A Proofs

A.1 Proof of Thm. 1

Suppose we split the actions into c cliques C_1, C_2, \ldots, C_c . First, let us consider the expected regret of the exponentially weighted forecaster ran over any such clique. Denoting the actions of the clique by $1, \ldots, n$, the forecaster works as follows: first, it initializes weights w_1, \ldots, w_n to be 1. At each round, it picks an action i with probability $w_i / \sum w_i$, receives the reward $g_i(t)$, and observes the noisy reward value $\hat{g}_j(t)$ for each of the other actions. It then updates $w_i = w_i \exp(\beta \hat{g}_i(t))$ (for some parameter $\beta \in (0, 1/b)$) for all $i = 1, \ldots, n$.

The analysis of this algorithm is rather standard, with the main twist being that we only observe unbiased estimates of the rewards, rather than the actual reward. For completeness, we provide this analysis in the following lemma.

Lemma 1. The expected regret of the forecaster described above, with respect to the actions in clique $|C_i|$ and under the optimal choice of the parameter β is at most $b\sqrt{\log(|C_i|)T}$.

Proof. We define the potential function $W_t = \sum_{j=1}^n w_j(t)$, and get that

$$\frac{W_{t+1}}{W_t} \leq \sum_{j=1}^n \frac{w_j(t)}{\sum_{l=1}^n w_l(t)} \exp(\beta \tilde{g}_j(t)).$$

For notational convenience, let $p_j(t) = \frac{w_j(t)}{\sum_{l=1}^n w_l(t)}$. Since $\tilde{g}_j(t) \leq b$, and $\beta \leq 1/b$, we have $\beta \tilde{g}_j(t) \leq 1$. Thus, we can use the inequality $\exp(x) \leq 1 + x + x^2$ (which holds for any $x \leq 1$), and get the upper bound

$$\sum_{j=1}^{n} p_j(t) \left(1 + \beta \tilde{g}_j(t) + 2\beta^2 \tilde{g}_j(t)^2 \right) = 1 + \beta \sum_{j=1}^{n} \tilde{g}_j(t) + \beta^2 \sum_{j=1}^{n} p_j(t) \tilde{g}_j(t)^2.$$

Taking logarithms and using the fact that $log(1 + x) \le x$, we get

$$\log\left(\frac{W_{t+1}}{W_t}\right) \leq \beta \sum_{j=1}^n p_j(t)\tilde{g}_j(t) + \beta^2 \sum_{j=1}^n p_j(t)\tilde{g}_j(t)^2.$$

Summing over all t, and canceling the resulting telescopic series, we get

$$\log\left(\frac{W_{T+1}}{W_{1}}\right) \leq \sum_{t=1}^{T} \left(\beta \sum_{j=1}^{n} p_{j}(t)\tilde{g}_{j}(t) + \beta^{2} \sum_{j=1}^{n} p_{j}(t)\tilde{g}_{j}(t)^{2}\right).$$
(6)

Also, for any fixed action *i*, we have

$$\log\left(\frac{W_{T+1}}{W_1}\right) \ge \log\left(\frac{w_i(T+1)}{W_1}\right) = \beta \sum_{t=1}^T \tilde{g}_i(t) - \log(n).$$

$$\tag{7}$$

Combining Eq. (6) with Eq. (7) and rearranging, we get

$$\sum_{t=1}^{T} \tilde{g}_i(t) - \sum_{t=1}^{T} \sum_{j=1}^{n} p_j(t) \tilde{g}_j(t) \le \frac{\log(n)}{\beta} + \beta \sum_{t=1}^{T} \sum_{j=1}^{n} p_j(t) \tilde{g}_j(t)^2.$$

Taking expectations on both sides, and using the facts that $\mathbb{E}[\tilde{g}_j(t)] = g_j(t)$ for all j, t, and $|\tilde{g}_j(t)| \le b$ with probability 1, we get

$$\sum_{t=1}^{T} g_i(t) - \sum_{t=1}^{T} \sum_{j=1}^{n} p_j(t) g_j(t) \le \frac{\log(n)}{\beta} + \beta b^2 T.$$

Thus, by picking $\beta = \sqrt{\log(n)/b^2T}$, we get that the expected regret is at most $b\sqrt{\log(n)T}$. \Box

Now, we define each such forecaster (one per clique C_i) as a meta-action, and run the EXP3 algorithm on the c meta-actions. By the standard guarantee for this algorithm (see corollary 3.2 in [4]), the expected regret incurred by that algorithm with respect to any fixed meta-action is at most $3b\sqrt{c \log(c)T}$. Combining this with Lemma 1, we get that the total expected regret of the ExpBan algorithm with respect to any single action is at most

$$\max_{i} b\sqrt{\log(|C_i|)T} + 3b\sqrt{c\log(c)T} \le b\sqrt{\log(k)T} + 3b\sqrt{c\log(k)T},$$

which is at most $4b\sqrt{\log(k)cT}$ since $c \ge 1$.

A.2 Proof of Thm. 2

To prove the theorem, we will need three lemmas. The first one is straightforward and follows from the definition of $\tilde{g}_j(t)$. The second is a key combinatorial inequality. We were unable to find an occurrence of this inequality in any previous literature, although we are aware of very special cases proven in the context of cyclic sums (see for instance [5]). The third lemma allows us to derive a more explicit bound by examining a particular choice of $\{s_i(t)\}_{i \in [k], t \in [T]}$.

Lemma 2. For all fixed t, j, we have

$$\mathbb{E}\left[\tilde{g}_j(t)\right] = g_j(t)$$

as well as

$$\mathbb{E}\left[\sum_{j=1}^{k} p_j(t)\tilde{g}_j(t)^2\right] \le b^2 \sum_{j=1}^{k} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)}$$

Proof. It holds that

$$\mathbb{E}\left[\tilde{g}_j^i(t)\right] = \sum_{i=1}^k p_i(t)\mathbb{E}[\tilde{g}_j(t) \mid \text{action i was picked}] = \sum_{i \in N_j(t)} p_i(t)\frac{g_j(t)}{\sum_{l \in N_j(t)} p_l(t)} = g_j(t).$$

As to the second part, we have

$$\mathbb{E}\left[\sum_{j=1}^{k} p_j(t)\tilde{g}_j(t)^2\right] = \sum_{i,j=1}^{k} p_j(t)p_i(t)\mathbb{E}\left[\tilde{g}_j(t)^2 \mid \text{action i was picked}\right]$$
$$\leq \sum_{j=1}^{k} \sum_{i \in N_j(t)} p_j(t)p_i(t)\frac{b^2}{\left(\sum_{l \in N_j(t)} p_l(t)\right)^2} = b^2 \sum_{j=1}^{k} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)}.$$

Lemma 3. Let G be a graph over k nodes, and let $\alpha(G)$ denote the independence number of G (i.e., the size of its largest independent set). For any $j \in [k]$, define N_j to be the nodes adjacent to node j (including node j). Let p_1, \ldots, p_k be arbitrary positive weights assigned to the node. Then it holds that

$$\sum_{i=1}^{k} \frac{p_i}{\sum_{l \in N_i} p_l} \leq \alpha(G).$$

Proof. We will actually prove the claim for any nonnegative weights p_1, \ldots, p_k (i.e., they are allowed to take 0 values), under the convention that if $p_j = 0$ and $\sum_{l \in N_j} p_i = 0$ as well, then $\sum_{i=1}^k p_i / \sum_{l \in N_i} p_i = 1$.

Suppose on the contrary that there exist some values for p_1, \ldots, p_k such that $\sum_{i=1}^k p_i / \sum_{l \in N_i} p_i > \alpha(G)$. Now, if p_1, \ldots, p_k are non-zero only on an independent set S, then

$$\sum_{i=1}^{k} \frac{p_i}{\sum_{l \in N_i} p_i} = \sum_{i \in S} \frac{p_i}{p_i} = |S|.$$

Since $|S| \leq \alpha(G)$, it follows that there exist some adjacent nodes r, s such that $p_r, p_s > 0$. However, we will show that in that case, we can only increase the value of $\sum_{i=1}^{k} p_i / \sum_{l \in N_i} p_i$ by shifting the entire weight $p_r + p_s$ to either node r or node s, and putting weight 0 at the other node. By repeating this process, we are guaranteed to eventually arrive at a configuration where the weights are non-zero on an independent set. But we've shown above that in that case, $\sum_{i=1}^{k} p_i / \sum_{l \in N_i} p_i \leq \alpha(G)$, so this means the value of this expression with respect to the original configuration was at most $\alpha(G)$ as well.

To show this, let us fix $p_r + p_s = c$ (so that $p_s = c - p_r$) and consider how the value of the expression changes as we vary p_r . The sum in the expression $\sum_{i=1}^k p_i / \sum_{l \in N_i} p_i$ can be split to 6 parts: when i = r, when i = s, when i is a node adjacent to s but not to r, when i is adjacent to r but not to s, when i is adjacent to both, and when i is adjacent to neither of them. Decomposing the sum in this way, so that p_r appears everywhere explicitly, we get

$$\begin{aligned} &\frac{p_r}{c+\sum_{l\in N_r\backslash r,s}p_l} + \frac{c-p_r}{c+\sum_{l\in N_j\backslash r,s}p_l} + \sum_{i:\{r,s\}\cap N_i=s}\frac{p_i}{c-p_r+\sum_{l\in N_i\backslash s}p_l} \\ &+\sum_{i:\{r,s\}\cap N_i=r}\frac{p_i}{p_r+\sum_{l\in N_i\backslash r}p_l} + \sum_{i:i\notin\{r,s\},r,s\subseteq N_i}\frac{p_i}{c+\sum_{l\in N_i\backslash\{r,s\}}p_l} + \sum_{i:\{r,s\}\cap N_i=\emptyset}\frac{p_i}{\sum_{l\in N_i}p_l}.\end{aligned}$$

It is readily seen that each of the 6 elements in the sum above is convex in p_r . This implies that the maximum of this expression is attained at the extremes, namely either $p_r = 0$ (hence $p_s = c$) or $p_r = c$ (hence $p_s = 0$). This proves that indeed shifting weights between adjacent nodes can only increase the value of $\sum_{i=1}^{k} p_i / \sum_{l \in N_i} p_i$, and as discussed earlier, implies the result stated in the lemma.

Lemma 4. Consider a graph G over nodes 1, ..., k, and let $\alpha(G)$ be its independence number. For any $j \in [k]$, define N_j to be the nodes adjacent to node j (including node j). Then there exist values of $s_1, ..., s_k$ on the k-simplex, such that

$$\frac{1}{\min_{j \in [k]} \sum_{l \in N_j} s_l} \le \alpha(G).$$
(8)

Proof. Let S be a largest independent set of G, so that $|S| = \alpha(G)$. Consider the following specific choice for the values of s_1, \ldots, s_k : For any j such that $j \in S$, let $s_j = 1/\alpha(G)$, and $s_j = 0$ otherwise. Suppose there was some node j such that $\sum_{l \in N_j} s_l = 0$. By the way we chose values for s_1, \ldots, s_k , this implies that node j is not adjacent to any node in S, so $S \cup \{j\}$ would also be an independent set, contradicting the assumption that S is a largest independent set. But since each value of s_l is either 0 or $1/\alpha(G)$, it follows that $\sum_{l \in N_j} s_l > 1/\alpha(G)$. This is true for any node j, from which Eq. (8) follows.

We now turn to the proof of the theorem itself.

Proof of Thm. 2. With the key lemmas at hand, most of the remaining proof is rather similar to the standard analysis for multi-armed bandits (e.g., [4]). We define the potential function $W_t = \sum_{j=1}^{k} w_j(t)$, and get that

$$\frac{W_{t+1}}{W_t} \leq \sum_{j=1}^k \frac{w_j(t)}{\sum_{l=1}^k w_l(t)} \exp(\beta \tilde{g}_j(t)).$$
(9)

We have that $\beta \tilde{g}_j(t) \leq 1$, since by definition of β and $\tilde{g}_j(t)$,

$$\beta \tilde{g}_j(t) \leq \frac{\beta b}{\sum_{l \in N_j(t)} p_l(t)} \leq \frac{\beta b}{\sum_{l \in N_j(t)} \gamma(t) s_l(t)} = \frac{\beta b}{\sum_{l \in N_j(t)} s_l(t)} \frac{\min_{j \in [k]} \sum_{l \in N_j(t)} s_l(t)}{\beta b} \leq 1.$$

Using the definition of $p_j(t)$ and the inequality $\exp(x) \le 1 + x + x^2$ for any $x \le 1$, we can upper bound Eq. (9) by

$$\sum_{j=1}^{k} \frac{p_j(t) - \gamma(t)s_j(t)}{1 - \gamma(t)} \left(1 + \beta \tilde{g}_j(t) + \beta^2 \tilde{g}_j(t)^2\right)$$

$$\leq 1 + \frac{\beta}{1 - \gamma(t)} \sum_{j=1}^{k} p_j(t)\tilde{g}_j(t) + \frac{2\beta^2}{1 - \gamma(t)} \sum_{j=1}^{k} p_j(t)\tilde{g}_j(t)^2.$$

Taking logarithms and using the fact that $\log(1+x) \leq x$, we get

$$\log\left(\frac{W_{t+1}}{W_t}\right) \leq \frac{\beta}{1-\gamma(t)} \sum_{j=1}^k p_j(t) \tilde{g}_j(t) + \frac{\beta^2}{1-\gamma(t)} \sum_{j=1}^k p_j(t) \tilde{g}_j(t)^2.$$

Summing over all t, and canceling the resulting telescopic series, we get

$$\log\left(\frac{W_{T+1}}{W_1}\right) \leq \sum_{t=1}^T \sum_{j=1}^k \frac{\beta}{1-\gamma(t)} p_j(t) \tilde{g}_j(t) + \sum_{t=1}^T \sum_{j=1}^k \frac{\beta^2}{1-\gamma(t)} p_j(t) \tilde{g}_j(t)^2.$$
(10)

Also, for any fixed action i, we have

$$\log\left(\frac{W_{T+1}}{W_1}\right) \ge \log\left(\frac{w_i(T+1)}{W_1}\right) = \beta \sum_{t=1}^T \tilde{g}_i(t) - \log(k).$$

$$(11)$$

Combining Eq. (10) with Eq. (11) and rearranging, we get

$$\beta \sum_{t=1}^{T} \tilde{g}_i(t) - \sum_{t=1}^{T} \sum_{j=1}^{k} \frac{\beta}{1 - \gamma(t)} p_j(t) \tilde{g}_j(t) \le \log(k) + \sum_{t=1}^{T} \sum_{j=1}^{k} \frac{\beta^2}{1 - \gamma(t)} p_j(t) \tilde{g}_j(t)^2.$$

Taking expectations on both sides, and using Lemma 2, we get

$$\beta \sum_{t=1}^{T} g_i(t) - \sum_{t=1}^{T} \sum_{j=1}^{k} \frac{\beta}{1 - \gamma(t)} p_j(t) g_j(t) \leq \log(k) + \sum_{t=1}^{T} \sum_{j=1}^{k} \frac{b^2 \beta^2}{1 - \gamma(t)} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)} dt$$

After some slight manipulations, and using the fact that $g_j(t) \in [0, 1]$ for all j, t, we get

$$\sum_{t=1}^{T} g_i(t) - \sum_{t=1}^{T} \sum_{j=1}^{k} p_j(t) g_j(t) \leq \sum_{t=1}^{T} \gamma(t) + \frac{\log(k)}{\beta} + \sum_{t=1}^{T} \frac{b^2 \beta}{1 - \gamma(t)} \sum_{j=1}^{k} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)}.$$

We note that $1/(1 - \gamma(t))$ can be upper bounded by 2, since by definition of $s_i(t)$,

$$\gamma(t) = \frac{\beta b}{\max_{a_1,\dots,a_k} \min_{j \in [k]} \sum_{l \in N_j(t)} a_l(t)} \le \frac{\beta b}{\min_{j \in [k]} \sum_{l \in N_j(t)} (1/k)} \le \beta bk \le 1/2.$$

Plugging this in as well as our choice of $\gamma(t)$ in the $\sum_t \gamma(t)$ term, and slightly simplifying, we get the upper bound

$$\sum_{t=1}^{T} g_i(t) - \sum_{t=1}^{T} \mathbb{E}[g_{i_t}(t)] \le \beta b^2 \left(\sum_{t=1}^{T} \frac{1}{\min_{j \in [k]} \sum_{l \in N_j(t)} s_l(t)} + 2 \sum_{j=1}^{k} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)} \right) + \frac{\log(k)}{\beta}.$$
(12)

Now, we recall that the $\{s_i(t)\}$ terms were chosen so as to minimize the bound above. Thus, we can upper bound it by any fixed choice of $\{s_i(t)\}$. Invoking Lemma 4, as well as Lemma 3, the theorem follows.

A.3 Proof of Thm. 3

The proof is very similar to the one of Thm. 2, so we'll only point out the differences.

Referring to the proof of Thm. 2 in Subsection A.2, The analysis is identical up to Eq. (12). To upper bound the terms there, we can still invoke Lemma 4. However, Lemma 3, which was used to upper bound $\sum_{j=1}^{k} p_j(t) / \sum_{l \in N_j(t)} p_l(t)$, not longer applies (in fact, one can show specific counter-examples). Thus, in lieu of Lemma 3, we will opt for the following weaker bound: Let $C_1, \ldots, C_{\bar{\chi}(G_t)}$ be a smallest possible clique partition of G_t . Then we have

$$\sum_{i=1}^{\bar{\chi}(G_t)} \sum_{j \in C_i} \frac{p_j(t)}{\sum_{l \in N_j(t)} p_l(t)} \le \sum_{i=1}^{\bar{\chi}(G_t)} \sum_{j \in C_i} \frac{p_j(t)}{\sum_{l \in C_i} p_l(t)} = \bar{\chi}(G_t).$$

Plugging this upper bound as well as Lemma 4 into Eq. (12), and using the fact that $\alpha(G_t) \leq \overline{\chi}(G_t)$ for any graph G_t , the result follows.

A.4 Proof of Theorem 4

Suppose that we are given a graph G with an independence number $\alpha(G)$. Let \mathcal{N} denote an independent set of $\alpha(G)$ nodes (i.e., no two nodes are connected). Suppose we have an algorithm \mathcal{A} with a low expected regret for every sequence of rewards. We will use this algorithm to form an algorithm for the standard multi-armed bandits problem (with no-side observations). We will then resort to the known lower bound for this problem, to get a lower bound for our setting as well.

Consider first a standard multi-armed bandits game on $\alpha(G)$ actions (with no side-observations), with the following randomized strategy for the adversary: the adversary picks one of the $\alpha(G)$ actions uniformly at random, and at each round, assigns it a random Bernoulli reward with parameter $1/2 + \epsilon$ (where ϵ will be specified later). The other actions are assigned a random Bernoulli reward with parameter 1/2. Roughly speaking, Theorem 6.11 of [6] shows that with this strategy and for $\epsilon = \Theta(\sqrt{\alpha(G)/T})$, the expected regret of any learning algorithm is at least $\Omega(\sqrt{\alpha(G)T})$.

Now, suppose that for the setting with side-observations, played over the graph G, there exists a learning strategy \mathcal{A} that achieves expected cumulative regret of at most $R_{\mathcal{A}}(T)$, for the graph G over T rounds, with respect to any adversary strategy. We will now show how to use \mathcal{A} for the standard multi-armed bandits game described above. To that end, arbitrarily assign the $\alpha(G)$ actions to the $\alpha(G)$ independent nodes in \mathcal{N} . We will then implement the following strategy \mathcal{A}' : whenever \mathcal{A} chooses one of the actions in \mathcal{N} , we choose the corresponding action in the multi-armed bandits problem and feed the reward back to \mathcal{A} (the reward of all neighboring nodes is 0, which we feed back to \mathcal{A} as well). Whenever \mathcal{A} chooses a node j not in \mathcal{N} , we use the next $|N_j \cap \mathcal{N}|$ rounds (where N_j is the neighborhood set of j) to do "pure exploration:" we go over all the neighbors of node j that belong to \mathcal{N} in some fixed order, and choose each of them once (since rewards are assumed stochastic the order does not matter). Nodes in $N_j \setminus \mathcal{N}$ are known to yield a reward of 0. The rewards of node j and all its neighbors are then fed to \mathcal{A} , as if they were side observations obtained in a single round by choosing a node not in \mathcal{N} . Since the rewards are chosen i.i.d., the distribution of these rewards is identical to the case where \mathcal{A} was really implemented with side-observations. We denote $R_{\mathcal{A}'}(T)$ as the expected regret of this strategy \mathcal{A}' , after T rounds.

We make the following observation: suppose A achieves an expected regret satisfying

$$R_{\mathcal{A}}(T) \le \sqrt{\alpha(G)}T$$

(we can assume this since our goal is to provide a lower bound which will only be smaller). Then the number of times \mathcal{A} chose actions outside \mathcal{N} must be smaller than $2\sqrt{\alpha(G)T}$. This is because whenever \mathcal{A} chooses an action not in \mathcal{N} it receives a reward of 0 while the highest expected reward is bigger than 1/2, so the expected per-round regret would increase by at least 1/2.

We apply \mathcal{A}' at each round, till \mathcal{A} is called T times. Let T' be the (possibly random) number of rounds which elapsed. It holds that $T' \geq T$, since we have the T' - T pure exploration rounds where \mathcal{A} is not called. In these exploration rounds, we pull arms in \mathcal{N} , so our expected regret in those rounds is at most ϵ . Moreover, by the observation above, the number of such rounds is at most $2\alpha(G)\sqrt{\alpha(G)T}$, since \mathcal{A} may choose an action outside \mathcal{N} at most $2\sqrt{\alpha(G)T}$ times, and this

follows by at most $|\mathcal{N}| = \alpha(G)$ pure exploration steps. In rounds where we do not do exploration steps, the expected per-round regret of \mathcal{A}' is the same as the expected per-round regret of \mathcal{A} . Overall, this implies that

$$R_{\mathcal{A}'}(T') \le R_{\mathcal{A}}(T) + 2\epsilon\alpha(G)\sqrt{\alpha(G)T}$$
(13)

Since the expected regret is monotone in the number of rounds, we can lower bound $R_{\mathcal{A}'}(T')$ by $R_{\mathcal{A}'}(T)$. Rearranging, we get

$$R_{\mathcal{A}}(T) \ge R_{\mathcal{A}'}(T) - 2\epsilon\alpha(G)\sqrt{\alpha(G)T}.$$

Now, \mathcal{A}' is a strategy for the standard multi-armed bandits setting, with a randomized adversary strategy which is identical to the one used to establish the lower bound of [6, Theorem 6.11]. Using this lower bound, by selecting $\epsilon = \sqrt{c_1 \alpha(G)/T}$ with $c_1 = 1/(8 \ln(4/3))$, we obtain

$$R_{\mathcal{A}}(T) \ge \sqrt{T\alpha(G)}c_2 - 2\sqrt{c_1}\alpha(G)^2, \tag{14}$$

where the first term of the right hand side comes from Page 168 in [6] and

$$c_2 = \frac{\sqrt{2} - 1}{\sqrt{32\ln(4/3)}}.$$

Since $T \ge 16\alpha(G)^3 c_1/c_2^2$, we have that $R_A(T) \ge \sqrt{T\alpha(G)}c_2/2$. Plugging in the values of c_1, c_2 above, the result follows.

Finally, we note that if the maximal degree of any node in G is bounded by d, then Eq. (13) can be improved to

$$R_{\mathcal{A}'}(T') \le R_{\mathcal{A}}(T) + 2\epsilon d\sqrt{\alpha(G)T},$$

since the number of pure-exploration steps following a call to \mathcal{A} is at most d rather than $\alpha(G)$. Repeating the analysis above, we get that Eq. (14) is replaced by

$$R_{\mathcal{A}}(T) \ge \sqrt{T\alpha(G)}c_2 - 2\sqrt{c_1}d\alpha(G).$$

This allows us to give the same lower bound, for any $T \ge 16\alpha(G)d^2c_1/c_2^2$, as opposed to $T \ge 16\alpha(G)^3c_1/c_2^2$ as before.