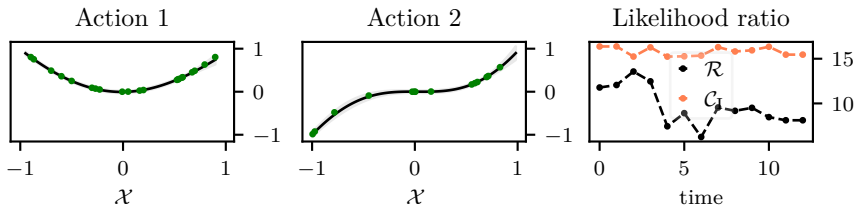
- Statistical test: $\mathcal{R} \geq \mathcal{T} \implies$ reject null hypothesis.
- Inference error:
 - **Type I**: wrong detection of regime change.
 - **Type II**: missed detection of regime change.

MTGP-LR Likelihood ratio

How to reset in non-noisy environments:

MTGP bandit with $\sigma = 0$, $P = 20$, $p = 10$, $\delta_I = 0.4$, $\beta = 0.6$

$f^{a_1} : x \mapsto x^2$, $f^{a_2} : x \mapsto x^3$.

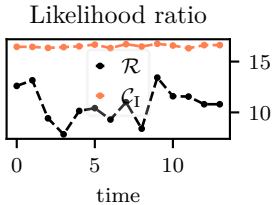
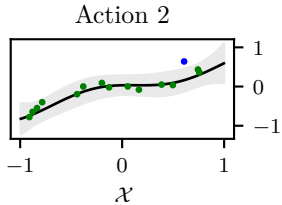
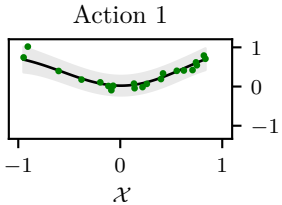
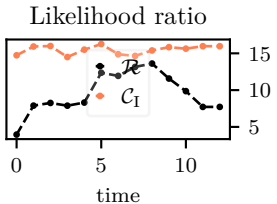
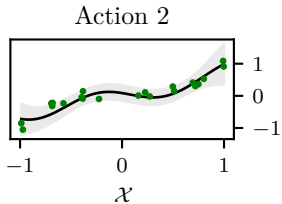
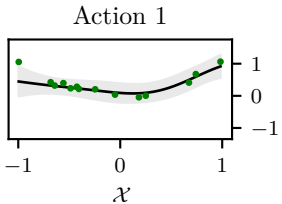


MTGP-LR Likelihood ratio

How to reset in noisy environments

MTGP bandit with $\sigma = 0.2$, $P = 20$, $p = 10$, $\delta_1 = 0.4$, $\beta = 0.6$

$f^{a_1} : x \mapsto x^2$, $f^{a_2} : x \mapsto x^3$.





Thank you for listening!

Any questions?

Why not RL

- bandit assumption: actions do not affect price dynamics or the predictive power of signals.
- Actions in the RL setting affect both rewards and states.
- RL training requires proper simulation and modeling of market impact.
- Impact is in the strategic layer.
- Possible to penalise reward with model-specific impact parameter.

Why not regression

- Predictive / descriptive power of features depends on the features values.
- Opposing signals can lead to cancellation of predictive power.
- Mixing (very) noisy signals in non stationary environments reduces statistical inference.

Computational complexity of GPs: many options

- Trailing period \implies maximum size for the matrix to inverse.
- Discretise the space of contextual features and group observations: reduces noise.
- Low-rank matrix approximation
 - Lanczos algorithm (iterative method to find the m most useful eigenvalues).
 - Control of the size of the low rank decomposition used for samples.

Pleiss, G., Gardner, J., Weinberger, K., Wilson, A. G. (2018, July). Constant-time predictive distributions for Gaussian processes.