# Sensor Acquisition with no Feedback

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

We propose a sensor acquisition problem (SAP) wherein sensors (and sensing tests) are organized into a cascaded architecture of increasingly informative and expensive tests and the goal is to choose a test with the optimal cost-accuracy tradeoff for a given instance. We consider the case where we obtain no feedback in terms of rewards for our chosen actions apart from test observations. Absence of feedback raises fundamentally new challenges since one cannot infer potentially optimal tests. We pose the problem in terms of competitive optimality with the goal of minimizing cumulative regret against optimally chosen actions in hindsight. In this context we introduce the notion of weak dominance and show that it is necessary and sufficient for realizing sub-linear regret. Weak dominance on a cascade supposes that a child node in the cascade has higher accuracy when its parent node makes correct predictions. When weak dominance holds we show that we can reduce SAP to a corresponding multi-armed bandit problem with side observations. Empirically we verify that weak dominance holds for many datasets.

## 1 Introduction

In many classification systems such as medical diagnosis and homeland security, sequential decisions are often warranted. For each instance, an initial diagnostic test is conducted and based on its results further tests maybe conducted. Tests have varying costs for acquisition, and these costs account for delay, throughput or monetary value [1]. Apart from these natural scenarios the problem also arises in the context of wireless communication systems, where a cascade of error-correcting decoders of increasing block lengths are designed to overcome channel noise.

Our goal is essentially a sensor acquisition problem (SAP), namely, to acquire the tests/sensors that achieves the optimal cost-accuracy tradeoff for that instance. We assume that the sensors/tests are organized into a diagnostic cascade architecture, where the ordering is based on costs/informativity of tests. Each stage in the cascade outputs a prediction of the underlying state of the instance (disease or disease-free, threat or no-threat etc.). We suppose that the classifiers (or predictors) corresponding to each node are part of the system and produce labeled outputs. This is often the case in diagnostic systems where a test ordering is a priori known and a report is produced by a human being or an automated mechanism corresponding to different sensor measurements. Thus our task in this paper is primarily to learn a decision rule to identify the collection of tests required for an instance.

Our problem can be framed as a version of a multi-armed bandit problem. Each arm of the bandit corresponds to a unique path from root to a node where the observation is a vector of outputs from

---

[1] As described in Trapeznikov et al. (2014) security systems utilize a suite of sensors/tests such as X-rays, millimeter wave imagers (expensive & low-throughput), magnetometers, video, IR imagers human search. Security systems must maintain a throughput constraint in order to keep pace with arriving traffic. In clinical diagnosis, doctors in the context of breast cancer diagnosis utilize tests such as genetic markers, imaging (CT, ultrasound, elastography) and biopsy. Sensors providing imagery are scored by humans. The different sensing modalities have diverse costs, in terms of health risks (radiation exposure) and monetary expense.

tests acquired along that path. Nevertheless, our problem is unconventional. Unlike a conventional bandit problem, where feedback (reward) is observed corresponding to each action, we do not get feedback of how well our action performed (either noisy or noiseless)[2].

Absence of reward information associated with chosen actions is fundamentally challenging since we cannot infer potential optimal actions. We pose the problem in terms of competitive optimality. In particular we consider a competitor who has the benefit of hindsight and can choose an optimal collection of tests for all the examples. Our goal is to choose an action for each instance so that the cummulative regret with respect to the competitor is sub-linear (and optimal).

We first provide negative results for the problem. We introduce the notion of weak dominance on tests. We show that weak dominance is fundamental, i.e., regardless of the algorithm, if this condition is not satisfied, we are left with a linear regret. On the other hand we develop UCB style algorithms that show that we can realize optimal regret (sub-linear regret) guarantees when the condition is satisfied. This leads to a sharp necessary and sufficient condition for learning under no feedback.

The weak dominance condition amounts to a stochastic ordering of the tests on the diagnostic cascade. Conceptually, the weak dominance condition says that the child node tends to be relatively more accurate when the parent is correct. Under weak dominance we show that the learner can partially infer losses of the stages. In particular, we reduce the SAP problem to a stochastic multi-armed bandit with side observations, where bandit arms are identified by the nodes of the cascade. The payoff of an arm is given by loss from the corresponding stage, and side observation structure is defined by the feedback graph induced by the cascade. Empirically we verify that weak dominance condition naturally holds for several datasets including breast-cancer and diabetes datasets. A stronger dominance condition is also shown to hold by design, namely, for error-correcting code cascades in the context of communication systems.

Related Work: Trapeznikov & Saligrama (2013)Seldin et al. (2014)

Structure of paper

## 2   Sensor Acquisition Problem

The learner has access to $K \geq 2$ sensors that are ordered in terms of their prediction efficiency. Specifically, we consider that the sensors form a cascade (order in which the sensors are selected is predetermined) and in each round the learner can sequentially select a subset of sensors in the cascade and stop at any depth.



Figure 1: Cascade of sensors

Let $\{Z_t, Y_t\}_{t>0}$ denote a sequence generated according to an unknown distribution. $Z_t \in \mathcal{C} \subset \mathcal{R}^d$, where $\mathcal{C}$ is a compact set, denotes a feature vector/context at time $t$ and $Y_t \in \{0, 1\}$ its binary label. We denote output/prediction of the $i^{th}$ sensor as $\hat{Y}_t^i$ when its input is $Z_t$. The set of actions available to the learner is $\mathcal{A} = \{1, \ldots, K\}$, where the action $k \in \mathcal{A}$ indicates acquiring predictions from sensors $1, \ldots, k$ and classifying using the prediction $\hat{Y}_t^k$.

The prediction error rate of the $i^{th}$ sensor is denoted as $\gamma_i := \Pr\{Y_t \neq \hat{Y}_t^k\}$. In this section we assume that the error rate does not depend on the context and postpone the treatment with contextual information to Section C. Further, the sensors are arranged such that the prediction error rate improves with depth in the cascade, i.e., $\gamma_{k-1} \geq \gamma_k$ for all $k > 2$. However, the learner incurs an extra cost of $c_k \geq 0$ to acquire output of sensor $k$ after acquiring output of sensor $k - 1$. The sensor cascade is depicted in the adjacent figure.

Let $H_t(k)$ denote the feedback observed in round $t$ from action $k$. Since we observe predictions of all the first $k$ senors by playing action $k$, we get $H_t(k) = \{\hat{Y}_t^1, \ldots, \hat{Y}_t^k\}$. The loss incurred in each round is defined in terms of the prediction error and the total cost involved. When the learner selects action $k$, loss is the prediction error of sensor $k$ plus sum of the costs incurred along path

---
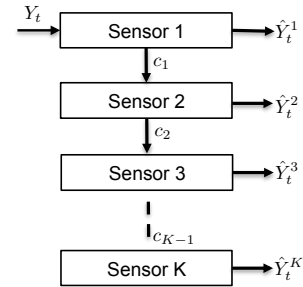
[2]This problem naturally arises in the surveillance and medical domains. We can perform a battery of tests on an individual in an airport but can never be sure whether or not he/she poses a threat.

$(c_1, \ldots, c_k)$. Let $L_t : \mathcal{A} \to \mathcal{R}_+$ denote the loss function in round $t$. Then,

$$L_t(k) = \mathbf{1}_{\{\hat{Y}_t^k \neq Y_t\}} + \sum_{j=1}^{k} c_j. \tag{1}$$

We refer to the above setup as Sensor Acquisition Problem (SAP) and denote it as $\psi = (K, \mathcal{A}, (\gamma_i, c_{i-1})_{i \in [K]})^3$. A policy $\pi^\psi = (\pi_1^\psi, \pi_2^\psi, \cdots)$ on $\psi$, where $\pi_t^\psi : \mathcal{H}_{t-1} \to \mathcal{A}$, gives action selected in each round using history $\mathcal{H}_{t-1}$ that consists of all actions and corresponding feedback observed before $t$. Let $\Pi^\psi$ denote set of policies on $\psi$. For any $\pi \in \Pi^\psi$, we compare its performance with respect to the optimal policy (single best action in hindsight) and define its expected regret as follows

$$R_T^\psi(\pi) = \mathbb{E}\left[\sum_{t=1}^{T} L_t(a_t)\right] - \min_{k \in A} \mathbb{E}\left[\sum_{t=1}^{T} L_t(k)\right], \tag{2}$$

where $a_t$ denotes the policy selected by $\pi_t$ in round $t$. The goal of the learner is to learn a policy that minimizes the expected total loss, or, equivalently, to minimize the expected regret, i.e.,

$$\pi^* = \arg\min_{\pi \in \Pi^\psi} R_T^\psi(\pi). \tag{3}$$

**Optimal action in hindsight:** For any $t$, we have

$$\mathbb{E}[L_t(k)] = \Pr\{Y_t \neq \hat{Y}_t^k\} + \sum_{j=1}^{k} c_j = \gamma_k + \sum_{j=1}^{k} c_j. \tag{4}$$

Let $k^* = \arg\min_{k \in \mathcal{A}} \gamma_k + \sum_{i<k} c_i$. Then the optimal policy is to play action $k^*$ in each round. If an action $i$ is played in any round then it adds $\Delta_k := \gamma_k + \sum_{i<k} c_i - (\gamma_{k^*} + \sum_{i<k^*} c_i)$ to the expected regret. Let $I_t$ denote the action selected in round $t$ and $N_k^\psi(s)$ denote the number of times action $k$ is selected till time $s$, i.e., $N_k^\psi(s) = \sum_{t=1}^{s} \mathbf{1}_{\{I_t=k\}}$. Then the expected regret can be expressed as

$$R_T^\psi(\pi) = \sum_{k \in \mathcal{A}} \mathbb{E}[N_k^\psi(T)]\Delta_k. \tag{5}$$

# 3 When is SAP Learnable?

In the SA-Problem feedback $H_t(\cdot)$ does not reveal any information about the true label $Y_t$ in any round $t$. Hence the loss values are not known, and we are in a hopeless situation where linear regret is unavoidable. In this section we explore conditions that lead to policies that are Hannan consistent Hannan (1957), i.e, a policy $\pi \in \Pi^\psi$ such that $R_T^\psi(\pi)/T \to 0$.

Let us consider $K = 2$ sensors and start with a simple condition that if sensor 1 predicts label 1 correctly, then sensor 2 also predicts it correctly[4], i.e.,

$$Y_t = 1 \text{ and } \hat{Y}_t^1 = 1 \implies \hat{Y}_t^2 = 1. \tag{6}$$

To fix ideas, we enumerate all possible 8 tuples $(Y, \hat{Y}^1, \hat{Y}^2)$ as shown in Table 3, and write probability of $i$th tuple $i = 1, 2, \cdots 8$ as $p_{i-1}$. From Table 3, we have $\gamma_1 = p_2 + p_3 + p_4 + p_5$ and $\gamma_2 = p_1 + p_3 + p_4 + p_6$, thus

$$\gamma_1 - \gamma_2 = p_2 + p_5 - p_1 - p_6. \tag{7}$$

---

[3]Note that $k \in \mathcal{A}$ implies that action $k$ selects all sensors $1, 2, \cdots, k$, not just sensor $k$. We set $c_0 = 0$

[4]Suppose we interpret label 1 as 'threat', the condition implies that if sensor 1 detects threat correctly, the better sensor 2 also detects it.

| $Y$ | $\hat{Y}^1$ | $\hat{Y}^2$ | $\Pr(Y, \hat{Y}^1, \hat{Y}^2)$ |
|---|---|---|---|
| 0 | 0 | 0 | $p_0$ |
| 0 | 0 | 1 | $p_1$ |
| 0 | 1 | 0 | $p_2$ |
| 0 | 1 | 1 | $p_3$ |
| 1 | 0 | 0 | $p_4$ |
| 1 | 0 | 1 | $p_5$ |
| 1 | 1 | 0 | $p_6$ |
| 1 | 1 | 1 | $p_7$ |

$$\Pr(\hat{Y}^1, \hat{Y}^2) = \begin{cases} p_1 + p_5 \text{ if } (\hat{Y}^1, \hat{Y}^2) = (0,1) \\ p_2 + p_6 \text{ if } (\hat{Y}^1, \hat{Y}^2) = (1,0) \\ p_0 + p_4 \text{ if } (\hat{Y}^1, \hat{Y}^2) = (0,0) \\ p_3 + p_7 \text{ if } (\hat{Y}^1, \hat{Y}^2) = (1,1) \end{cases} \tag{8}$$

Since we only observe feedbacks $(\hat{Y}_t^1, \hat{Y}_t^2)$ and not the true labels $Y_t$, only marginal probabilities $\Pr(\hat{Y}^1, \hat{Y}^2)$ as given in (8) can be estimated but not $\Pr(Y, \hat{Y}^1, \hat{Y}^2)$. The following example demonstrate marginals do not unambiguously decide optimal action.

Set $c = 0.35$ and consider the following two cases: 1) $p_2 = 1/2, p_1 = 1/4 - 1/40, p_5 = 1/4 + 1/40$ and 2) $p_2 = 1/2, p_1 = 1/4 - 3/40, p_5 = 1/4 + 3/40$. From condition (6) we have $p_6 = 0$. Also, set $p_0 = p_4 = p_3 = p_7 = 0$ in both the cases. We get $\gamma_1 - \gamma_2 = 0.3$ in the first case, whereas $\gamma_1 - \gamma_2 = 0.4$ in the second case. From 4, optimal action is 1 in the first case, whereas it is 2 in the second case. However, for both the cases the marginals $\Pr(\hat{Y}^1, \hat{Y}^2)$ are the same for all pairs $(\hat{Y}^1, \hat{Y}^2)$. Since we only observe $\Pr(\hat{Y}^1, \hat{Y}^2)$, the two cases cannot be distinguished and linear regret is unavoidable.

Next, assume that if sensor 0 predicts label 0 correctly, then sensor 2 also predicts it correctly, i.e.,

$$Y_t = 0 \text{ and } \hat{Y}_t^1 = 0 \implies \hat{Y}_t^2 = 0. \tag{9}$$

Similar to the previous example, we can argue that this conditions is not sufficient to achieve sub-linear regret. Now assume that both (6) and (9) hold. Then, $p_1 = p_6 = 0$ and we get $\gamma_1 - \gamma_2 = p_2 + p_5$. Since $p_5 + p_2 = \Pr(\hat{Y}^1 \neq \hat{Y}^2)$, it can be estimated from observations $(\hat{Y}_t^1, \hat{Y}_t^2)$, and thus hope for an Hannan consistent policy. In the following we assume that (6) and (9) hold and refer to it as dominance condition. For the case of $K > 2$ sensors, it is given as follows:

**Assumption 1 (Dominance Condition)** *If sensor $i$ predicts correctly, all the sensors in the subsequent stages of the cascade also predict correctly, i.e.,*

$$\hat{Y}_t^i = Y_t \rightarrow \hat{Y}_t^j \quad \forall j > i \geq 1. \tag{10}$$

### 3.1 Robustness of Dominance Condition

For notional convenience rewrite $\gamma_1 - \gamma_2 = p_1 + p_2 + p_5 + p_6 - 2(p_1 + p_6) := p_{12} - 2\delta$, where $p_{12} := \Pr(Y^1 \neq Y^2)$ is the probability that sensors disagree and $\delta := \Pr(Y^2 \neq Y | Y^1 = Y)$ is the conditional probability that sensor 2 is incorrect given that sensor 1 is correct. We can estimate $p_{12}$ from feedback $(\hat{Y}_t^1, \hat{Y}_t^2)$, but $\delta$ cannot be estimated.

**Theorem 1** *For SA-Problem with $K = 2$, a Hannan consistent policy exists if and only if $c \notin [p_{12} - 2\delta, p_{12}]$.*

**Proof:** Under dominance condition $\delta = 0$, thus actions 1 is optimal if $p_{12} < c$, otherwise action 2 is optimal. Suppose dominance condition is violated, i.e., $\delta > 0$, but decisions are made assuming dominance condition holds (i.e., using estimates of $p_{12}$ only), then the optimal action is correctly identified provided $\delta$ is such that $p_{12} - 2\delta < c \Rightarrow p_{12} < c$ or $p_{12} - 2\delta > c \Rightarrow p_{12} > c$. Now, notice that the latter implication is always true. So, whenever action 2 is optimal, violation of dominance condition does not miss the optimal action. However, the first implication holds if and only if $c \notin [p_{12} - 2\delta, p_{12}]$.

In the following we establish that under the dominance condition efficient algorithms for a SAP problem can be derived from algorithms on a suitable stochastic multi-armed bandit problem. We first recall the stochastic multi-armed bandit setting and the relevant results.

4

## 4    Background on Stochastic Multi-armed Bandits with Side Observations

A stochastic multi-armed bandit (MAB), denoted as $\phi := (K, (\nu_k)_{1 \leq k \leq K})$, is a sequential learning problem where number of arms $K$ is known and each arm $i \in [K]$ gives rewards drawn according to an unknown distribution $\nu_k$. Let $X_{i,n}$ denote the random reward from arm $i$ in its $n$th play. For each arm $i \in [K]$, $\{X_{i,t} : t > 0\}$ are independently and identically (i.i.d) distributed and for all $t > 0$, $\{X_{i,t}, i \in [K]\}$ are independent. We note that in the standard MAB setting the learner observes only reward from the selected arm in each round and no information from the other arms is revealed. However, in many applications playing an arm reveals information about the other arms which can be exploited to improve learning performance. Let $\mathcal{N}_i$ denote neighborhood of $i$ such that playing arm $i$ reveals rewards of all arms $j \in \mathcal{N}_i$. Given a set of neighborhood $\{\mathcal{N}_i, i \in [K]\}$, let $\phi_G := (K, (\nu_k)_{1 \leq k \leq K}, G)$ denote a MAB with side-information graph $G = (V, E)$, where $|V| = K$ and $(i, j) \in E$ if $j \in \mathcal{N}_i$. The side-observation graph is known to the learner and remains fixed during the play. To avoid cluttering, we henceforth drop subscript $G$ in $\phi_G$ and it should be clear from context if side-observations exists or not.

Let $\Pi^\phi$ denote a set of polices on $\phi$ that maps the past history into an arm in each round.. If the learner knows $\{\nu_k\}_{k \in [K]}$, then the optimal policy is to play the arm with highest mean. Given a policy $\pi \in \Pi^\phi$, its performance is measured with respect to the optimal policy and is defined in terms of expected cumulative regret (or simply regret) as follows ( only reward from the arm played contribute to the regret and not that from the side-observations): Let $\pi$ selects arm $i_t$ in round $t$. After $T$ rounds, its regret is

$$R_T^\phi(\pi) = T\mu_{i^*} - \sum_{t=1}^{T} \mu_{i_t}, \tag{11}$$

where $\mu_i = \mathbb{E}[X_{i,n}]$ denotes mean of distribution $\nu_i$ for all $i \in [K]$ and $i^* = \arg\max_{i \in [K]} \mu_i$. Let $N_i^\phi(t) = \sum_{s=}^{t} \mathbf{1}\{i_s = i\}$ denote the number of pulls of arm $i$ till time $t$. Then, the regret of policy $\pi$ can be expressed

$$R_T^\phi(\pi) = \sum_{i=1}^{K} (\mu_{i^*} - \mu_i)\mathbb{E}[N_i^\phi(T)].$$

The goal is to learn a policy that minimizes the regret.

Buccapatnam et al. (2014) establish that any policy $\pi \in \Pi^\phi$ where side observation graph is such that $i \in \mathcal{N}_i$ for all $i \in [K]$ satisfies

$$\liminf_{T \to \infty} R_T^\phi(\pi)/\log T \geq \eta(G) \tag{12}$$

where $\eta(G)$ is the optimal value of the following linear optimization

$$\text{LP1}: \quad \min_{\{w_i\}} \sum_{i \in [K]} (\mu_{i^*} - \mu_i)w_i$$

$$\text{subjected to } \sum_{j \in \mathcal{N}_i} w_i \geq 1/D(\mu_i \| \mu_{i^*}) \text{ and } w_i \geq 0 \text{ for all } i \in [K], \tag{13}$$

$D(\mu_i \| \mu_{i^*})$ here denotes the Kullback-Leibler divergence between $\nu_i$ and $\nu_{i^*}$. When $\mathcal{N}_i = \{i\}$ for all $i \in [K]$, it reduces to the classical lower bound of $\sum_{i \neq i^*}(\mu_{i^*} - \mu_i)/D(\mu_i \| \mu_{i^*})$ established in Lai & Robbins (1985). Further, Buccapatnam et al. (2014) also gave an UCB based strategy, named UCB-LP, that explores arms at a rate in proportion to the size of their neighborhood. Specifically, UCB-LP plays arms in proportions to the values $\{z_i^*, i \in [K]\}$ computed from the following linear optimization which is a relaxation of LP1.

$$\text{LP2}: \min_{\{z_i\}} \sum_{i \in [K]} z_i \text{ subjected to } \sum_{j \in \mathcal{N}_i} z_i \geq 1 \text{ and } z_i \geq 0 \text{ for all } i \in [K] \tag{14}$$

The regret of UCB-LP is upper bounded by

$$\mathcal{O}\left(\sum_{i \in [K]} z_i^* \log T\right) + \mathcal{O}(K\delta), \tag{15}$$

5

177 where $\delta = \max_{i \in [K]} |\mathcal{K}_i|$ and $\{z_i^*\}$ are the optimal values of $LP2$.

**Definition 1 (Domination number Buccapatnam et al. (2014))** *Given a graph $G = (V, E)$, a sub-
set $W \subset V$ is a dominant set if for each $v \in V$ there exists $u \in W$ such that $(u, v) \in E$. The size of
the smallest dominant set is called weak domination number and is denoted as $\xi(G)$.*

Since any dominating set is a feasible solution of LP2, we get $\sum_{i \in [K]} z_i^* \leq \xi(G)$, and the regret of
UCB-LP is $\mathcal{O}(\xi(G) \log T)$.

## 5 Regret Equivalence

In this section we establish that under the dominance condition SAP is 'regret equivalent' to an
instance of MAB with side-information and the corresponding algorithm for MAB can be suitably
imported to solve SAP efficiently.

**Definition 2 (Regret Equivalence)** *Consider a SAP problem $\psi := (K, \mathcal{A}, (\gamma_i, c_{i-1})_{i \in [K]})$ and a
bandit problem with $\phi_G := (N, (\nu_i)_{i \in [N]}, G)$ side-information graph $G$. We say that $\psi$ is regret-
equivalent to $\phi_G$ if given a policy $\pi$ for problem $\psi$, one can come up with a policy $\pi'$ that uses $\pi$,
such that the regret of $\pi'$ on any instance of $\phi_G$ is the same as the regret of $\pi$ on some corresponding
instance of $\psi$, and vice versa.*

In the following we first consider the SAP with 2 sensors and then the general case with more than 2
sensors. The 2 sensors case helps to draw comparison with the well studied apple tasting problem
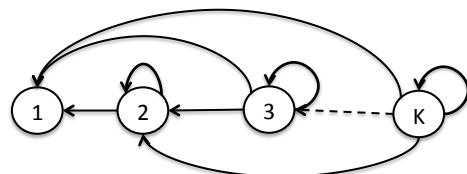and understand role of the dominance condition.

### 5.1 SAP with two sensors

In the SAP with only two actions, the feedback from action $i = 1$ reveals no information about
the loss incurred in that round. However feedback after action $i = 2$ reveals (partial) information
about the loss of both actions. Suppose feedback is such that predictions of the sensors disagree,
i.e., $\hat{Y}_t^1 \neq \hat{Y}_t^2$ after action 2. The dominance condition then implies that the only possible events are
$\hat{Y}_t^1 \neq Y_t$ and $\hat{Y}_t^2 = Y_t$. I.e., the true label is that predicted by sensor-2, hence loss incurred is just $c$
(prediction loss is zero). Suppose predictions of the sensors agree, i.e., $\hat{Y}_t^1 = \hat{Y}_t^2$, then the dominance
condition implies that either predictions of both are correct or both are incorrect. Though the true
loss is not known in this case, the learner can infer some useful knowledge: in round $t$, if prediction
of both the sensors agree, then the difference in losses of the actions is $L_t(2) - L_t(1) = c > 0$.
And if predictions of the sensors disagree, then dominance assumption implies that $L_t(1) = 1$ and
$L_t(2) = c$ or $L_t(2) - L_t(1) = c - 1 < 0$. Thus, each time learner plays action 2, he gets to know
whether or not he was better off by selecting the other action. This setup sounds similar to the standard
apple tasting problem Helmboat et al. (2000) ], but differs in terms of the information structure when
action 2 is played.

**Apple tasting problem:** In the apple tasting problem, a learner gets a sequence of apples and some
of them can be rotten. In each round, the learner can either accept or reject an apple. If an apple is
accepted, the learner tastes it and incurs a penalty if it is rotten. If apple is rejectsed, he still incurs
the penalty if it is rotten, but do not get to observe its quality. The goal of the learner is to taste more
good apples. The SAP setting is a more general version than the apple tasting problem–in any round,
actions 1 reveals no loss values. Action 2 reveals only partial information about the losses, but not
the exact losses as in the apple tasting problem. However, we next show that the partial information
is enough to achieve optimal performance.

### 5.2 SAP with more than two actions

In the SAP with two sensors, only action 2 provides information about the losses. In the case
with $K > 2$ sensors, by playing an action $k$, we can obtain information about the losses of
all sensors $l < k$ by recursively applying the dominance condition between pair of sensors.
Further, any information provided by action $k > 2$ is
contained in that provided by all actions $k' \geq k-$ if

6

action $k$ is played in round $t$, then we observe predictions $\{\hat{Y}_t^1, \hat{Y}_t^2, \cdots, \hat{Y}_t^i\}$ which includes the observed predictions of all actions $k' \leq i$. This side-observation can be represented by a directed graph $G^S = (V, E)$, where $|V| = K$ and $E = \{(i, j) : i1 < i \leq j \leq K\}$. Note that $G^S$ has self loops for all nodes except for node 1. The nodes in $G^S$ represents actions of the SAP and an edge $(i, j) \in E$ implies that actions $i$ provides information about action $j$. The side-observation graph for the SAP is shown in Figure (2).

**Theorem 2** *Let the dominance condition (10) holds. Then SAP $\psi$ with $K \geq 2$ is regret equivalent to a MAB with side-observation graph $G^S$.*

**Proposition 1 (SAP regret lower bound)** *Let $\pi$ be any policy on SAT with 2 sensors such that it pulls the suboptimal arm only sub polynomial many times, i.e., $\mathbb{E}[N_i^\psi(T)] = o(T^a)$ for all $a > 0$ and $i \neq i^*$. Then,*

$$\liminf_{T \to \infty} R_T^\psi(\pi)/\log T \geq \kappa \ where \tag{16}$$

$$\kappa = \min_{\{w_i\}} \sum_{i \in [K]} (\mu_{i^*} - \mu_i) w_i$$

$$subjected\ to \sum_{ji} w_i \geq 1/D(\mu_i + \sum_{j<i} c_j || \mu_{i^*})\ for\ all\ i \in [K] \tag{17}$$

$$w_i \geq 0\ for\ all\ i \in [K]$$

**Proposition 2 (K-SAT regret upper bound)** *Let $\pi'$ denote a policy on a $K$-armed stochastic bandit where mean of arm $i > 1$ is $\gamma_1 - \gamma_i - \sum_{j<i} c_j$ and arm 1 has a fixed reward of value zero, and the side-observation graph is $G^S$. Then, the regret of a policy $g_1(\pi)$ for the SAT problem obtained from mapping (28) is upper bounded as*

$$R_T^\psi(g(\pi)) \leq \mathcal{O}(\xi(G^S) \log T + K^2) \tag{18}$$

*when $\pi' =$ UCB-LP Buccapatnam et al. (2014).*

## References

Abbasi-Yadkori, Yasin, Pál, Dávid, and Szepesvári, Csaba. Improved algorithms for linear stochastic bandits. In *Proceeding of Advances in Neural Information Processing Systems (NIPS)*, pp. 2312–2320, 2011.

Buccapatnam, S., Eryilmaz, A., and Shroff, N. B. Stochastic bandits with side observation on networks. In *Proceeding of Sigmetrics*, 2014.

Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. In *Proceeding of Conference on Learning Theory, COLT*, Helsinki, Finland, July 2008.

Hannan, J. Approximation to bayes risk in repeated plays. *Contributions to the Theory of Games*, 3: 97–139, 1957.

Helmboat, D. P., Littlestone, PN, and Long, P.M. Apple tasting. *Journal of Information and Computation*, 161(2):85–139, 2000.

Lai, Tze Leung and Robbins, Herbert. Asymptotically efficient adaptive allocation rules. *Journal of Advances in applied mathematics*, 6(1):4–22, 1985.

Li, L., Wei, C., Langford, J., and Schapire, R. E. A contextual-bandit approach to personalized news article recommendation. In *Proceeding of International Word Wide Web conference, WWW*, NC, USA, April 2010.

Mackay, David. J. C. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.

Rusmevichientong, Paat and Tsitsiklis, John N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.

Seldin, Y., Bartlett, P., Crammer, K., and Abbasi-Yadkori, Y. Prediction with limited advice and multiarmed bandits with paid observations. In *Proceeding of International Conference on Machine Learning, ICML*, pp. 208–287, 2014.

Trapeznikov, K. and Saligrama, V. Supervised sequential classification under budget constraints. In *Proceeding of International Conference on Artificial Intelligence and Statistics, AISTATS*, pp. 235–242, 2013.

Trapeznikov, K., Saligrama, V., and Castanon, D. A. Mulit-stage classifier design. *Machine Learning Journal*, 39:1–24, 2014.

Xu, Z., Kusner, M., Chen, M., and Weinberger, K. Q. Cost-sensitive tree of classifiers. In *Proceeding of International Conference on Machine Learning, ICML*, pp. 133–141, 2013.

## A Proof of Theorem ??

Consider a 1-armed stochastic bandit problem where arm with constant reward has value $c$ and the arm that gives stochastic reward has mean value $\gamma_1 - \gamma_2$. Given an arbitrary policy $\pi = (\pi_1, \pi_2, \cdots \pi_t)$ for the SAP, we obtain a policy for the bandit problem from $\pi$ as follows: Let $H_{t-1}$ denote the history, consisting of all arms played and the corresponding rewards, available to policy $\pi_{t-1}$ till time $t-2$. Let $a_{t-1}$ denote the action selected by the bandit policy in round $t-1$ and $r_{t-1}$ the observed reward. Then, the next action $a_t$ is obtained as follows:

$$a_t = \begin{cases} \pi_t(H_{t-1} \cup \{1, \emptyset\}) \text{ if } a_{t-1} = \text{fixed rewad arm} \\ \pi_t(H_{t-1} \cup \{2, r_{t-1}\}) \text{ if } a_{t-1} = \text{stochastic arm} \end{cases} \tag{19}$$

Conversely, let $\pi' = \{\pi'_1, \pi'_2, \cdots\}$ denote an arbitrary policy for the 1-armed bandit problem. we obtain a policy for the SAP as follows: Let $H'_{t-1}$ denote the history, consisting of all actions played and feedback, available to policy $\pi'_{t-1}$ till time $t-1$. Let $a'_{t-1}$ denote the action selected by the SAP policy in round $t-1$ and observed feedback $F_t$. Then, the next action $a'_t$ is obtained as follows:

$$a'_t = \begin{cases} \pi'_t(H'_{t-1} \cup \{1, c\}) \text{ if } a'_{t-1} = \text{action 1} \\ \pi'_t(H'_{t-1} \cup \{2, \mathbf{1}\{\hat{Y}^1_t \neq \hat{Y}^2_t\}\}) \text{ if } a_{t-1} = \text{actions 2.} \end{cases} \tag{20}$$

We next show that regret of $\pi$ on the SAP is same as that of derived policy on the 1-armed bandit, and regret of $\pi'$ on the 1-armed bandit is same as regret of the derived policy on SAP. We first argue that any policy on the SAP problem with 2 actions needs the information if whether the predictions of sensors match or not whenever action 2 is played. The following observation is straightforward.

**Lemma 1** *Let dominance condition holds. Then,* $\Pr\{\hat{Y}^1_t \neq \hat{Y}^2_t\} = \gamma_1 - \gamma_2$.

$$\Pr\{\hat{Y}^1_t \neq \hat{Y}^1_t\} = \Pr\{\hat{Y}^1_t = Y_t, \hat{Y}^2_t \neq Y_t\} + \Pr\{\hat{Y}^2_t = Y_t, \hat{Y}^1_t \neq Y_t\} \tag{21}$$

$$= \Pr\{\hat{Y}^2_t = Y_t, \hat{Y}^1_t \neq Y_t\} \text{ from assumption (10)} \tag{22}$$

$$= \Pr\{\hat{Y}^1_t \neq y_t\} \Pr\{\hat{Y}^2_t = Y_t | \hat{Y}^1_t \neq Y_t\} \tag{23}$$

$$= \Pr\{\hat{Y}^1_t \neq Y_t\} \left(1 - \Pr\{\hat{Y}^2_t \neq Y_t | \hat{Y}^1_t \neq Y_t\}\right) \tag{24}$$

$$= \Pr\{\hat{Y}^1_t \neq Y_t\} \left(1 - \frac{\Pr\{\hat{Y}^2_t \neq Y_t, \hat{Y}^1_t \neq Y_t\}}{\Pr\{\hat{Y}^1_t \neq Y_t\}}\right) \tag{25}$$

$$= \Pr\{\hat{Y}^1_t \neq Y_t\} - \Pr\{\hat{Y}^2_t \neq Y_t\} \text{ by contrapositve of (10)} \tag{26}$$

From Lemma 1, mean of the observations $Z_t := \mathbf{1}\{\hat{Y}^1_t \neq \hat{Y}^2_t\}$ from action 2 in the SAP is a sufficient statistics to identify the optimal arm. Thus, any SAP only needs to know $Z_t$ in each round, and $Z_t$ are i.i.d with mean $\gamma_1 - \gamma_2$. Our mapping of policies is such that any poilcy for SAP (1-armed bandits) and the derived policy on the 1-armed bandit (SAP) play the sub-optimal arm same number of times. For the sake of simplicity assume that action 1 is optimal for SAP ($\gamma_1 > \gamma_2 + c$) and let a policy $\pi$ on SAP plays it $N_1(T)$ number if times. Then, we have

$$R^\psi_T(\pi) = \Delta_i \mathbb{E}[N^\psi_1(T)] = (\gamma_1 - \gamma_2 - c)\mathbb{E}[N_1(T)]$$

Let $f(\pi)$ denote the policy for the 1-armed bandit obtained using the mapping (19). Now, for the 1-armed bandit, where the arm with stochastic rewards is optimal, we have

$$R^\phi_T(f(\pi)) = (\mu_2 - \mu_1)\mathbb{E}[N_1(T)] = (\gamma_1 - \gamma_2 - c)\mathbb{E}[N^\phi_1(T)]$$

Thus the regret of $\pi$ on the SAP problem and that of $f(\pi)$ on the 1-armed bandit are the same. We can argue similarly for the other case.

## B Proof of Theorem 2

Consider a $K$-armed stochastic bandit problem where rewards distribution $\nu_i$ has mean $\gamma_1 - \gamma_i - \sum_{j<i} c_j$ for all $i > 1$ and arm 1 gives a fixed reward of value 0. The arms have side-observation

9

structure defined by graph $G^S$. Given an arbitrary policy $\pi = (\pi_1, \pi_2, \cdots \pi_t)$ for the SAP, we obtain a policy for the bandit problem with side observation graph $G^S$ from $\pi$ as follows: Let $H_{t-1}$ denote the history, consisting of all arms played and the corresponding rewards, available to policy $\pi_{t-1}$ till time $t - 2$. In round $t - 1$, let $a_{t-1}$ denote the arm selected by the bandit policy, $r_{t-1}$ the corresponding reward and $o_{t-1}$ the side-observation defined by graph $G_S$ excluding that from the first arm. Then, the next action $a_t$ is obtained as follows:

$$a_t = \begin{cases} \pi_t(H_{t-1} \cup \{1, \emptyset\}) \text{ if } a_{t-1} = \text{arm } 1 \\ \pi_t(H_{t-1} \cup \{i, r_{t-1} \cup o_{t-1}\}) \text{ if } a_{t-1} = \text{arm i} \end{cases} \tag{27}$$

Conversely, let $\pi' = \{\pi'_1, \pi'_2, \cdots\}$ denote an arbitrary policy for the $K$-armed bandit problem with side-observation graph. we obtain a policy the SAP as follows: Let $H'_{t-1}$ denote the history, consisting of all actions played and feedback, available to policy $\pi'_{t-1}$ till time $t - 2$. Let $a'_{t-1}$ denote the action selected by the SAP policy in round $t - 1$ and observed feedback $F_t$. Then, the next action $a'_t$ is obtained as follows:

$$a'_t = \begin{cases} \pi'_t(H'_{t-1} \cup \{1, 0\}) \text{ if } a'_{t-1} = \text{action } 1 \\ \pi'_t(H'_{t-1} \cup \{i, \mathbf{1}\{\hat{Y}^1_t \neq \hat{Y}^2_t\} \cdots \mathbf{1}\{\hat{Y}^1_t \neq \hat{Y}^i_t\}\}) \text{ if } a_{t-1} = \text{action i.} \end{cases} \tag{28}$$

We next show that regret of a policy $\pi$ on the SAP problem is same as that of the policy derived from it for the $K$-armed bandit problem with side information graph $G^S$, and regret of $\pi'$ on the $K$-armed bandit with side information graph $G^S$ is same as that of the policy derived from it for the SAP.

Given a policy $\pi$ for the SAP problem let $f_1(\pi)$ denote the policy obtained by the mapping defined in (27). The regret of policy $\pi$ that plays actions $i$, $N_i(T)$ times is given by

$$R^\psi_T(\pi) = \sum_{i=1}^K \left[ \left( \gamma_i + \sum_{j<i} c_j \right) - \left( \gamma_{i^*} + \sum_{j<i^*} c_j \right) \right] \mathbb{E}[N^\psi_i(T)] \tag{29}$$

$$\tag{30}$$

Now, regret of regret policy $f_1(\pi)$ on the $K$-armed bandit problem with side information graph $G^S$

$$R^{\phi_G}_T(f_1(\pi)) = \sum_{i=1}^K \left[ \left( \gamma_1 - \gamma_{i^*} - \sum_{j<i^*} c_j \right) - \left( \gamma_1 - \gamma_i - \sum_{j<i} c_j \right) \right] \mathbb{E}[N^{\phi_G}_i(T)] \tag{31}$$

which is same as $R^\phi_T(\pi)$. This concludes the proofs.

# C Extension to context based prediction

In this section we consider that the prediction errors depend on the context $X_t$, and in each round the learner can decide which action to apply based on $X_t$. Let $\gamma_i(X_t) = \Pr\{\hat{Y}^1_t \neq \hat{Y}^2_t | X_t\}$ for all $i \in [K]$. We refer to this setting as Contextual Sensor Acquisition Problem (CSAP) and denote it as $\psi_c = (K, \mathcal{A}, \mathcal{C}, (\gamma_i, c_i)_{i \in [K]})$.

Given $x \in \mathcal{C}$, let $L_t(a|x)$ denote the loss from action $a \in \mathcal{A}$ in round $t$. A policy on $\phi^c$ maps past history and current contextual information to an action. Let $\Pi^{\psi_c}$ denote set of policies on $\psi_c$ and for any policy $\pi \in \Pi^{\psi_c}$, let $\pi(x_t)$ denote the action selected when the context is $x_t$. For any sequence $\{x_t, y_t\}_{t>0}$, the regret of a policy $\pi$ is defined as:

$$R^{\phi_c}_T(\pi) = \sum_{t=1}^T \mathbb{E}\left[L_t(\pi(x_t)|x_t)\right] - \sum_{t=1}^T \min_{a \in \mathcal{A}} \mathbb{E}\left[L_t(a|x_t)\right]. \tag{32}$$

As earlier, the goal is to learn a policy that minimizes the expected regret, i.e., $\pi^* = \arg\min_{\pi \in \Pi^{\psi_c}} \mathbb{E}[R^{\psi_c}_T(\pi)]$.

In this section we focus on CSA-problem with two sensors and assume that sensor predictions errors are linear in the context. Specifically, we assume that there exists $\theta_1, \theta_2 \in \mathcal{R}^d$ such that $\gamma_1(x) = x'\theta_1$ and $\gamma_2(x) + c = x'\theta_2$ for all $x \in \mathcal{C}$, were $x'$ denotes the transpose of $x$. By default all vectors are column vectors. In the following we establish that CSAP is regret equivalent to a stochastic liner bandits with varying decision sets. We first recall the stochastic linear bandit setup and relevant results.

## C.1 Background on Stochastic Linear Bandits

In round $t$, the learner is given a decision set $D_t \subset \mathcal{R}^d$ from which he has to choose an action. For a choice $x_t \in D_t$, the learner receives a reward $r_t = x_t'\theta^* + \epsilon_t$, where $\theta^* \in \mathcal{R}^d$ is unknown and $\epsilon_t$ is random noise of zero mean. The learner's goal is to maximize the expected accumulated reward $\mathbb{E}\left[\sum_{t=1}^{T} r_t\right]$ over a period $T$. If the leaner knows $\theta^*$, his optimal strategy is to select $x_t^* = \arg\max_{x \in D_t} x'\theta^*$ in round $t$. The performance of any policy $\pi$ that selects action $x_t$ at time $t$ is measured with respect to the optimal policy and is given by the expected regret as follows

$$R_T^L(\pi) = \sum (x_t^*)'\theta^* - \sum x_t'\theta^*. \tag{33}$$

The above setting, where actions sets can change in every round, is introduced inAbbasi-Yadkori et al. (2011) and is a more general setting than that studied in Dani et al. (2008); Rusmevichientong & Tsitsiklis (2010) where decision set is fixed. Further, the above setting also specializes the contextual bandit studied in Li et al. (2010). The authors in Abbasi-Yadkori et al. (2011) developed an 'optimism in the face of uncertainty linear bandit algorithm' (OFUL) that achieves $\mathcal{O}(d\sqrt{T})$ regret with high probability when the random noise is $R$-sub-Gaussian for some finite $R$. The performance of OFUL is significantly better than $ConfidenceBall_2$ Dani et al. (2008), $UncertainityEllipsoid$ Rusmevichientong & Tsitsiklis (2010) and $LinUCB$ Li et al. (2010).

**Theorem 3** *Consider a CSA-problem with $K = 2$ sensors. Let $\mathcal{C}$ be a bounded set and $\gamma_i(x) + c_i = x'\theta_i$ for $i = 1, 2$ for all $x \in \mathcal{C}$. Assume $x'\theta_1, x'\theta_2 \in [0\ 1]$ for all $x \in \mathcal{C}$. Then, equivalent to a stochastic linear bandit.*

## C.2 Proof of Theorem 3

Let $\{x_t, y_t\}_{t \geq 0}$ be an arbitrary sequence of context-label pairs. Consider a stochastic linear bandit where $D_t = \{0, x_t\}$ is a decision set in round $t$. From the previous section, we know that given a context $x$, action 1 is optimal if $\gamma_1(x) - \gamma_2(x) - c < 0$, otherwise action 2 is optimal. Let $\theta := \theta_1 - \theta_2$, then it boils down to check if $x'\theta - c < 0$ for each context $x \in \mathcal{C}$.

For all $t$, define $\epsilon_t = \mathbf{1}\{\hat{Y}_t^1 \neq \hat{Y}_t^2\} - x_t'\theta$. Note that $\epsilon_t \in [0\ 1]$ for all $t$, and since sensors do not have memory, they are conditionally independent given past contexts. Thus, $\{\epsilon_t\}_{t>0}$ are conditionally $R$-sub-Gaussian for some finite $R$.

Given a policy $\pi$ on a linear bandit we obtain next to play for the CSAP as follows: For each round $t$ define $a_t \in \mathcal{C}$ and $r_t \in \{0, 1\}$ such that $a_t = 0$ and $r_t = 0$ if action 1 is played in that round, otherwise set $a_t = x_t$ and $r_t = \mathbf{1}\{\hat{y}_t^1 \neq \hat{y}_t^1\}$. Let $\mathcal{H}_t = \{(a_1, r_1) \cdots (a_{t-1}, r_{t-1})\}$ denote the past actions and corresponding rewards observed till time $t - 1$. In round $t$, after observing context $x_t$, we transfer $((a_{t-1}, r_{t-1}), D_t)$, where $D_t = \{0, x_t\}$. If $\pi$ outputs $0 \in D_t$ as the optimal choice, we play action 1, otherwise we play action 2.

Conversely, suppose $\pi'$ denote a policy for the CSAP problem we select action to play from decision set $D_t = \{0, x_t\}$ as follows. For each round $t$ define $a_t' \in 1, 2$ and $r_t' \in \mathcal{R}$ such that $a_t' = 1$ and $r_t' = \emptyset$ if 0 is played otherwise set $a_t' = 2$ and $r_t' = x_t'\theta^* + \epsilon_t$ if $x_t$ is played. Let $\mathcal{H}_t' = \{(a_1', r_1') \cdots (a_{t-1}', r_{t-1}')\}$ denote the past actions and corresponding rewards observed till time $t - 1$. In round $t$, after observing set $D_t$, we transfer $((a_{t-1}', r_{t-1}'), x_t)$ to policy $\pi'$. If $\pi$ outputs action 1 as the optimal choice, we play action 0, otherwise we play $x_t$.