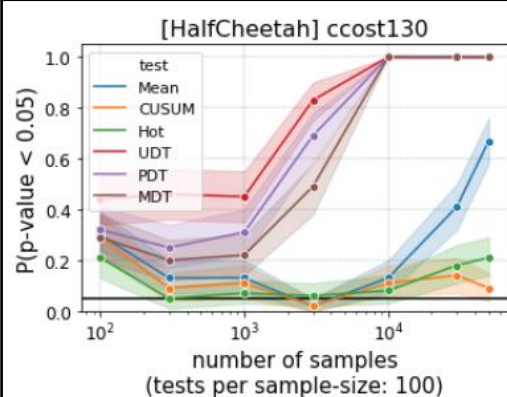


Problem

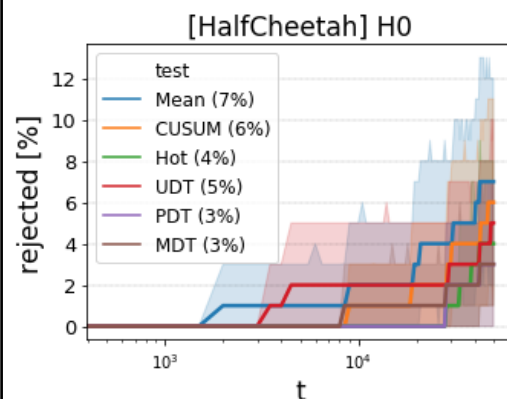
- Reinforcement Learning: learn a policy to make decisions in various states
- Reasonable real-world pipeline: train \rightarrow freeze \rightarrow test \rightarrow market
- What if the world changes / the agent falters for any reason?
- Re-exploring policies online is often forbidden / unsuccessful**
 - E.g. in autonomous driving: pass control to human driver
 - E.g. in insulin injector: send patient to the doctor
- Key is noticing performance degradation ASAP**
- Common assumptions (Markov, known model, etc.) do not hold

Framework and solution approach

- Focus on rewards (no need to learn a states-dependent model)
- Episodic framework: i.i.d episodes of length T
 - Rewards within episodes: **NOT [independent, identically-distributed, or Markov]**
 - Reference dataset: "valid" recorded episodes (e.g. the test period before marketing)
- Pseudo-algorithm: for every test-point (e.g. several times per episode):**
 - If the rewards of the last few episodes are small wrt the reference data – raise a flag**
- How to summarize the rewards (i.e. what is the **test-statistic**)?
- What is "small" (i.e. how to choose the **test-threshold**)?

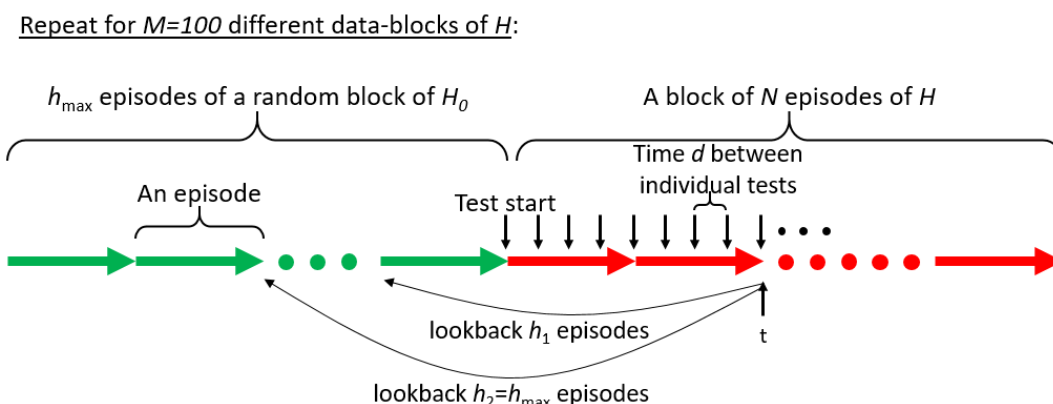


Due to sub-optimal handling of noise, standard tests sometimes become worse when getting more data



Threshold tuning: all 6 tests were successfully tuned to yield ~5% false-alarms within 50 episodes

Sequential test for a scenario H



Test statistic

- Naïve: mean reward (e.g. over last few episodes)
- Our suggestion: **weighted mean $W \cdot r$** ($W := \mathbf{1}^T \Sigma^{-1}$)
 - $\Sigma \in R^{T \times T}$ is the covariance matrix, estimated from the reference data
 - Intuition for independent rewards: $w_t = 1/\sigma_t^2 \Rightarrow$ look where it's less noisy
- Theorems:
 - Uniform degradation over time-steps + normal rewards \Rightarrow **optimal test** (max power)
 - Without normality \Rightarrow still better than simple mean
 - Advantage over simple mean increases with heterogeneity of Σ 's eigenvalues
- Experiments – how long it takes to detect degradation of rewards in modified environment?
 - Our test variants (UDT,PDT,MDT) usually win** simple mean, CUSUM and Hotelling

Test-threshold tuning

- Threshold-tuning by bootstrap usually relies on i.i.d data samples (for re-sampling)
- BFAR:** Bootstrap for False-Alarm Rate control that **can be applied to episodic (non-i.i.d) data**

