



Figure 1: Comparison with correlated prior sampling for the $K_{4,4}$ matching problem, with $c = 0, 0.5, 1$.

1 We want to thank the reviewers for their useful and positive feedbacks. We answer their comments/questions
 2 in the following paragraphs, that will be incorporated in the paper.

3 **Correlated vs independent prior** We briefly discussed the use of a correlated prior in future work and
 4 also in footnote 3, mentioning that the policy would perform better than using an independent prior. We ran
 5 additional empirical comparisons to assess this, plotting the results in Figure 1. As expected, the correlated
 6 prior policy is better (when outcomes are correlated). This motivates the theoretical study of such policy for
 7 future work. However, some new challenges are raised, as we see in the following paragraph.

8 **Importance of using a factorized prior in our analysis** A factorized prior allows us to bound the
 9 filtered regret against the event $\mathfrak{S}_t(Z) \wedge \mathfrak{T}_t(Z)$ (see p. 13, beginning of step 4). More precisely, we count the
 10 number of rounds needed for $\mathfrak{S}_t(Z) \wedge \neg\mathfrak{T}_t(Z)$ to occur for the q -th time. Under such event, the arms of Z
 11 are all observed, i.e., the corresponding priors are updated and counters are thus lower bounded by q . The
 12 importance of having a factorize prior resides in bounding the expected number of rounds for $\mathfrak{S}_t(Z) \wedge \neg\mathfrak{T}_t(Z)$
 13 to occur. Indeed, this is (essentially) equal to $\mathbb{E}[1/\mathbb{P}[\neg\mathfrak{T}_t(Z)|\mathfrak{S}_t(Z)]]$. Since $\mathfrak{S}_t(Z)$ and $\mathfrak{T}_t(Z)$ are independent,
 14 this is further equal to $\mathbb{E}[1/\mathbb{P}[\neg\mathfrak{T}_t(Z)]]$, allowing us to conduct our analysis by showing that this last quantity
 15 is exponentially decreasing in q (so summable over q). To the best of our knowledge, it is unknown how to
 16 get such a bound when $\mathfrak{S}_t(Z)$ and $\mathfrak{T}_t(Z)$ are not independent.

17 **Technical contributions for cts-beta** Although the gain in the upper bound for the beta prior might be
 18 considered as marginal, we want to stress on the asymptotic (quasi) optimality of the new bound (considering
 19 that the $\log^2(m)$ factor is negligible compared to n). We also fixed some technical issues in the proof of the
 20 previous bound; as a consequence, we believe that, although our work on CTS-BETA might appear somewhat
 21 incremental over Wang and Chen [2018], it brings essential clarifications to the literature.

22 **Technical contributions for cts-gaussian** Although the CTS-GAUSSIAN policies¹ are natural and essen-
 23 tially not new, our contribution is in their analyses, that are non-incremental. In particular, the stochastic
 24 dominance method is completely new, and allows us to convert correlated outcomes into independent Gaussian
 25 ones. Independence is crucial to be able to factorize the expectation $\mathbb{E}[1/\mathbb{P}[\neg\mathfrak{T}_t(Z)]]$. To the best of our
 26 knowledge, asymptotic quasi optimal analysis only exists for a restrictive class of action spaces (matroid) or
 27 outcomes distributions (independent). Dealing with both a general action space and outcome distribution is
 28 challenging, and represents our main contribution.

29 **Motivation for sub-Gaussian outcomes** In the same way as boundedness generalizes to sub-Gaussianity
 30 in 1d, we have that if \mathbf{X} is a.s. in a compact \mathcal{K} , it is \mathbf{C} -sub-Gaussian, with \mathbf{C} built from the John’s ellipsoid
 31 of \mathcal{K} . In this case, D_i is computed with a linear maximization over \mathcal{A} (see footnote 4). In particular,
 32 $\mathcal{K} = B_{\ell_\infty}(0, 1)$ gives $D_i = m$, and $\mathcal{K} = B_{\ell_2}(0, 1)$ gives $D_i = 1$. We can also use other structures on the
 33 outcomes to have D_i , such as negative dependence (as in our shortest path experiments).

34 **Comparison to previous work** We will add a comparison to Thompson Sampling for the MNL Bandit
 35 in the revised version, mentioning that the use of a common posterior distribution boosts the probability
 36 that the samples are all optimistic, and can thus greatly improve the constant term in our bound. However,
 37 as we saw, obtaining a quasi optimal gap dependent bound for such correlated sampling is an open question.

¹ $\beta > 1$ is an artefact of the analysis and can in practice be taken equal to 1 (as we did in our experiments).