# Pre-Analysis Plan:
# Increasing the uptake of long-acting reversible contraceptives among adolescents and young women in Cameroon

Susan Athey, Vitor Hadad, Julian Jamison, Berk Özler, and Luca Parisotto

March 23, 2021

**Abstract**

This pre-analysis plan describes an adaptive experiment proposing to test an integrated behavioral science approach to increase the uptake of modern contraceptives, in particular long-acting reversible contraceptives (LARCs) among reproductive-age females in Cameroon, including adolescents who may be unmarried and/or nulliparous. It aims to do so with the aid of a tablet-based "app," or a job-support tool, designed for use by health providers counseling women on family planning. The adaptive experiment aims to find the optimal combination of subsidies for contraceptive methods and counseling approaches for different sub-groups among the clients of a women and children's hospital in Yaoundé, Cameroon to increase the uptake of LARCs while also improving the quality of care.

## 1 Introduction

Nearly half of all pregnancies worldwide are considered unwanted or mistimed (Bearak et al., 2018). Unintended pregnancies are associated with a range of poor outcomes for women of childbearing age. Foremost among these in low-income countries (LIC) is maternal mortality. In sub-Saharan Africa, the maternal mortality ratio was 542 in 2017 (MMR, or the number of maternal deaths per 100,000 live births (World Health Organization, 2019). In Cameroon, which had an MMR of 529 in 2017 and a total fertility rate of about 4.6 in 2018, a fifth of all births were unwanted or considered mistimed by the mother. For adolescents, this figure was one half (DHS, 2020). Unintended pregnancies not only cause immediate welfare losses due to miscarriages and abortions, but they can also lead to school dropouts, early marriages, and reduce women's empowerment (see Baird et al. 2019, 2011, and citations within). In addition to the large welfare loss for women, unintended pregnancies also cause poor outcomes for their children: they are strongly predictive of short inter-pregnancy intervals, which are in turn positively associated with babies being born prematurely, at low birth weight, or small for their gestational age (DeFranco et al., 2015; Frost et al., 2014).

There exist effective modern contraceptive methods to prevent unintended pregnancies. However, about a quarter of women in low- and middle-income countries (LMIC) and half of adolescents, who report wanting to avoid a pregnancy, also report not using a contraceptive method (see Sully et al. 2020). Furthermore, despite the efficacy of modern contraception, many unintended pregnancies occur in the same month when contraception is being used (see Finer and Henshaw 2006 for evidence from the U.S.). Therefore, both increasing the uptake of modern contraceptives and optimizing

1

the individual choice of contraceptive method among users is important to reduce unintended pregnancies. Contraceptive methods differ in their typical use effectiveness and only a small portion of contraceptive users adopt methods with the lowest failure rates. While approximately 18 women out of 100 using the male condom will get pregnant within the next 12 months, this number declines to six to nine among those using short-acting reversible contraceptives (SARCs, i.e. the pill and injectable), and less than one among the women using long-acting reversible contraceptives (LARCs, i.e. the intra-uterine device, IUD, or implant) (Trussell, 2011). In the U.S., between 2017-19, 26% of women aged 15-49 were using a SARC (including male condoms) but only 10% were using a LARC (Daniels and Abma, 2020). The problem is worse in LMICs, where only a minority of women of childbearing age use any modern method of contraception, most of whom rely on male condoms. In Cameroon, the setting for our study, less than 5% of married women were using the oral contraceptive pill or the injectable (Depo-Provera), while 3.5% were using LARCs in 2018 (DHS, 2020).

**Interventions**   The tablet-based decision-support tool, or *the app*, was designed for use by service providers conducting family planning (FP) counseling sessions. It effectively serves to structure the counseling session, thus ensuring that clients consistently receive high quality consultations. The structure of the counseling session as guided by the job-support tool is not fundamentally different than standard best-practice (see Appendix section A.1 for a more detailed description). During a structured discussion, *the app* records the clients' goals, fertility plans, needs, and preferences regarding contraceptive methods, as well as their medical and birth history, blood pressure, weight, and other relevant risk-factors. The elicitation of these preferences is to help the client, with guidance from the provider, to find a suitable a contraceptive method that fits her context.[1]

Counseling clients adequately on contraceptive method choice and use is critical to ensuring that individuals' needs are respected and met (Holt et al., 2017). While there is justifiably a focus on increasing the uptake of long-acting contraceptive methods, concerns have been raised over potentially negative consequences of this focus on respect for the clients' autonomy and decision-making. Therefore, in the established approach to counseling, the client is given information about all contraceptive methods and asked to choose which method she would prefer to discuss. Generally, the provider is expected to not provide guidance or advice during this informed decision-making process by the client alone. The main innovation that *the app* brings to counseling sessions is a small but important paradigm shift with respect to shared decision-making. *The app* takes all the answers by the client into consideration and uses an internal algorithm to rank methods from most suitable to least suitable for the client's specific context. The provider then recommends the client that they discuss methods in this order of suitability until they find a method that the client would like to adopt (more details on the ranking algorithm in Appendix section A.4). During the experiment, we will test this approach by randomly allocating clients to the *status-quo* – which we will refer to as *Side-by-Side (SBS)* – or the *ranked recommendation* counseling style – referred to as *sequential (SEQ)* (see Section 2 for more details on treatment arms and Appendix section A.2 for a more detailed description of the consultation style under the two regimes).

Cost is also an important barrier to the adoption of contraceptives, perhaps especially for adolescents

---

[1]That is not contraindicated as per the assessment of the client's medical eligibility (more detail on the medical eligibility criteria in Appendix section A.3).

who may not independently have the means to adopt methods that are not free or very cheap.[2] Clients who receive a FP consultation by a provider in one of the hospital 's participating services will thus be offered randomly varying discounts, or randomly varying prices, for LARCs — i.e. the (copper) IUD and Implant.[3] Clients will be randomized into one of four price categories for LARCs {High: CFA 4,000/$7.30 , Mid: CFA 2,000/$3.60 , Low: CFA 1,000/$1.80 , and Free}.[4] It is important to note here that even the High LARC price constitutes at least a 20% discount relative to the normal prices charged at the hospital prior to the study period, where a copper IUD would cost CFA 5,000/$9.1 and the Implant CFA 5,250/$9.5. During the pilot phase of this adaptive experiment, we also offered randomly varying price discounts for the pill and the injectable. As these discounts did not change take-up rates of SARCs or LARCs, we forgo such discounts for SARCs in the main study. SARCs are thus offered at their regular pre-experiment prices of CFA 1,250/$2.30 for the injectable (to be renewed in three months) and CFA 1,500/$2.70 for three monthly cycles of the pill. The price discounts are revealed towards the end of the consultation, after clients have made their initial method choice - unless they inquire about prices earlier on in the counseling session.[5] The prices for all methods are revealed at once: this means that the clients can always reconsider their decision to adopt a method after the full set of prices are revealed by *the app*.

**Sample**   The study takes place at the Hôpital Gynéco-Obstétrique et Pédiatrique de Yaoundé (HGOPY), located in Yaoundé, the capital city of Cameroon. All family planning consultations taking place within the confines of the hospital were conducted by a trained provider using *the app*. Our sample is drawn from the universe of women who receive a family planning consultation. Clients who receive a family planning consultation with the app are a diverse set of women with different needs with respect to family planning. During the pilot phase, approximately half of the clients consulted presented themselves directly at the Family Planning unit seeking family planning services, including: simply to receive information, adopt a new method, renew their current method, switch to another method, manage side effects of their current method, or discontinue their current method. The remaining clients who were counseled mainly presented at the maternity and gynecology wards: some had just given birth; some were pregnant and receiving ante-natal services; some might have returned to the hospital post-partum for a check-up or for their babies to receive vaccinations; yet others might have presented with a gynecological problem. This group of clients are especially important to target due to the importance of post-partum adoption of LARCs to improve birth spacing.

The sample of women who visit the hospital is equally diverse in characteristics. Nevertheless, on

---

[2]While large changes in prices of contraceptive methods were found to have little impact on contraceptive use in Indonesia during the 1998 financial crisis (McKelvey et al., 2012), more recent studies indicate that women, unmarried young women in particular, may be more responsive to contraceptive prices (Rau et al., 2017; Lindo and Packham, 2017). In the U.S., adolescents who receive comprehensive counseling and face no cost barriers preferentially select LARC methods and continue to use them long-term (Mestad et al., 2011). In Colorado, U.S., funding provided to Title X clinics on the condition that they stock and provide free LARCs to low-income women not only caused many clinics, which had never offered these methods before, to start providing them for free, but also decreased teen birth rates substantially, with the effect being largest in the poorest counties of the state (Lindo and Packham, 2017).

[3]The LNG IUD was not available at the study hospital at the time of the study.

[4]At the proximate exchange rate of CFA 550 per USD

[5]This decision was made after careful deliberation with the healthcare providers at the hospital and based on established best-practices and ethical norms.

3

average, the client base at HGOPY is wealthier and better educated than the average Cameroonian woman. For example, comparing clients enrolled in the study during the pilot period to estimates obtained from the 2018 round of the DHS survey, about 40% of women enrolled in the study have tertiary education compared to 13% of women from Yaounde or 4% in all of Cameroon. However, the clients enrolled in the study are similar to the average Cameroonian woman in terms of contraceptive use. Prior to receiving the consultation, about 3.5% were using a LARC and 3.9% were using a SARC in the enrolled sample, compared to 5.4% and 3.5% using a LARC and 5% and 4.8% using a SARC in Yaoundé and the whole of Cameroon, respectively.

**Data**    There are two primary data sources for the study. The first consists of data collected by the app during the consultation. The app functions like a data collection tool in that the clients' answers during the consultation are recorded and uploaded to a central server. The application was built to be integrated with the hospital's own administrative and health record systems and thus collects a rich set of baseline characteristics, including demographics characteristics, relevant medical history, and birth history. Then, as clients are taken through the consultation the app records the client's answers to the questions asked by the provider. This results in a rich database that includes the client's fertility preferences, their experience with contraception (i.e. whether they are currently using a method and if so what their experience with side-effects are, etc.), whether and why they seek to adopt a specific method, their relevant medical history to check for contraindications, and finally which method - if any - they adopted at the end of the consultation. For each method considered the app also records the reason why a client did not want to adopt it, and if the client chose a method but did not adopt it immediately, the reasons for this decision.

This data is stored on a HIPAA compliant secure server managed by the hospital, which can be accessed only by authorised hospital staff.[6] The authorized hospital staff shares anonymized client data with the study team for the analysis, as per the IRB approval granting authorization for the study team to review de-identified electronic health records from the family planning consultation. This allows us to construct and analyze primary outcomes during the study period without any consent bias. Increasingly, a case is being made by researchers and public health officials that access to de-identified, minimally risky electronics health records data should be made available to researchers without the need to seek consent from each patient. The primacy of respect for patient autonomy is uncontroversial: patients, of course, have the right to decide whether and which procedures they wish to undergo, what methods they wish to adopt, and whether they wish to participate in a research protect. However, as is the case in this study, a strong argument can be made for the use of information that is already stored by the provider when (a) that information, properly anonymized, can have significant public health and/or biomedical value; and (b) when the only risk to the client is a breach of privacy; and (c) that risk is minimal (e.g. see Porsdam Mann et al. 2016).[7]

---

[6] The dedicated server is hosted by Amazon Web Services (AWS) and is compliant with various assurance and certification programs, including the Health Insurance Portability and Accountability Act (HIPAA) of the USA, the EU Data Protection Directive, Privacy Acts of Australia and New Zealand, among others. Therefore, it is compliant with the highest standards of security to protect the privacy and confidentiality of clients. The server will only be accessed by personnel authorized by the HGOPY administration through a username and password.

[7] We argue that the likelihood of a privacy breach is minimal in our context for two reasons. First, the sensitive data are already kept in two locations: in physical family planning registers and in the hospital's own administrative records. The additional probability of a breach of privacy because of this study is very small. In fact, the study may

The second source of data consists of three waves of follow-up phone surveys. The follow-up interviews, each of which are expected to take 15 minutes or less to conduct by phone, are planned at two weeks, 16 weeks, and 52 weeks after the initial visit – when the client is counseled for the first time by a provider at HGOPY using the app. The short follow-up interviews will include questions about client satisfaction with the counseling session; side effects of the chosen methods; renewal, switching, or discontinuation of the chosen methods; and 12-month pregnancy status. Before starting the counseling sessions, the provider will invite the client to participate in the follow-up study by reading them a notice of information, go over the informed consent form, and obtain their signature if they agree to participate in the study.[8] Proper care is taken to ensure the clients clearly understand that their willingness to participate in the study has no bearing on the quality or the range of services they will receive at HGOPY, nor will it affect their likelihood of receiving discounts for their chosen contraceptive method.

**Ethics** The study protocols were approved by Cameroon's national ethics committee for human subjects research, the Comite National d'Ethique de la Recherche pour la Sante Humaine (CNERSH) – decision No. 2019/08/1183/CE/CNERSH/SP – and received administrative authorization from the Ministry of Health's (MinSante) Division of Health Operations Research (DROS) – decision No. D30-760/L/MINSANTE/SG/DROS. The protocols were also approved by the implementing hospital's own IRB – decision No. 780/CIERSH/DM/2018.[9]. The protocols cover the full set of study procedures and methodology including, but not limited to: data management and information security, enrollment criteria, consent procedures, and treatment of adverse reactions.

The project has benefited from close collaboration with Cameroonian health practitioners and researchers at all stages of development. The study objectives and the decision-support tool were developed by a multi-disciplinary working group formed in Cameroon and comprised of nurses, doctors, and researchers, public health and adolescent health specialists from the Department of Family Health in the Ministry of Health (DSF/MinSanté), and public health and economics researchers focusing on adolescent health from the World Bank. The tool has been extensively tested by Ob-Gyns, medical doctors, nurses, and other health providers conducting family planning counseling, and their feedback has been incorporated to improve the framework adopted by the tool. The study was then developed in cooperation with the implementing hospital's chief of research:

reduce this risk by eliminating the need for double entry of sensitive PII data; requiring dedicated and more secure servers for these data; and encouraging better data handling practices by hospital staff. Access to the job-support tool on each tablet requires a pre-assigned user ID and password. Each provider is asked to transmit data to the dedicated and secure server daily. Transmitted data are automatically deleted from the tablet once transmission is successful. Second, the dedicated and encrypted server is only accessible by personnel authorized by the HGOPY administration through a username and password and off-limits to the rest of the study team.

[8] The target population of our study is females aged 15-49, who present at HGOPY seeking family planning counseling. This means that a small percentage of the clients seeking family planning services will be under the age of 18. While some of these clients will be married and, hence, emancipated, others will be minors. For a 'minor adolescent,' we will seek her assent to participate in our study following the same protocol described above. In addition, we will explain to her that since she is legally a minor, we would like to seek her parents' or guardians' permission for her participation in our study. If the minor adolescents agrees to participate and the parents/guardians are at HGOPY, then their consent will be sought in person. If the parents are not present at the hospital, then the provider will ask the minor adolescent to provide a contact number for the parents/guardians to obtain their verbal consent by phone.

[9] The study protocols, drawn from the documentation submitted for ethics review, are available at: https://www.socialscienceregistry.org/trials/3514

Prof. Dohbit Sama, who is a co-principal investigator.

**Pilot**    The pilot phase of this study was conducted between December 2019 and January 2021. Originally, the pilot phase was envisioned to be shorter, but the declaration of a worldwide pandemic in early 2020 disrupted the normal flow of events. While the pilot continued at the study site, HGOPY, which is a women's and children's hospital that continued its operations throughout the pandemic, progress on various aims of the pilot phase was slower due to a smaller number of patients arriving at HGOPY and the full study team's inability to be on site (due to travel restrictions).

Despite the pandemic, however, the pilot phase successfully achieved its aims, which are briefly summarized here. **First**, the pilot phase aimed to train the nurse counselors at the family planning (FP) unit of HGOPY to become comfortable using the tablet-based app for every client receiving FP counseling. As the nurses were directly involved in the creation of the tablet-based app, this was not a major challenge. However, as the study aimed to trial two different counseling approaches, both of which are embedded into the app and can be randomly assigned for each client, they had to become fully comfortable with each counseling approach: the status quo counseling regime and the new recommendation regime. This meant that the study team initially forwent random assignment of these approaches but, instead, assigned each regime for weeks at a time – until the nurses were fully comfortable with the protocols for each approach. During this period, the nurses gave constant feedback to the study team on various aspects of using the app, which helped the study team tweak the app to make it more user-friendly as a job-support tool.

During this period, it also became clear that training only the three nurses in the FP unit was sufficient neither from the perspective of efficient operations at HGOPY nor the necessary sample size for the study. For HGOPY, it was becoming clear that a lot of clients, including those that arrived at the hospital for pre-natal care, gynecological problems, deliveries, or early childhood vaccinations, were leaving the hospital without receiving FP counseling and adopting a suitable method for birth control. For the study team, the number of clients counseled at the FP unit meant that the adaptive experiment would have to take a long time to recruit the number of clients needed to have a study with sufficient power. It was also clear that the clients who presented at HGOPY explicitly seeking FP counseling and services were a significantly different group than those who presented at HGOPY for other reasons. Therefore, the **second** task for the study team was to recruit seven more nurses from the maternity and gynecology wards and train them in FP counseling and administration. This was a harder task, as these skills are not routinely acquired in school prior to becoming a practicing nurse. So, unlike the nurses in the FP department, who had years of experience counseling patients and administering (and removing) contraceptives, nurses from other departments needed a week-long training course (created specifically for this purpose by HGOPY) to be able to counsel patients, followed by on-the job training (supervised by nurses from the FP unit). After this training, the number of clients counseled by a nurse using the app from a department outside of FP started to slowly increase: by the end of the pilot phase, the number of clients counseled at the FP unit were roughly equal to that counseled elsewhere in HGOPY. This helped the hospital counsel many more women, who would have left the hospital without having been counseled and increased the number of clients counseled daily. As we mention below, many of these marginal clients adopted long-acting methods, which is expected to reduce the incidence of unintended pregnancies. In addition, the nurses from other departments saw this as a great

opportunity for professional development: being able to provide a complementary service to their clients in their own units was clearly welcome and many of them asked to receive the training that HGOPY delivered in early 2020.

**Third**, the pilot phase allowed the study team to optimize a number of study protocols. These included, among others, how to handle return clients, the frequency of syncing the data from the tablets to the cloud, anonymization/de-identification of the data on the hospital's servers before they could be shared with the study team, and the nurse counselors' exact forms of interaction with the app platform. For example, the nurses were initially able to conduct 'practice counseling sessions' with each other or with friends and family that they later deleted from the tablet before uploading cases to the server. After a number of months, once all 10 study nurses were fully comfortable with the app, the study team disabled the ability of providers to delete cases from the tablet, allowing the study team to track the universe of cases – regardless of their completion status or outcome.

**Fourth**, in addition to providing training and optimizing study protocols, the data from the pilot phase have been useful in revising the design of the trial and putting together this pre-analysis plan. Once the nurses were trained, the app revised, and study protocols optimized, we were able to experiment with the main components of our proposed study design – namely providing price discounts and alternative counseling strategies, randomly assigned at the individual level in a static (as opposed to adaptive as proposed for the main study), factorial design setting. We have discussed main findings from this phase, which was between dates March 2020 to January 2021, elsewhere. But, the main takeaways that influenced study design and the construction of this pre-analysis plan were the following:

1. Price discounts for short-acting methods (i.e. the pill and the injectable, SARCs for short) did not affect any of the primary outcomes of interest – neither directly nor in interaction with the price discounts for long-acting methods. Therefore, we have decided to forgo price discounts for SARCs in the adaptive study. These short-acting methods come at a substantially lower upfront cost to the clients (although they are much more expensive per month of protection provided) and, hence, discounts on their relatively low regular prices proved ineffective in increasing uptake.

2. Seeking out women who came to HGOPY seeking services other than FP to offer them counseling proved very effective. This group of women have a lower likelihood of adopting a method after being counseled, but our interventions are highly effective in increasing the uptake of LARCs and reducing the chances of the client leaving the hospital with no method. This finding led to the hospital to revise its protocols, so that providers in various departments would ask clients whether they would like to receive a free FP counseling session before they leave HGOPY.

3. The data from the pilot phase also allowed us to observe potential sources of heterogeneity in impacts. In particular, we noticed heterogeneity by age (teenager vs. adult), by department at which the client presented, and by whether the client already had a method in mind to adopt or was open to discuss all alternatives. These data allowed us to create a policy that we want to evaluate at the end of the adaptive trial, as well as helped us specify our pre-analysis plan to analyze heterogeneity of impacts.

# 2   Experiment design

In this section we will present our experiment design in detail. Sections 2.1.3 to 2.1.5 define relevant quantities, Section 2.2 formalizes our experiment objectives, and Section 2.3 shows treatments are assigned to fulfill these objectives.

## 2.1   Overview and definitions

### 2.1.1   Timeline

Our experiment is divided into the following phases.

1. "Pilot" phase: March 2020 to Feb 2021. This phase is already concluded and was discussed in the introduction. Its data will be used to inform the simulations in this pre-analysis plan and parameter choices in subsequent phases.

2. "Learning" phase: March 2021 to May 2021 (expected). Data will be collected adaptively. These data will be used to learn a desirable treatment assignment "policy", that is, a method for subsidizing and recommending LARC methods that, in expectation, would result in a smaller number of unwanted pregnancies in a cost-effective manner.

3. "Evaluation" phase: May 2021 to December 2021 (expected). Data will be collected non-adaptively, but it will mostly follow the policy learned in the previous phase. These data will be used to evaluate the learned policy.

4. "Follow-up" phase: January 2022 to December 2022. No new clients are recruited into the study. Phone-based follow-up interview are conducted.

When we refer to the "main experiment," which is adaptive, we are referring to phases 2-4 above.

### 2.1.2   Notation

For convenience, here we briefly collect most symbols used in this paper. Their precise definitions are given in the next subsections.

Data are numbered according to the each client's arrival time $t$. A client's observable characteristics, or "contexts", are denoted by $X_t \in \mathcal{X}$ (see Section 2.1.3). The assigned treatment is denoted $W_t \in \mathcal{W}$ (see Section 2.1.4). Once the treatment is assigned we observe several outcomes of interest, such as the contraceptive method $M_t$ chosen by the client (observed immediately through administrative data), whether they discontinued their current method (observed at a four-month follow-up through a phone survey), and whether or not they had an unintended pregnancy $G_t$ (observed 12 months after the counseling session through a follow-up phone survey). The complete list of outcomes is on Table 1 and discussed in detail in Section 2.1.5.

| Variable | Meaning |
|----------|---------|
| $M_t$ | Adopted method (categorical) |
| $L_t$ | Adopted method is LARC (binary) |
| $N_t$ | Adopted method is none/condoms (binary) |
| $C_t$ | Cost of subsidy (real-valued) |
| $F_t$ | Failure rate of adopted method (real-valued) |
| $G_t$ | Unwanted pregnancy within a year (binary) |

Table 1: Notation for different outcomes of interest.

In order to characterize our objective we will make use of potential outcome notation (Imbens and Rubin, 2015). For example, $M_t(w)$ is a random variable representing the contraceptive choice of client $t$ if they were to be assigned arm $w$. In reality can only observe the potential outcome associated with the treatment arm $W_t$ to which the client was actually assigned; e.g., $M_t \equiv M_t(W_t)$.

The time index $t$ represents the chronological order of client arrival to their first counseling session. To simplify some of the algorithmic notation below, the numbering includes the 886 observations in the pilot data, so that the adaptive experiment effectively starts at time t = 887.

Finally, by a "policy" we mean a deterministic mapping from contexts to treatments. Policies will be often denoted by the letter $\pi$.

### 2.1.3   Contexts

The "app" and our follow-up surveys collect additional information about each client (see separate supporting documentation), but our experimental design and the simulations will depend on only four covariates. Following the adaptive experimentation literature, we call these covariates "contexts". The four covariates are: the client's age in years ("age"); a binary indicator representing whether they came to the hospital specifically for a family planning counseling session ("own initiative") or were there for other purposes; an indicator for a delivery (live birth or stillbirth after 28 weeks of pregnancy) within the three months leading up to the session ("recent pregnancy"); and a binary indicator for whether they had a particular method in mind when they came to the counseling ("method in mind").[10] However, as we'll discuss in Section 2.3, we will not use "method in mind" or "own initiative" for personalization, due to the fact that a client could potentially misrepresent her contraceptive preferences in order to get a different discount. We still need information on "method in mind" because, as discussed in the next section, the set of available treatments is different (and the app functions differently) depending on their answer to this question.

---

[10]The exact wording of the question is: "Is there a specific method you absolutely want to adopt?" This question asks the clients to state a strong interest in a specific method. Furthermore, the providers ask a few further questions to make sure the client knows and understand their preferred method, and to potentially dispel any misconceptions about this or other methods.

### 2.1.4 Treatment arms

Each client will be assigned one of four LARC prices and one of two recommendation styles ("view"), as shown on Table 2. However, only clients that do not have a method in mind when arriving at the consult receive a randomized "view", because clients who do have a method in mind at the time of the consultation deterministically receive a "sequential" view with their preferred contraceptive method displayed at the top. This restriction avoids forcing clients, who might simply want to renew their current method or know exactly the method they would like to adopt, to avoid a lengthy discussion about the relative pros and cons of alternative methods via the "side-by-side" style of status quo counseling. Therefore, depending on whether the client "has a method in mind" or not, the number of arms can be eight (4 "LARC" × 2 "view") for some clients and four (4 "LARC" × 1 "view") for others.

Relative to the pilot experiment, LARC prices have been modified by removing the "small" price of 150. The description of the "view" arm remains unchanged from the pilot.

| Treatment arm | Available values |
|---|---|
| LARC prices | 0, 1000, 2000, 4000 |
| View* | Side-by-side (SBS), Sequential (SEQ) |

Table 2: Treatment arms during the main experiment. SARC prices are fixed at their full prices. [*] Only randomized for clients who did *not* have a method in mind.

For the remainder of the paper, we will denote treatment arms by the tuple (LARC price, view). For example, $(2000, \text{SEQ})$ represents the arm that assigns LARC prices of 2000 CFA and recommends contraceptives using the "sequential" view.

### 2.1.5 Outcomes

We are interested in studying the effect of treatment on several outcomes of interest (Table 1). Let's define them more precisely here. We present all the definitions below in terms of potential outcomes to emphasize that their distributions depend on arm assignment but, of course, we can only observe the outcome associated with the arm to which the client was assigned.

The variable $M_t(w)$ is a categorical variable denoting the contraceptive method selected by the client upon being offered treatment arm $w$,

$$M_t(w) \in \{\text{Implant, IUD, Pill, Injectable, None}\}. \tag{1}$$

For convenience we also define an indicator for LARC adoption,

$$L_t(w) := \begin{cases} 1 & \text{if client adopted IUD or implant} \\ 0 & \text{otherwise.} \end{cases} \tag{2}$$

10

Given a client's adopted contraceptive, we can associate with that choice a prediction for their probability of an unwanted pregnancy within the next year, which we call "failure" and denote by

$$F_t(w) := \begin{cases} 0.05\% & \text{if client adopted implant} \\ 0.8\% & \text{if client adopted IUD} \\ 6\% & \text{if client adopted pill} \\ 9\% & \text{if client adopted injectable} \\ 25\% & \text{otherwise} \end{cases} \tag{3}$$

The expected failure rates in (3) are obtained using the 12-month unintended pregnancy rates under typical use (Hatcher, 2007, Table 3.2). The 25% rate for those leaving the hospital without adopting any of the four methods offered by HGOPY, is an estimate of the authors for the risk of an unintended pregnancy during the next 12 months. The same estimates are 18% for those using male condoms and 85% for those using no birth control method. As the clients who adopt none of the four modern methods offered by HGOPY leave the hospital with a free packet of condoms and are counseled and encouraged to use them, we chose a value that is close to the unintended pregnancy rate with the typical use of male condoms. Ideally, we would prefer to use probabilities that take the client's characteristics into account (e.g., young clients may be more prone to fail to follow the daily pill regimen and experience higher failure rates), to our knowledge such estimates are not available.

The next variable is the amount of subsidy that will be disbursed towards LARCs for each client, which represents a "cost" from the perspective of the experiment designer. This quantity depends on the LARC price offered to the client (which in turn depends on the arm assigned to them), and on the contraceptive that the client adopted. If the client does not adopt a LARC, then the cost is zero. Otherwise, the cost will be the difference between the maximum LARC price (4000 CFA) and the price implied by the assigned arm. That is,

$$C_t(w) := \begin{cases} 4000 - (\text{assigned LARC price}) & \text{if client adopted implant or IUD} \\ 0 & \text{otherwise} \end{cases} \tag{4}$$

Finally, the pregnancy indicator variable

$$G_t(w) := \begin{cases} 1 & \text{if client has an unwanted pregnancy within a year of counseling} \\ 0 & \text{otherwise} \end{cases} \tag{5}$$

is only observed one year after the client's arrival, via follow-up interviews.

### 2.1.6 Policies

A "policy" is a deterministic mapping from contexts to treatments. The "control policy" $\pi^{\text{ctrl}}$ does not provide any contraceptive subsidies and assigns the side-by-side "view" where available:

$$\pi^{\text{ctrl}}(x) := \begin{cases} (4000, \text{SEQ}) & \text{if method in mind} = 1 \\ (4000, \text{SBS}) & \text{otherwise.} \end{cases} \tag{6}$$

11

We call "fixed policies" those that (almost) do not vary with contexts: they always assign the same price, and assign either only the "sequential" view or a mixture of "sequential" and "side-by-side" depending only on whether the client had a method in mind. For example:

$$\pi^{2000,\ \text{SBS}}(x) = \begin{cases} (2000, SBS) & \text{if method in mind} = 0 \\ (2000, SEQ) & \text{otherwise} \end{cases} \tag{7}$$
$$\pi^{2000,\ \text{SEQ}}(x) \equiv (2000, SEQ)$$

Other fixed policies for the remaining prices are defined similarly. Note that the control policy (6) is a type of fixed policy with full prices, i.e., $\pi^{\text{ctrl}} \equiv \pi^{4000,\ \text{SBS}}$.

We will be interested in understanding the performance of the following pre-specified policy,

$$\pi^{\text{pre}} = \begin{cases} (0, SEQ) & \text{if } \texttt{recent delivery} = 1 \text{ or age} < 20 \\ (4000, SEQ) & \text{otherwise.} \end{cases} \tag{8}$$

Finally, the goal of the "learning" phase will be to estimate a desirable policy (we elaborate on the specific objective in the next section). This learned policy will be denoted by $\hat{\pi}$.

## 2.2 Experiment objectives

In a nutshell, the objective of this experiment can be succinctly stated as *to learn a treatment assignment policy minimizes unwanted pregnancies in a cost-effective manner, without making too many mistakes during data-collection.* Let's unpack this sentence.

**Learning objective**   During the "learning" phase of the experiment, we will adaptively collect data to so as to learn a policy that satisfies two desiderata. First, this learned policy should assign arms to minimize average number of unwanted future pregnancies (5); but since pregnancies are not immediately observed, we instead minimize the expected contraceptive failures (3) as a proxy. Second, the learned policy should be cost-effective, in the sense that it provides subsidies according to the price-sensitivity of each client.

We characterize characterize these competing aims by defining a random variable that represents our "loss," given observables. This loss is a weighted sum of the failures $F_t(w)$ and costs $C_t(w)$:

$$loss_t(w) := \alpha_t \cdot F_t(w) + \lambda \cdot C_t(w). \tag{9}$$

The first term in (9) represents our need minimize unwanted pregnancies. The factor $\alpha_t > 0$ is a deterministic function of the client's age. It equals $c$ for the clients that are 15 years old (the youngest age allowed in the experiment), and decays linearly until 20 years old, at which point it equals 1.[11] The value of $c$ is a tuning parameter whose effect will be discussed in Section 2.3. By default, we set $c = 2$.

---

[11]That is, $\alpha_t := \max\{c - (\text{Age}_t - 15) \cdot (c-1)/5, 1\}$.

The second term in (9) represents a need to minimize cost where possible. The parameter $\lambda > 0$ controls the trade-off between minimizing costs and failures. We set this parameter to $\lambda = 1.3 \times 10^{-5}$; see the discussion in Appendix C.1 for more details.

**Regret objective**   One additional requirement of our experiment will be that, as much as possible, we minimize the number of "mistakes" during the experiment. That is, that we keep in mind the goals of pregnancy- and cost-minimization even during data collection. This is attained by acting adaptively during the "learning" phase, and by concentrating data collection towards policies of interest in the "evaluation" phase. The adaptive data collection will use a modified version an algorithm called Exploration Sampling, which is explained in detail in Section 2.3.1.

## 2.3   Treatment assignment

The goal of this section is to provide a high-level overview of how treatments are chosen during the "learning" and "evaluation" phases. The mathematical details are in Appendix C.

### 2.3.1   "Learning" phase

Using data gathered during the pilot, before the "learning" phase starts we select an appropriate partition of the covariate space into "regions". These regions do not change during the experiment, but within each region, we will proceed as if we had a non-contextual adaptive experiment. That is, we will sequentially change the region-wide assignment probabilities to each arm depending on past data collected in that region. Assignment probabilities are updated in batches of about 100 observations.[12]

**Finding regions**   We begin by dividing the pilot data into two subsets, $S_1^{\text{pilot}}$ and $S_2^{\text{pilot}}$. Using the first subset $S_1^{\text{pilot}}$ as input, we fit a particular type of decision tree. We then partition the covariate space by identifying each "leaf" of the decision tree with a different region. This step is described in detail in Appendix C.3.

The decision tree used to partition the covariate space is the output of the software `policytree` (Sverdrup et al., 2020), which is based on the offline policy learning algorithm of Zhou et al. (2018) and Athey and Wager (2021):

$$\widehat{\varphi} = \arg \max_{\varphi \in \Phi(d)} -\widehat{\mathbb{E}} \left[ loss_t(\varphi(\tilde{X}_t)) \right] \tag{10}$$

where $\widehat{\mathbb{E}}[\cdot]$ is an unbiased estimate of the loss (9) based on doubly-robust estimation (see Appendix C.2) that uses data in $S_1^{\text{pilot}}$, $\tilde{X}_t$ is a random vector of covariates that does not include "method-in-mind" (see Section 2.1.3), and $\Phi(d)$ is the set of trees of depth $d$. In practice we select $d = 2$, so

---

[12]In the actual experiment deployment we expect there to be minor deviations from 100 observations due to the logistics of updating assignment probabilities in the app.

that $\widehat{\varphi}$ maps contexts to at most four different values (i.e., the four "leaves" of a depth-two tree). This mapping induces a data-driven partition of the covariate space, that is:

$$R_j := \{x \in \mathcal{X} : \widehat{\varphi}(x) = j\}. \tag{11}$$

Once we have these regions, we will act as if we had a non-contextual adaptive experiment within each region. This will be described shortly. Note that, as the name suggests, `policytree` is usually used to find a treatment assignment policy directly – i.e,. to find a mapping between contexts and rewards that approximately minimizes the average loss. However, through simulations we have found that, in the context of this specific adaptive experiment, we are able to find a better treatment assignment policy by simply using this algorithm as a preliminary step to divide the covariate space into regions, and then discarding the information about which arm $\widehat{\varphi}$ suggests should be assigned. To avoid confusion, we will call the function $\widehat{\varphi}$ a "partitioning function" and reserve the term "policy" for the final treatment assignment rule obtained at the end of the "learning" phase.



Figure 1: Overview of treatment assignments during the "learning" phase. At the beginning of the phase, we partition of the covariate space into four regions. Within each region, we run a separate (non-contextual) adaptive experiment. Probabilities are via the Exploration Sampling algorithm (Kasy and Sautmann, 2021). At the end of the "learning" phase we assign the arm that had the highest Exploration Sampling probability within each region.

**Computing assignment probabilities** At the beginning of each batch $m$ of about 100 observations, within each region we update arm assignment probabilities via a modified Exploration Sampling (ES) algorithm (Kasy and Sautmann, 2021). This is described more precisely in Appendix C.5, but an overview follows. For each region $j$:

1. Using a non-informative Bayesian model with Normal prior and Normal likelihood, compute

the posterior probability $\tilde{e}_t(w, j)$ that arm $w$ maximizes the average negative loss[13] (9) within region $j$, for all available arms. This is done using the remaining pilot data $S_2^{\text{pilot}}$ and all the data collected during the "learning" phase so far. Unavailable arms (see Section 2.1.4) get probability zero.

2. Compute "exploration sampling probabilities" by taking the product of each probability with its complement and then re-normalizing so that the probabilities sum to one; i.e.,[14]

$$e_t(w, j) \propto \tilde{e}_t(w, j)(1 - \tilde{e}_t(w, j)). \tag{12}$$

3. We impose a minimum "floor" on assignment probabilities, so that every available arm gets at least some probability $\bar{p}_0^{\text{all}}/K_t$ of being selected, where $K_t$ is the number of available arms for client $t$ (see Section 2.1.4). This introduces some additional exploration, but more importantly it ensures that, if needed, we are still able to analyze these data using methods that require positive assignment probabilities everywhere (i.e., methods based on inverse-probability weighting). We set $\bar{p}_0^{\text{all}} = .2$.

**Freezing the policy**  At the end of the "learning" phase, within each region we select the arm that the highest exploration sampling probability or, equivalently, the arm that had the highest posterior probability of being optimal. That is, for any context $x$ falling in region $j$, we set:[15]

$$\hat{\pi}(x) = \arg\max_w e_{\tilde{T}+1, j}(w), \tag{13}$$

where $\tilde{T}$ is the last period of the "learning" phase.

### 2.3.2 Evaluation phase

In this phase we assign treatments non-adaptively, but also non-uniformly. The main goal of this phase is to collect data to accurately evaluate the learning policy $\hat{\pi}$ and compare it to the control policy (6). Therefore, in this phase we assign arms according to this non-adaptive rule:

$$W_t = \begin{cases} \hat{\pi}(X_t) & \text{with prob. } (1 - \bar{p}_1^{\text{all}})\bar{p}^{\text{opt}} + \bar{p}_1^{\text{all}}/K_t, \\ \pi^{\text{ctrl}}(X_t) & \text{with prob. } (1 - \bar{p}_1^{\text{all}})(1 - \bar{p}^{\text{opt}}) + \bar{p}_1^{\text{all}}/K_t, \\ \text{any other arm} & \text{with prob. } \bar{p}_1^{\text{all}}/K_t. \end{cases} \tag{14}$$

That is, for each new client, a $\bar{p}_1^{\text{all}}$ fraction of the time we select treatment arms uniformly, so any treatment arm has probability at least $\bar{p}_1^{\text{all}}/K_t$, where $K_t$ is the number of arms available. In the remaining $1 - \bar{p}_1^{\text{all}}$ fraction of the time we either select the arm suggested by the optimal policy $\hat{\pi}(X_t)$

---

[13]These are sometimes called Thompson Sampling probabilities.

[14]For the case of three or more arms, this is equivalent to selecting the arm with highest probability of being optimal according to our posterior (for the case of two arms, exploration sampling forces sampling with equal probability).

[15]We have also experimented with other selection methods, such as selecting the policy with highest sample mean or posterior mean, and found very similar results to those described here.

or the one suggested by the control policy $\pi^{\mathrm{ctrl}}(X_t)$, with the share of the former being controlled by the parameter $\bar{p}^{\mathrm{opt}}$.

The rationale for this particular parametrization is the following. Our primary hypotheses involve the optimal and control policies, so we will select them a majority $(1 - \bar{p}_1^{\mathrm{all}})$ fraction of the time, assuming that $\bar{p}_1^{\mathrm{all}}$ is a small number. This extra randomization $\bar{p}_1^{\mathrm{all}}$ is so that we are also able to test secondary hypotheses that have to do with other treatment arms, such as those related to the pre-specified policy (8). Finally, from a regret-minimizing perspective it will be advantageous to select the optimal policy as much as we can, and we control the share of optimal policy via the parameter $\bar{p}^{\mathrm{opt}}$. We set $\bar{p}_1^{\mathrm{all}} = .2$ and $\bar{p}^{\mathrm{opt}} = 0.8$.

### 2.3.3    Follow-up phase

Once the counseling session is over, we will conduct phone-based follow-up interviews.[16]

- Within two weeks after the session, we will ask the clients about their experience with their counseling session, using a module that is designed to get at the quality of service they received, and ask about their level of satisfaction with the services they received at HGOPY.

- After a 16-week period we will ask clients about their satisfaction with the current contraceptive and whether they have decided to discontinue the method or switch to a different one.

- After 52 weeks we will ask the clients about any unwanted pregnancies and thus gather data on the realization of $G_t$. We expect somewhere between a 10-20% attrition rate for this period.

## 2.4    Post-experiment Analysis

Once the experiment is over we will conduct a series of analyses using the data collected during the "evaluation" phase.

### 2.4.1    Comparing optimal policy against control

Our primary hypotheses compare the policy $\hat{\pi}$ that was obtained at the end of the "learning" phase with the control policy $\pi^{\mathrm{ctrl}}$ defined in (6).

**Main effects**    We are primarily interested in testing if the learned policy $\hat{\pi}$ improves upon the control policy in terms of the probability that the client takes up a LARC $L_t$ or leaves without any

---

[16]We will update the list of questions asked in the during the consultation and follow-up surveys in a separate supplementary appendix.

modern contraceptive $N_t$. Therefore our main hypotheses are:

$$H_0 : \mathbb{E}\left[L_t(\hat{\pi}(X_t))\right] \leq \mathbb{E}\left[L_t(\pi^{\mathrm{ctrl}}(X_t))\right] \qquad H_a : \mathbb{E}\left[L_t(\hat{\pi}(X_t))\right] > \mathbb{E}\left[L_t(\pi^{\mathrm{ctrl}}(X_t))\right],$$
$$H_0 : \mathbb{E}\left[N_t(\hat{\pi}(X_t))\right] \geq \mathbb{E}\left[N_t(\pi^{\mathrm{ctrl}}(X_t))\right] \qquad H_a : \mathbb{E}\left[N_t(\hat{\pi}(X_t))\right] < \mathbb{E}\left[N_t(\pi^{\mathrm{ctrl}}(X_t))\right]. \tag{15}$$

The empirical counterparts of the expectations in (15) and (16) that are used for hypothesis testing are based on augmented inverse propensity-weighted (AIPW) estimators (see Appendix C.2).

In addition, as a secondary hypothesis we will test whether the frequency of unwanted pregnancies one year after the consultation is smaller under the optimal policy than under control.

$$H_0 : \mathbb{E}\left[G_t(\hat{\pi}(X_t))\right] \geq \mathbb{E}\left[G_t(\pi^{\mathrm{ctrl}}(X_t))\right] \qquad H_a : \mathbb{E}\left[G_t(\hat{\pi}(X_t))\right] < \mathbb{E}\left[G_t(\pi^{\mathrm{ctrl}}(X_t))\right] \tag{16}$$

Finally, recall from Section 2.3.3 that in our follow-up surveys we will also monitor other quantities of interest, such as whether the client was satisfied with the interview, and whether they discontinued the method (and if so, when and why). Although these variables will provide us with additional descriptive evidence about the client's trajectory after the consultation, at this point we do not expect to include them in our set of primary or secondary hypotheses.

**Heterogeneous treatment effects** We will test whether particular subgroups react differently to our learned policy compared to control.

One important subgroup of interested are adolescent females. We will perform three hypothesis tests with null hypothesis of the form

$$H_0 : \mathbb{E}\left[L_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[L_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} \geq 20\right]$$
$$H_0 : \mathbb{E}\left[N_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[N_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} \geq 20\right] \tag{17}$$
$$H_0 : \mathbb{E}\left[G_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[G_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{age} \geq 20\right]$$

In addition, we will perform a two-sided test for the null hypothesis that the effect of switching from control to the learned policy is equal for teenagers and adults:

$$H_0 : \mathbb{E}\left[L_t(\hat{\pi}(X_t)) - L_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[L_t(\hat{\pi}(X_t)) - L_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} \geq 20\right]$$
$$H_0 : \mathbb{E}\left[N_t(\hat{\pi}(X_t)) - N_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[N_t(\hat{\pi}(X_t)) - N_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} \geq 20\right] \tag{18}$$
$$H_0 : \mathbb{E}\left[G_t(\hat{\pi}(X_t)) - G_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} < 20\right] = \mathbb{E}\left[G_t(\hat{\pi}(X_t)) - G_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{age} \geq 20\right].$$

Next, since our treatment assignment policy $\hat{\pi}$ acts differently depending on the region of the covariate space (11), it is natural to test whether the optimal policy is able to improve upon control in each region. That is, for each region $j$,

$$H_0 : \mathbb{E}\left[L_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{Region}\ j\right] = \mathbb{E}\left[L_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{Region}\ j\right]$$
$$H_0 : \mathbb{E}\left[N_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{Region}\ j\right] = \mathbb{E}\left[N_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{Region}\ j\right] \tag{19}$$
$$H_0 : \mathbb{E}\left[G_t(\hat{\pi}(X_t)) \,\big|\, \mathrm{Region}\ j\right] = \mathbb{E}\left[G_t(\pi^{\mathrm{ctrl}}(X_t)) \,\big|\, \mathrm{Region}\ j\right].$$

We may also run similar hypothesis tests for subgroups discovered by a data-driven algorithm such as causal trees (Athey and Imbens, 2016).

### 2.4.2 Other policies

We will also test the hypotheses (15) and (16) replacing the learned policy $\hat{\pi}$ by the pre-specified policy $\pi^{\text{pre}}$ described in (8).

## 3 Numerical experiments

In this section, we show how simulations based on the pilot data (Section 3.1) helped us select our algorithm's tuning parameters (Section 3.2) and quantify the predicted benefits from our adaptive design (3.3).

### 3.1 Simulation design

The data-generating process used in all simulations is detailed in Appendix B, but at a high level we proceed as follows.

At the beginning of each simulation we draw a sample with replacement from the pilot data, and use this sample to estimate the probability that a client will adopt each contraceptive based on their contexts and treatments $\hat{p}(x, w)$. More specifically, we use a penalized multinomial logistic regression model on a set of polynomials of the four covariates and treatment arms. Also, in order to better explore the space of possible outcome models, this model is fit on a sample with replacement of the pilot data, and we also randomly perturb the amount of L2-penalization. We also use the same sample above to fit the joint distribution of contexts $\hat{p}(x)$.

During the simulation, for each incoming client, we first draw covariates $X_t$ from the fitted distribution of contexts $\hat{p}(x)$; next, we decide on a treatment assignment $W_t$ as explained in Section 2.3; finally, we draw a contraceptive from the fitted multinomial model $M_t \sim \hat{p}(X_t, W_t)$.

Given the selected contexts, treatment and contraceptive, most other outcomes are deterministic. For example, if we know the contraceptive choice, we know whether or not it is a LARC ($L_t$).

The only outcome that remains to be determined is the pregnancy indicator $G_t$. For this variable, we use the probabilities in (3) as ground truth. For example, (3) implies that a client that adopted the pill has a 6% chance of incurring an unwanted pregnancy, hence for this client we draw $G_t \mid M_t \sim$ Bernoulli(0.06). Moreover, as mentioned before we expect an attrition rate up to 20% for the actual pregnancy rates $G_t$. To model this attrition mechanism, at each simulation we draw $p_{\text{attrit}} \in \{.1, .15, .2\}$ uniformly at random and then replace $p_{\text{attrit}}$ of our data on $G_t$ by missing values.

## 3.2 Tuning parameters

Our experiment design in Section 2 relied on a set of tuning parameters whose values must the decided by the experimenter. In this subsection, we validate our choices for these tuning parameters via simulations. At a high level our procedure is as follows:

1. Use the pilot data to estimate the underlying data-generating process.

2. Simulate the main experiment, learn and evaluate an optimal policy as in Section 2.3.

3. Compute and store:
   - The average value of the optimal policy relative to control the control policy (6) and the amount of power we have to test our main (15) and secondary hypotheses (16);
   - The average regret during the experiment for each outcome variable (i.e., the loss incurred during the main experiment, relative to an oracle policy (9)).

4. Select parameters that satisfy our experiment objectives as described in Section 2.2.

The list of parameters that need to be tuned and their selected ("default") values on Table 3. The default values will be used in the actual experiment.

| Parameter | Symbol | Default |
|---|---|---|
| Total length of experiment (including pilot, learning and evaluation phases) | T | 1800 |
| Fraction of "learning" phase relative to "learning" plus "evaluation" | None | .5 |
| Failure weight for younger clients | $c$ | 2. |
| Depth of partition function $\widehat{\varphi}$ | $d$ | 2 |
| Floor parameter on assign. prob. for all arms during "evaluation" phase | $\bar{p}_1^{\text{all}}$ | .2 |
| Share of $\hat{\pi}(X_t)$ during "evaluation" when not acting uniformly at random | $\bar{p}^{\text{opt}}$ | .8 |

Table 3: Tuning parameters. Default values are used during the actual experiment.

In order to better understand the role of each parameter, we will vary one set of parameters at a time, and see how they affect our estimates at the end of the experiment.

Before showing results for the optimal policy, however, it is instructive to understand the behavior of "fixed" policies. Figure 2 shows the value of different fixed policies. The most notable takeaway is that policies that fully subsidize LARCs for everyone attain a LARC adoption rate of around 50% for an average per-person cost of about 2000 CFA, compared to a LARC adoption rate 25% for the control policy for a cost of zero. In addition, the effect of "sequential" versus "side-by-side" view is essentially the same on average.

Figure 3 shows how results change as we vary the total length of the experiment. In particular we see that at the chosen length of $T = 1800$ we are well-powered to test our main hypotheses (15), though we may not be powered to test the secondary hypothesis (16).

19

Similarly, Figure 4 shows how results vary as we fix the total length of the experiment to $T = 1800$, but change the fraction of the main experiment data that is devoted to the "learning" phase (as opposed to the "evaluation" phase). Allowing for a longer "learning" phase leads to a policy that performs better in terms of contraceptive adoptions and therefore pregnancy prevention, but because the "evaluation" phase is shorter the standard error is smaller, decreasing power. We set the default value to .5, indicating that we spend 50% of the main experiment in the "learning" phase, and the remaining 50% in the "evaluation" phase.

Figures 5 and 6 show how our results change as we vary how many observations are assigned according to the learned policy during the "evaluation" phase, i.e., as we vary $\bar{p}_1^{\text{all}}, \bar{p}^{\text{opt}}$. Neither of these parameters has an effect in the value of the learned policy, but they affect how accurately we can evaluate it.

Figure 7 shows that increasing the failure weight for adolescents from the default $c = 1$ slightly decreases the fraction of clients leaving with no contraceptive, but it also increases cost. This likely happens because a higher adolescent weight forces the selection of a more expensive treatment in regions where there are many teenagers.

Finally, Figure 8 shows what would happen if our partition function $\widehat{\varphi}$ (see Section 2.3.1) were a tree of maximal depth 1, 2 or 3; i.e., at most two, four or eight leaves. There seem to be large gains from going from depth 1 to 2, in particular in terms of cost (a decrease from 1450 to 1300), but beyond that the benefits are not clear. Since larger leaves also allow us to have more information about our chosen subgroups, we set the default depth to 2 (four leaves).

Figure 2: Value of "fixed" policies in our simulations (i.e., policies that do not personalize by context). Left column shows average true value of each policy across simulations, in terms of each key outcome (row). Right column shows standard error of that value at $T = 1800$, when all tuning parameters are set to their default values.

Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

Figure 3: Varying the **total length of the experiment**, keeping the relative size between "learning" and "evaluation" rates fixed.

Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

Figure 4: Varying **relative length of the learning phase**, as a fraction of the length of the "learning" plus "evaluation" phase, keeping other parameters fixed at their default values (in particular, keeping the total experiment length at 1800). Default value is 0.5.

Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

Figure 5: Varying the **share of observations assigned according to learned policy** during the "evaluation" phase. This parameter has no effect in the value of the learned policy, but it affects how accurately we can evaluate it.

Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

Figure 6: Varying the **assignment probability floor during evaluation phase**. This parameter has no effect in the value of the learned policy, but it affects how accurately we can evaluate it. Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

Figure 7: Varying the **adolescent weight** parameter $c$.
Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

26

Figure 8: Varying the **depth of the partition function** $\widehat{\varphi}$. Default value is $d = 2$. Results are averages are across over 10,000 simulations, each using a data-generating process drawn as explained in Section B. Confidence intervals are empirical 95% intervals across simulations.

## 3.3 Benefits of an adaptive experiment design

To assess the statistical benefits of adaptivity in our simulations, we take the simulation results obtained by running adaptive experiments as explained in Section 2, and compare them to what the policy that we would have obtained via the following non-adaptive experiment. We fix the covariate regions as explained in Section (2.3), but during the "learning" phase we simply assign arms uniformly at random. At the end of the "learning" phase we select the arm that has the highest Exploration Sampling probability[17] within each region.

Table 4 compares the learned policy when data are collected adaptively to the one obtained when data are collected non-adaptively. On average, the former is superior in terms of LARC and SARC adoptions, at a minor increase in cost. It's also important to note that, the former is also superior in terms of the (negative) loss (9), indicating that adaptivity is producing better policy according to the criterion it is trying to optimize.

|  |  | Adaptive | Non-Adaptive |
|---|---|---|---|
| LARC adoption | Mean (Avg. True Value) | 0.5069 | 0.5019 |
|  | Std. Err. | 0.0002 | 0.0002 |
| SARC adoption | Mean (Avg. True Value) | 0.0520 | 0.0501 |
|  | Std. Err. | 0.0001 | 0.0001 |
| No method | Mean (Avg. True Value) | 0.4411 | 0.4480 |
|  | Std. Err. | 0.0002 | 0.0002 |
| Pregnancies | Mean (Avg. True Value) | 0.1153 | 0.1169 |
|  | Std. Err. | 0.0000 | 0.0000 |
| Cost (CFA) | Mean (Avg. True Value) | 1563.4114 | 1492.8360 |
|  | Std. Err. | 1.4684 | 1.3415 |
| Failure | Mean (Avg. True Value) | 0.1153 | 0.1169 |
|  | Std. Err. | 0.0000 | 0.0000 |
| Negative Loss | Mean (Avg. True Value) | -0.1414 | -0.1421 |
|  | Std. Err. | 0.0000 | 0.0000 |

Table 4: Average out-of-sample performance comparison between policy learned via an adaptive experiment vs. a non-adaptive one.

Table 5 compares the standard errors associated with our estimates of the difference in value between the learned and control policies, for each outcome. The size of the standard errors tend to be smaller when data are collected following our method.

---

[17]Selecting the arm with highest sample mean or posterior mean has qualitatively the same effect.

|  |  | Adaptive | Non-Adaptive |
|---|---|---|---|
| LARC adoption | Mean (Avg. Std. Error) | 0.0531 | 0.0730 |
|  | Std. Err. | 0.0000 | 0.0000 |
| SARC adoption | Mean (Avg. Std. Error) | 0.0272 | 0.0347 |
|  | Std. Err. | 0.0000 | 0.0000 |
| No method | Mean (Avg. Std. Error) | 0.0549 | 0.0745 |
|  | Std. Err. | 0.0000 | 0.0000 |
| Pregnancies | Mean (Avg. Std. Error) | 0.0473 | 0.0612 |
|  | Std. Err. | 0.0000 | 0.0000 |
| Cost (CFA) | Mean (Avg. Std. Error) | 97.7035 | 180.3143 |
|  | Std. Err. | 0.0587 | 0.1253 |
| Failure | Mean (Avg. Std. Error) | 0.0131 | 0.0179 |
|  | Std. Err. | 0.0000 | 0.0000 |
| Negative Loss | Mean (Avg. Std. Error) | 0.0134 | 0.0177 |
|  | Std. Err. | 0.0000 | 0.0000 |

Table 5: Average comparison of standard errors around estimates of the difference between the learned and control policies.

In addition, Table 6 shows the average performance of each key outcome *during* the experiment. Both during the "learning" and "evaluation" phases, collecting data adaptively leads to higher LARC adoption and fewer clients leaving without adopting anything, which in turn translates into a small decrease in pregnancies.

|  |  | Non-Adaptive | | Adaptive | |
|---|---|---|---|---|---|
|  |  | Learning | Evaluation | Learning | Evaluation |
| LARC adoption | Mean | 0.4170 | 0.4171 | 0.4543 | 0.4438 |
|  | Std. Err. | 0.0001 | 0.0002 | 0.0001 | 0.0002 |
| SARC adoption | Mean | 0.0512 | 0.0512 | 0.0495 | 0.0529 |
|  | Std. Err. | 0.0001 | 0.0001 | 0.0001 | 0.0001 |
| No method | Mean | 0.5318 | 0.5317 | 0.4962 | 0.5033 |
|  | Std. Err. | 0.0001 | 0.0002 | 0.0001 | 0.0002 |
| Pregnancies | Mean | 0.1376 | 0.1377 | 0.1289 | 0.1309 |
|  | Std. Err. | 0.0001 | 0.0001 | 0.0001 | 0.0001 |
| Cost | Mean | 1063.3515 | 1062.6476 | 1246.5276 | 1194.7537 |
|  | Std. Err. | 0.4120 | 0.4675 | 0.4862 | 0.9345 |
| Failure | Mean | 0.1377 | 0.1377 | 0.1287 | 0.1308 |
|  | Std. Err. | 0.0000 | 0.0000 | 0.0000 | 0.0000 |

Table 6: Average performance comparison during the experiment.

# 4 Abridged mock output

In this section, we use output data from a *single* simulation to illustrate the key results that we will present after our actual experiment concluded. The parameters were all set to the default parameters displayed on Table 3. The data-generating process was the same as used in Section B.

## 4.1 Learned policy

**Behavior of treatment assignment**  Figure 9 shows the evolution of assignment probabilities during the experiment, for the four learning batches (named as "Learning-1" through "Learning-4" and during the evaluation phase) and across the four regions. Note that during the learning phase, because we are using exploration sampling probabilities, no arm can receive more than probability one half. In addition, because we are imposing a lower bound on assignment probabilities, every (available) arm receives positive probability.

**Behavior of the learned policy**  Figure 10 shows the learned policy in this particular simulation.



Figure 10: Learned policy in this particular simulation. Recall that the tree structure is fixed, and therefore so are the regions. What is learned at the end of the "learning" phase are the arm assignments within each region.

Figure 9: Evolution of assignment probabilities at the beginning of each batch of the learning phase and during the evaluation phase.

Figure 11: Average covariate per subgroup defined by the region or final arm assignment according to the learned policy.

## 4.2 Policy values

Table 7 compares the expected value of several outcomes of interest under three policies: the estimated optimal policy that was obtained at the end of the "learning" phase, the control policy (6) and the policy that was pre-specified in (8).

| Outcome | Policy | Optimal | Control | Pre-specified |
|---|---|---|---|---|
| Adopted LARC | Mean | 0.443 | 0.236 | 0.346 |
| | Std. Error | 0.029 | 0.044 | 0.040 |
| Adopted Nothing | Mean | 0.493 | 0.714 | 0.543 |
| | Std. Error | 0.029 | 0.047 | 0.049 |
| Adopted SARC | Mean | 0.064 | 0.050 | 0.111 |
| | Std. Error | 0.014 | 0.023 | 0.057 |
| Cost | Mean | 1163.440 | 0.000 | 732.666 |
| | Std. Error | 91.373 | 0.000 | 117.416 |
| Failure | Mean | 0.129 | 0.182 | 0.145 |
| | Std. Error | 0.007 | 0.011 | 0.009 |
| Pregnancy | Mean | 0.101 | 0.223 | 0.175 |
| | Std. Error | 0.018 | 0.048 | 0.060 |

Table 7: Estimates of the expected value of different outcomes for the learned policy ("optimal") $\hat{\pi}$, the control policy $\pi^{\text{ctrl}}$ (6), and for the pre-specified policy (8).

| | Adopted LARC | | Adopted Nothing | | Adopted SARC | |
|---|---|---|---|---|---|---|
| | Mean | Std. Error | Mean | Std. Error | Mean | Std. Error |
| (0, SEQ) | 0.513 | 0.063 | 0.395 | 0.067 | 0.092 | 0.057 |
| (0, SBS) | 0.396 | 0.109 | 0.611 | 0.112 | -0.007 | 0.015 |
| (1000, SEQ) | 0.557 | 0.114 | 0.423 | 0.111 | 0.020 | 0.030 |
| (1000, SBS) | 0.381 | 0.092 | 0.532 | 0.095 | 0.087 | 0.062 |
| (2000, SEQ) | 0.411 | 0.084 | 0.554 | 0.084 | 0.034 | 0.012 |
| (2000, SBS) | 0.367 | 0.109 | 0.514 | 0.115 | 0.119 | 0.079 |
| (4000, SEQ) | 0.307 | 0.057 | 0.649 | 0.058 | 0.044 | 0.017 |
| (4000, SBS) | 0.236 | 0.044 | 0.714 | 0.047 | 0.050 | 0.023 |
| | Cost | | Failure | | Pregnancy | |
| | Mean | Std. Error | Mean | Std. Error | Mean | Std. Error |
| (0, SEQ) | 2051.616 | 251.087 | 0.108 | 0.015 | 0.097 | 0.035 |
| (0, SBS) | 1583.213 | 436.006 | 0.153 | 0.027 | 0.135 | 0.089 |
| (1000, SEQ) | 1669.544 | 340.606 | 0.108 | 0.028 | 0.133 | 0.080 |
| (1000, SBS) | 1142.906 | 277.344 | 0.139 | 0.022 | 0.221 | 0.096 |
| (2000, SEQ) | 822.965 | 167.778 | 0.142 | 0.021 | 0.087 | 0.035 |
| (2000, SBS) | 734.037 | 217.070 | 0.137 | 0.027 | 0.303 | 0.117 |
| (4000, SEQ) | 0.000 | 0.000 | 0.166 | 0.014 | 0.229 | 0.083 |
| (4000, SBS) | 0.000 | 0.000 | 0.182 | 0.011 | 0.223 | 0.048 |

Table 8: Estimates of the expected value of "fixed" policies for each outcome.

## 4.3 Hypothesis tests

**Learned policy vs control policy**   Table 9 displays the results of our main hypothesis tests (15) regarding contraceptive adoption, as well as our secondary hypothesis regarding pregnancies (16). As expected, we can easily detect a large difference in contraceptive adoption and are able to reject our main hypotheses, while changes in pregnancies are harder to detect but have the right sign.

|                 | Mean   | Std. Error | p-values    | Count |
|-----------------|--------|------------|-------------|-------|
| Adopted LARC    | 0.208  | 0.050      | 0.0***      | 457   |
| Adopted Nothing | -0.221 | 0.052      | 0.0***      | 457   |
| Pregnancy       | -0.122 | 0.050      | 0.0076***   | 417   |

Table 9: Testing the average performance of the learned policy versus control for key outcomes.

Table 10 suggests that contraceptive adoption effect is strong enough that we are often able to detect it even if we divide the data according to the covariate regions. Similarly, Table 11 suggests that the effect is more strongly detectable in adult subgroups, reflecting the fact that we have fewer teenagers in our sample.

|                 | Region | Mean   | Std. Error | p-values    | Count |
|-----------------|--------|--------|------------|-------------|-------|
| Adopted LARC    | A      | 0.291  | 0.116      | 0.006***    | 112   |
| Adopted LARC    | B      | 0.139  | 0.062      | 0.0118**    | 117   |
| Adopted LARC    | C      | 0.241  | 0.110      | 0.014**     | 116   |
| Adopted LARC    | D      | 0.162  | 0.105      | 0.0617*     | 112   |
| Adopted Nothing | A      | -0.311 | 0.119      | 0.0046***   | 112   |
| Adopted Nothing | B      | -0.127 | 0.074      | 0.0433**    | 117   |
| Adopted Nothing | C      | -0.224 | 0.113      | 0.0238**    | 116   |
| Adopted Nothing | D      | -0.225 | 0.108      | 0.0187**    | 112   |
| Pregnancy       | A      | -0.096 | 0.081      | 0.1193      | 103   |
| Pregnancy       | B      | -0.271 | 0.114      | 0.0085***   | 107   |
| Pregnancy       | C      | -0.070 | 0.082      | 0.1956      | 107   |
| Pregnancy       | D      | -0.045 | 0.119      | 0.353       | 100   |

Table 10: Testing the performance of the learned policy against the control policy with respect to average key outcomes in each covariate region.

|  | Age Subgroup | Mean | Std. Error | p-values | Count |
|---|---|---|---|---|---|
| outcome | | | | | |
| Adopted LARC | Teenagers | 0.441 | 0.232 | 0.0289** | 31 |
| Adopted LARC | Adults | 0.191 | 0.051 | 0.0001*** | 426 |
| Adopted Nothing | Teenagers | -0.328 | 0.260 | 0.1033 | 31 |
| Adopted Nothing | Adults | -0.213 | 0.053 | 0.0*** | 426 |
| Pregnancy | Teenagers | -0.106 | 0.166 | 0.2612 | 27 |
| Pregnancy | Adults | -0.123 | 0.052 | 0.0095*** | 390 |

Table 11: Testing the performance of the learned policy against the control policy with respect to average key outcomes for each age-based subgroup.

Finally, Table 12 tests whether the effect of switching from control to optimal policies is different each age subgroup. The null hypotheses are specified in (18). In this simulation, we were not able to detect this difference.

|  | Mean | Std. Error | p-values |
|---|---|---|---|
| outcome | | | |
| Adopted LARC | 0.250 | 0.238 | 0.2933 |
| Adopted Nothing | -0.115 | 0.265 | 0.6647 |
| Pregnancy | 0.017 | 0.174 | 0.9217 |

Table 12: Testing the performance of the learned policy against the control policy with respect to average key outcomes for each age-based subgroup. The point estimate ("Mean" column) is the estimates of the difference between the left- and right-hand sides of (18).

**Pre-specified policy vs control policy**   Table 13 compares the manually pre-specified policy $\pi^{\mathrm{pre}}$ to the control policy $\pi^{\mathrm{ctrl}}$. The effects are in the same direction as in the learned policy, but smaller in size and, hence, harder to detect.

|  | Mean | Std. Error | p-values | Count |
|---|---|---|---|---|
| Adopted LARC | 0.111 | 0.051 | 0.0147** | 457 |
| Adopted Nothing | -0.171 | 0.060 | 0.0021*** | 457 |
| Pregnancy | -0.049 | 0.072 | 0.251 | 417 |

Table 13: Testing the average performance of the pre-specified policy versus control.

# References

Athey, S. and Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360.

Athey, S. and Wager, S. (2021). Policy learning with observational data. *Econometrica*, 89(1):133–161.

Baird, S., McIntosh, C., and Özler, B. (2011). Cash or condition? evidence from a cash transfer experiment. *The Quarterly journal of economics*, 126(4):1709–1753.

Baird, S., McIntosh, C., and Özler, B. (2019). When the money runs out: Do cash transfers have sustained effects on human capital accumulation? *Journal of Development Economics*, 140:169 – 185.

Bearak, J., Popinchalk, A., Alkema, L., and Sedgh, G. (2018). Global, regional, and subregional trends in unintended pregnancy and its outcomes from 1990 to 2014: estimates from a bayesian hierarchical model. *The Lancet Global Health*, 6(4):e380–e389.

Daniels, K. and Abma, J. C. (2020). Current contraceptive status among women aged 15–49: United states, 2017–2019. NCHS Data Brief No. 388, US Center for Disease Control.

DeFranco, E. A., Seske, L. M., Greenberg, J. M., and Muglia, L. J. (2015). Influence of interpregnancy interval on neonatal morbidity. *American Journal of Obstetrics and Gynecology*, 212(3):386.e1 – 386.e9.

DHS (2020). Republique du cameroun: Enquete demographique et de sante 2018. Technical report, Institut National de la Statistique du Cameroun (INS), and ICF, Yaounde, Cameroun.

Finer, L. B. and Henshaw, S. K. (2006). Disparities in rates of unintended pregnancy in the united states, 1994 and 2001. *Perspectives on sexual and reproductive health*, 38(2):90–96.

Frost, J. J., Sonfield, A., Zolna, M. R., and Finer, L. B. (2014). Return on investment: A fuller assessment of the benefits and cost savings of the us publicly funded family planning program. *The Milbank Quarterly*, 92(4):696–749.

Hadad, V., Hirshberg, D. A., Zhan, R., Wager, S., and Athey, S. (2019). Confidence intervals for policy evaluation in adaptive experiments. *arXiv preprint arXiv:1911.02768*.

Hatcher, R. A. (2007). *Contraceptive technology*. Ardent Media.

Holt, K., Dehlendorf, C., and Langer, A. (2017). Defining quality in contraceptive counseling to improve measurement of individuals' experiences and enable service delivery improvement. *Contraception*, 96(3):133–137.

Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.

Kasy, M. and Sautmann, A. (2021). Adaptive treatment assignment in experiments for policy choice. *Econometrica*, 89(1):113–132.

Lindo, J. M. and Packham, A. (2017). How much can expanding access to long-acting reversible contraceptives reduce teen birth rates? *American Economic Journal: Economic Policy*, 9(3):348–76.

McKelvey, C., Thomas, D., and Frankenberg, E. (2012). Fertility regulation in an economic crisis. *Economic Development and Cultural Change*, 61(1):7–38.

Mestad, R., Secura, G., Allsworth, J. E., Madden, T., Zhao, Q., and Peipert, J. F. (2011). Acceptance of long-acting reversible contraceptive methods by adolescent participants in the contraceptive choice project. *Contraception*, 84(5):493–498.

Porsdam Mann, S., Savulescu, J., and Sahakian, B. J. (2016). Facilitating the ethical use of health data for the benefit of society: Electronic health records, consent and the duty of easy rescue. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2083):20160130.

Rau, T., Sarzosa, M., and Urzúa, S. S. (2017). The children of the missed pill. Technical report, National Bureau of Economic Research.

Russo, D. J., Van Roy, B., Kazerouni, A., Osband, I., Wen, Z., et al. (2018). A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96.

Sully, E., Biddlecom, A., Darroch, J., Riley, T., Ashford, L., Lince-Deroche, N., et al. (2020). Adding it up: investing in sexual and reproductive health, 2019. *New York: Guttmacher Institute.*

Sverdrup, E., Kanodia, A., Zhou, Z., Athey, S., and Wager, S. (2020). policytree: Policy learning via doubly robust empirical welfare maximization over trees. *Journal of Open Source Software*, 5(50):2232.

Trussell, J. (2011). Contraceptive failure in the united states. *Contraception*, 83(5):397–404. 21477680[pmid].

World Health Organization (2019). Trends in maternal mortality 2000 to 2017: estimates by who, unicef, unfpa, world bank group and the united nations population division. Technical report.

Zhou, Z., Athey, S., and Wager, S. (2018). Offline multi-action policy learning: Generalization and optimization. *arXiv preprint arXiv:1810.04778.*

# A The job-support tool

## A.1 What does tablet-based job-support tool do?

Counseling clients on family planning and contraceptive methods is not straightforward. There are more than 10 methods that can be considered by a client. Each modern method, especially hormonal ones, can cause various side effects, which can vary from person to person. The clients' contexts, such as how long she would like to wait before getting pregnant, her birth history, her previous experience with some of the methods, her preferences for side effects, are relevant as to what methods will be most suitable for her. Finally, depending on her birth history, breastfeeding status, medical history, blood pressure, and medicines she is taking at the time of the visit, some methods may be contra-indicated for the client.

An experienced family planning counselor, whether it is a medical doctor or a nurse – is trained to consider all these issues during a family planning counseling session. Since all of this information is difficult to consider and impart from memory, health providers often have job-aids, such as the medical eligibility criteria wheel or the quick reference chart of the WHO, the Balanced Counseling Strategy Plus Toolkit of the Population Council, checklists to be reasonably sure that the client is not pregnant, cue cards describing each contraceptive method, etc. These job aids are often available in counseling room in the form of posters and booklets in print form. Referring to these job aids during the counseling session is necessary, but often takes time and leaves room for provider error. The tablet-based job-support tool, into which all this information had been carefully programmed, was designed as a substitute for these charts, posters, toolkits, and reference books.[18]

The app serves to effectively structure the counseling session, and ensures that clients consistently receive high quality consultations. The app records the clients' answers to a series of questions eliciting their goals, fertility plans, needs, and preferences regarding contraceptive methods, as well as their medical and birth history, blood pressure, weight, and other relevant risk-factors. These questions will help the client, under the guidance of the provider, find a suitable a contraceptive method that is not contraindicated as per the assessment of the client's medical eligibility (more detail on the medical eligibility criteria in Appendix section A.3).

The structure of the counseling session as guided by the job-support tool is not fundamentally different than the standard practice –worldwide and at HGOPY. In other words, the tool does not require any new knowledge or training on part of the provider. It simply is a job-aid that allows her to conduct the counseling session more efficiently. The principles of counseling remain the same and the provider simply needs to familiarize herself with the tablet and the flow of the guided session to be able to take advantage of this tool in counseling clients. The process of family planning counseling using the job-support tool consists of three main sections, followed by a concluding section:

---

[18]Examples of similar smartphone-based job-support tools exist elsewhere. For example, in Nigeria, an initiative by the International Committee of the Red Cross, Swiss Tropical and Public Health Institute, and the Adamawa State Primary Health Care Agency developed a smartphone-based tool, called the Algorithm for the Management of Childhood Illnesses is an electronic upgraded version of the more commonly used IMCI (Integrated Management of Childhood Illnesses), which improves both preventive efforts and curative care for children under 5.

1. Introduction:

   (a) Welcome the client, explain the purpose of the session (talk about her life and goals, healthy families, pregnancy spacing, safe sex, and contraceptive methods), and clarify that the session is private and confidential.

   (b) Collect basic demographic information (age, marital status, education, primary activity, religion, and neighbourhood)

   (c) Discuss client's plans for having children in the future, how long she would like to wait before getting pregnant, how many more children she would like to have, and healthy birth spacing

   (d) Cover her birth history and establish her breastfeeding status

   (e) Conduct a pregnancy check

2. Consultation and needs assessment:

   (a) Discuss current method of birth control used by the client, if any. Discuss her experience with the method, how long she has been using it, and assess whether she would like to continue or switch

   (b) Discuss any methods that she might be worried about. Any methods she has in mind that she is curious about.

   (c) Clarify any questions or misconceptions the client might have about any contraceptive methods

   (d) Ask her about her preferences regarding side effects concerning:

       i. Increased bleeding and cramping,
       ii. Decreased bleeding, spotting, and amenorrhea, and
       iii. Weight gain

   (e) Obtain her medical history to avoid the adoption of contra-indicated methods. Take blood pressure and measure the height and weight of the client.

3. Method choice and follow-up:

   (a) Depending on the intervention condition, either ask the client to choose the modern method she would like to discuss first OR ask her whether she would like to discuss the method that is recommended by the tool as being the most suitable for her needs. The tool excludes methods that are contra-indicated from the list of methods presented to the client. When the client makes her choice as to which method she would like to discuss, the provider presents neutral, evidence-based, and understandable information on the effectiveness and the side effects of that method. The provider does so with the help of printed and laminated cue cards, each of which contain the necessary information for a different method.

   (b) Answer client's questions and concerns about the method being discussed. Listen to the client carefully and counsel her individually, based on her needs assessment.

   (c) Ask the client whether she would like to adopt the method or discuss another method. Discuss next preferred method, and so on, until the client decides to adopt or leave with no method.

(d) Adoption of chosen method, as appropriate, along with documentation of consent for adopting a modern method, such as the pill, injectable, implant, or IUD.

(e) Provide information on method use and follow-up mechanisms for switching or discontinuing selected method.

4. Conclusion:

(a) Discuss the importance of dual protection from sexually transmitted diseases and provide the client with a package of condoms.

(b) Schedule the next appointment, as appropriate.

(c) Provide the client information about the study and seek her informed consent to participate in the study.

It is important to reemphasize that the process of counseling a client for modern contraceptives with the help of the tablet-based job-support tool is very similar to what the provider would have done in the absence of the tablet.

## A.2 Consultation style: "sequential" vs. "side-by-side" views

The study aims to test the proposed paradigm shift in FP counseling – i.e. going from discussing all modern methods and letting the client state the method she would like to discuss first, to having the job-support tool recommending which method to discuss first based on the information elicited from the client during the session. The *method choice* happens after the providers have discussed the clients' past experiences with contraception and have elicited their fertility goals, their preferences towards side effects, and their medical eligibility – i.e. the elements necessary to make an informed choice for a contraceptive method to adopt. Keeping in mind that all clients receive a consultation using the app, the two different paradigms are compared by randomly varying the protocol providers are asked to follow when assisting clients during the method choice phase, and is reflected through the instructions and options displayed.

At the method choice section, the app randomly assigns each client to one of two regimes:

1. **Side-by-side recommendation:** In this regime, the job-support tool displays all available modern contraceptive methods that have not been ruled out by the client or contraindicated due to medical eligibility. The available modern methods are presented as unranked (i.e. as if each method is equally suitable for the client) and the providers will provide basic information on all available methods (in order of the methods displayed, which is randomized). The basic information covers what the method is, how it is used and its duration (capsule placed under the skin in the arm for 3-5 years or pill taken daily, etc.), and its typical use effectiveness. This quick description is expected to take about 30-60 seconds for each method. The provider will then ask the client to indicate which method they would like to discuss in more depth first. The provider will then use the relevant cue card to discuss the method in question in more detail to inform the client of all the relevant information the client needs to know before

adoption. After going through this information, the client can choose to either adopt this method or discuss another method (of her choice). This process is repeated until a decision is made.

2. **Sequential recommendation:** In this regime, the tablet will first display the method that is deemed most suitable for the client given her preferences, as described in section A.4, and ask them if they would like to hear about it. If the client answers 'no,' then the next highest ranked recommendation is displayed, and the provider asks the client if she would like to discuss this method, the process is repeated until the client decides to discuss one of the recommended methods (the app displays non-modern methods if the client does not want to discuss any one of the modern methods). If the client answers 'yes,' that they would like to hear about a modern method, then the provider uses the appropriate cue card to explain to the client the relevant information they need to know about the method in question. The client can then decide whether to adopt this method or discuss the next method recommended by the app. Again, this process is repeated until a decision is made.

## A.3   Medical eligibility

Medical eligibility and contra-indications

We use the U.S. Centre for Disease Control and Prevention's recommendations for "U.S. Medical Eligibility Criteria for Contraceptive Use, 2016" to determine methods that are contraindicated for various conditions the clients may have. The main considerations are the following:

- Recent delivery of a baby
- Breastfeeding
- Unexplained vaginal bleeding
- Current blood pressure/history of hypertension
- Risk factors such as older age (>35), smoking, diabetes
- Medications such as TB drugs, Barbiturates, and Antiretroviral drugs

A large number of medical eligibility rules relating to the conditions above are programmed into the algorithm and when a condition is satisfied, the method is ruled out and excluded from rankings. When this happens, the job-support tool displays a message at the top of the method choice section that certain methods are being excluded due to medical eligibility criteria, so that the provider can explain the client why she is not being given the option to discuss that method.

## A.4   Method rankings

The tablet-based job-support tool takes client preferences regarding how long to wait before becoming pregnant and the importance of various side effects into account, in order to produce method

rankings.

The algorithm uses three key criteria to rank the suitability of contraceptive methods for a client:

1. How long they would like to wait before becoming pregnant, and

2. How strongly they feel about avoiding the following three categories of side effects:

    (a) Increased bleeding and cramping,

    (b) Decreased bleeding, spotting, and amenorrhea, and

    (c) Weight gain

3. Typical use effectiveness of each method

Based on the client's answers to a number of questions eliciting their preferences, the algorithm creates a score for each method. In the "ranked recommendation" intervention condition, the highest ranked method (that is not eliminated by the algorithm due to contraindications or client wishes – please see below) is suggested to the client first, followed by the next highest ranked method, and so on until the client decides.

For example, if the client feels strongly about minimizing the chances of all three categories of side effects and would like to wait more than one year before getting pregnant, the ranking of methods (from most suitable to least suitable) is as follows: IUD, pill, (lactational amenorrhea method or LAM), implant, and the injectable. The LAM method is included in the rankings if a client has (i) given birth in the past six months; (ii) is fully breastfeeding; and (iii) has not menstruated since birth; and excluded otherwise. Please note that in this example, the pill, which is a short-acting method with a typical use effectiveness much lower than the implant and the injectable, is ranked above both methods because of the client's preferences regarding side effects. The algorithm uses evidence on the average side effects of each method from the existing peer-reviewed literature. Similarly, the typical use effectiveness data used by the algorithm comes from peer-reviewed literature.

In contrast, consider a client who wishes to have no more children and does not care about any of the side effects. The method rankings for such an individual is: IUD, implant, injectable=(LAM), and the pill. The reader will now notice that because the client is interested in avoiding pregnancies altogether and is not concerned with side effects, long-acting methods are ranked higher, while the pill, which has the mildest expected side effects, is ranked at the bottom. Please also note that the algorithm sometimes produces identical scores for two or more methods that result in a tie in the rankings. In such cases, the client is told that two (or more) methods are equally suitable for her and the tablet uses an internal random number generator to decide the ordering of the tied methods for discussion.

Conversely, in some cases, a client may have arrived at the facility almost certain of the method that she would like to adopt. For example, a client may have a friend or relative who is very satisfied with a certain method. Or, the client may have experience with a method in the past, say before her last pregnancy, and would like to return to using that method. To ascertain such cases, the provider asks the client in the "Consultation and needs assessment" phase, whether she

has any method in mind. If she specifies a method and indicates that she would not like to discuss any other methods, the specified method is moved to the top of the rankings and a default ranking based on typical use effectiveness is used for the remaining methods.

# B    Data-generating process

**Covariates** $X_t$    We estimate the marginal probability of covariates by fitting a joint normal distribution via maximum likelihood and then using inverse probability sampling to draw the covariates during the simulation. This procedure above ensures that the covariates are correlated but have the correct support and marginal distributions. In detail,

1. Using the pilot data, compute estimates of the mean $\hat{\mu}$, covariance matrix $\widehat{\Sigma}$ of the relevant contexts. Also, compute an estimate of their inverse marginal CDF $\hat{F}_j^{-1}(\cdot)$ for each relevant context $j$.

2. During the simulation, at each period $s$:

   (a) Draw the vector $Z_s \sim \mathcal{N}(\hat{\mu}, \widehat{\Sigma})$ (jointly).
   (b) Let $X_{s,j} = F_j^{-1}(Z_{s,j})$ for each $j$.

**Contraceptive adoption** $M_t$    We model the adoption probabilities via a multinomial logistic regression model that is fitted on a set of features of contexts and treatments. However, in order to ensure that we explore a larger extent of the set of possible outcome models, we "perturb" the model in two ways: by fitting on a bootstrap sample of pilot data, and by randomly penalizing less than suggested by cross-validation. That is, at the beginning of each simulation we go through the following steps:

1. Draw a sample with replacement of contexts, treatments and adopted contraceptives from pilot data $(X_t^*, W_t^*, M_t^*)_{t=1}^{n_{pilot}}$.

2. Construct a polynomial feature vector using contexts and treatments, as follows. First, we include a third-order polynomial of continuous variables (age, LARC prices), then we include the binary variables linearly ("method in mind", "own initiative", "view") , and finally include product interactions between the continuous and binary variables. The final feature vector is

$$\phi(x, w)$$
$$= [1, age, age^2, age^3, price, price^2, price^3,$$
$$method\_in\_mind, recent\_pregnancy, own\_initiative, view,$$
$$age \times method\_in\_mind,\ age^2 \times method\_in\_mind,\ age^3 \times method\_in\_mind,$$
$$price \times method\_in\_mind,\ price^2 \times method\_in\_mind,\ price^3 \times method\_in\_mind,$$
$$age \times recent\_pregnancy,\ age^2 \times recent\_pregnancy,\ age^3 \times recent\_pregnancy,$$
$$price \times recent\_pregnancy,\ price^2 \times recent\_pregnancy,\ price^3 \times recent\_pregnancy,$$
$$age \times own\_initiative,\ age^2 \times own\_initiative,\ age^3 \times own\_initiative,$$
$$price \times own\_initiative,\ price^2 \times own\_initiative,\ price^3 \times own\_initiative$$
$$age \times view,\ age^2 \times view,\ age^3 \times view,$$
$$price \times view,\ price^2 \times view,\ price^3 \times view]$$

(20)

3. For a grid of regularization parameters, fit a multinomial logistic regression of adopted contraceptives $M_t^*$ on constructed features $\phi(X_t^*, W_t^*)$ with L1-regularization. However, do not penalize the main effects of treatments (i.e., set their penalization to zero).[19] Compute the optimal regularization parameter via K-fold cross-validation.

4. Output the logistic model whose regularization parameter is closest to $r$ times the optimal parameter, where $r$ is drawn uniformly at random from $\{0.1, 0.5, 1\}$. This "under-penalization" allows us to explore more complex data-generating processes than observed in the pilot data.

**Other outcome variables**   Given contraceptive method $M_t$ above and the treatment $W_t$ that is revealed during the experiment, most other variables such as failures $F_t$ and cost $C_t$ are decided as explained in Section (2.1.5).

## C   Learning phase details

### C.1   Selecting the cost multiplier $\lambda$

Our selection of the parameter $\lambda = 1.3 \times 10^{-5}$ used in the definition of the loss (9) was motivated in large part by our objectives to increase the adoption of modern contraceptives while keeping costs low wherever possible, but was also informed by other competing concerns that were not easily accommodated in our framework.

---

[19]This estimation step the `glmnet` function from the `glmnet` R package. To estimate a multinomial logistic regression, we set `family` to `multinomial`. To avoid penalizing the treatment main effects, we set `penalty.factor` equal to 0 for coefficients associated with $price, price^2, price^3$ and $view$, and 1 otherwise. Other parameters are as in `glmnet` defaults. See https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html.

Recall from Section (2.2) that a higher value of $\lambda$ leads to policies that are better at minimizing costs, and therefore pay out fewer subsidies. We believe that the small drop in LARC uptake at around $\lambda = 1.3 \times 10^{-5}$ (relative to $\lambda = 0$) would be commensurate with the average decrease in costs. However, different values of $\lambda$ also lead to different partitions of the covariate space $\widehat{\varphi}$ (Section 2.2) and different initial assignment probabilities (Section 2.3). As we'll see below in Figure 12, even though alternative values of $\lambda$ in the neighborhood of our chosen $\lambda = 1.3 \times 10^{-5}$ yielded similar results in the aggregate, we selected this value because it not only yielded a sensible partition of the covariate space $\widehat{\varphi}$, but it also produced initial assignment probabilities that were "less aggressive" (i.e. lower) in terms of prices offered to adolescents. Our experience is that the latter concern (initial prices offered to adolescent females and young women at the start of the adaptive experiment) was important for the nurses conducting the counseling sessions.

Figure 12 shows how the performance of the learned policy changes as we vary $\lambda$, in terms of the average value of key outcomes. This is done for a policy learned using adaptively-collected data and for one learned by assigning arms uniformly at random during the "learning" phase (as in Section 3.3).

An important takeaway from Figure 12 is that the negative loss for a policy learned under adaptivity is at least as high (i.e. less bad) as the one for a policy learned without it. This implies that the former is able to uniformly improve upon the latter in terms of the loss (9) across all values of $\lambda$ tested here. Also, note that the graphs seem to jump discontinuously at certain points. This happens because the structure of the learned partition function $\widehat{\varphi}$ changes with $\lambda$ at these points, causing the learned policy to behave quite differently.



Figure 12: True value of the policy learned at the end of the "learning" phase, depending on the cost multiplier $\lambda$. Numbers are aggregated across over 10000 simulations. Errors bars are 95% confidence intervals around the average value.

## C.2 Constructing doubly-robust scores

Here we explain how to obtain doubly robust scores for several outcomes of interest. These will be used below when we described how to estimate the partition function $\widehat{\varphi}$ from (10) (Section C.3) and how to evaluate policies (Section C.4). Throughout we assume access to some dataset $S$.

**Adoption probabilities**   Recall from Section 2.1.5 that most outcomes are essentially deterministic functions of the adopted contraceptive $M_t$ and (possibly) assigned treatment. Therefore, our first step will be to estimate the probability that a client will adopt each contraceptive given their observable contexts and treatment These estimates will then be used in next steps to construct unbiased estimates of other outcomes of interest.

Because we will make use of a debiasing step that will be shown shortly, we do not require that these estimated probabilities be well-specified. In fact they could all be set to a constant, but better predictions will translate into smaller variances and better policies. We call these "plugin" estimates. However, in order to preserve desirable statistical properties of the data, when predicting the adoption probabilities for the $t^{th}$ client, we can only use observations up to time $t-1$. This is accomplished as follows.

Fix a time grid of timepoints $\{\tau_j\}_{j=1}^{N_\tau}$ with $\tau_1 = 1$ and $\tau_{N_\tau} = |S|$. Let $\hat{p}(m \mid x, w; H^{s-1})$ denote the estimated probability that a client will choose contraceptive $m$ given their observed characteristics $x$ and an arm $w$, fitted using the past data $H^{s-1}$. Use this time grid to split the current data into "sections". We use $N_\tau = 20$ throughout.

Next, at each timepoint $\tau_j$, estimate $\hat{p}(m \mid x, w; H^{\tau_j - 1})$ and use this model to predict adoption probabilities for clients with arrival times in the next "section", between $\tau_j$ and $\tau_{j+1} - 1$. The probability model is estimated by fitting a regularized multinomial logistic regression of $M_t$ on contexts $X_t$, treatments $W_t$ and their multiplicative interactions, without penalizing the coefficients associated with treatment main effects. For observations in the first section, default to $\hat{p}(m \mid X_t, w; H^0) \equiv 1/5$ for all methods $m$.

**Failure and cost**   Using the adoption probabilities estimated in the previous step, we compute implied estimates of conditional expected failure and cost given contexts and arms. As shown in equation (21), we compute these quantities by averaging across contraceptive methods,

$$
\begin{aligned}
\hat{\mu}_{t,w}^{\text{fail}} &:= \sum_{\text{Contraceptive } m} \hat{p}(m \mid X_t, w; H^{\tau(t)-1}) \cdot \text{failure}(m) \\
\hat{\mu}_{t,w}^{\text{cost}} &:= \sum_{\text{Contraceptive } m} \hat{p}(m \mid X_t, w; H^{\tau(t)-1}) \cdot \text{cost}(w, m),
\end{aligned}
\tag{21}
$$

where, somewhat overloading notation, $\tau(t)$ refers to the "section" to which observation $t$ belongs, and failure and cost are deterministic mappings as described in (3) and (4). For example, if $m = \text{pill}$, then $\text{failure}(m) = 0.06$.

The expressions (21) may be biased due to model misspecification. Therefore, we augment each estimate with a debiasing term. The resulting object is known as an augmented inverse propensity-weighted (AIPW) score. Letting $F_t$ and $C_t$ denote the actual (observed) failure rates and disbursed costs,

$$
\begin{aligned}
\widehat{\Gamma}_{t,w}^{\text{fail}} &:= \hat{\mu}_{t,w}^{\text{fail}} + \frac{\mathbb{I}\{W_t = w\}}{e_{t,w}} \left( F_t - \hat{\mu}_{t,w}^{\text{fail}} \right) \\
\widehat{\Gamma}_{t,w}^{\text{cost}} &:= \hat{\mu}_{t,w}^{\text{cost}} + \frac{\mathbb{I}\{W_t = w\}}{e_{t,w}} \left( C_t - \hat{\mu}_{t,w}^{\text{cost}} \right),
\end{aligned}
\tag{22}
$$

where $e_{t,w}$ are the propensity scores, or the assignment probabilities of arm $w$ to client $t$. The additive term on the right ensures that unconditional expectations are unbiased; i.e., $\mathbb{E}\left[\widehat{\Gamma}_{t,w}^{\text{fail}}\right] = \mathbb{E}\left[\text{failure}(M_i)\right]$ and $\mathbb{E}\left[\widehat{\Gamma}_{t,w}^{\text{cost}}\right] = \mathbb{E}\left[\text{cost}(w, M_i)\right]$. For more on the role of AIPW scores in policy estimation, see Athey and Wager (2021).

**Loss**  Doubly-robust scores for the loss (9) can be readily computed from (22),

$$
\widehat{\Gamma}_{t,w}^{\text{loss}} := \alpha_t \cdot \widehat{\Gamma}_{t,w}^{\text{fail}} + \lambda \cdot \widehat{\Gamma}_{t,w}^{\text{cost}}.
\tag{23}
$$

**Pregnancies**  To compute doubly-robust scores for pregnancies, first estimate the probability of pregnancy given contexts and treatments,

$$
\hat{\mu}_{t,w}^{\text{preg}} = \widehat{\mathbb{P}}\left[ G_t = 1 \,\big|\, X_t, W_t = w, H^{\tau(t)-1} \right].
\tag{24}
$$

where again the model is estimated via L2-regularized logistic regression of pregnancy outcomes $G_t$ on contexts $X_t$, treatments $W_t$ and their interactions, without penalizing coefficients associated with treatment main effects, and using only past data. Next, compute

$$
\widehat{\Gamma}_{t,w}^{\text{preg}} := \hat{\mu}_{t,w}^{\text{preg}} + \frac{\mathbb{I}\{W_t = w\}}{e_{t,w}} \left( G_t - \hat{\mu}_{t,w}^{\text{preg}} \right).
\tag{25}
$$

## C.3  Partition function

To obtain the partitioning function $\widehat{\varphi}$ defined in (10), first construct doubly-robust scores (22) as explained in Section C.2 using the subset $S_1^{pilot}$ obtained by splitting the pilot sample. Next, estimate a decision tree $\varphi$ that minimizes the average loss (23),

$$
\widehat{\varphi} := \arg\min_{\varphi \in \Phi(d)} \frac{1}{|S_1^{\text{pilot}}|} \sum_{t \in S_1^{\text{pilot}}} \widehat{\Gamma}_{t,\varphi(\tilde{X}_t)}^{\text{loss}}
\tag{26}
$$

The solution to (26) is found via the R package `policytree` (Sverdrup et al., 2020) based on Zhou et al. (2018). The maximal depth of the tree is set to $d = 2$. Recall from 2.1.3 that $\tilde{X}_t$ are contexts that do not include "method in mind".

## C.4   Policy evaluation

In order to evaluate any policy, use the entire dataset (from pilot and main experiment) to construct doubly-robust scores (22) as explained in Section C.2. Then, for example, the average failures induced by any policy $\pi$ that does not depend on data collected during the "evaluation" phase can be estimated in an unbiased manner as

$$\widehat{\mathbb{E}}\left[F_t(\pi(X_t))\right] := \frac{1}{|S^{\text{eval}}|} \sum_{t \in S^{\text{eval}}} \widehat{\Gamma}^{\text{fail}}_{t,\pi(X_t)}, \tag{27}$$

where $S^{\text{eval}}$ are the observations collected during the "evaluation" phase. The average performance of the policy $\pi$ with respect to other outcomes is estimated in a similar manner, replacing the failure doubly-robust scores in (27) for the ones associated with the outcome of interest.

The only exception to the above is the pregnancy outcome, which is different in two ways. First, we do not collect data on pregnancies for observations in the pilot, so doubly-robust scores for pregnancies can only be obtained for observations in the main experiment. Second, even during the main experiment we only observe observations where there was no attrition, so evaluation is done on the possibly smaller subset $\tilde{S}^{\text{eval}}$ of observations in the "evaluation" for which we do have data on pregnancies.

## C.5   Assignment probabilities

Let's describe in detail how probabilities are computed at the beginning of a "learning" phase batch $m$. To simplify notation, let $S_m$ denote all the available data including the pilot data subset $S_2^{\text{pilot}}$ that was not used to estimate the partition function $\widehat{\varphi}$, as well as previous data from the learning phase. Also, let

$$I_{j,w} = \{t \in \mathbb{N} : t \in S_m, \widehat{\varphi}(X_t) = j, \text{ and } w \in W_t\}, \tag{28}$$

that is, the indices for available observations whose covariates fall in estimated region $j$ and were assigned to arm $w$.

**Thompson Sampling probabilities**   Given a Bayesian model of the data-generating process, we call the posterior probability that arm $w$ is the best arm its "Thompson Sampling" (TS) probability. Several bandit algorithms (see e.g., Russo et al., 2018) use these probabilities to assign arms in adaptive experiments whose objective is to minimize the mistakes during the experiment. Here, our goal is not only to minimize the number of mistakes *during* the experiment, but also to learn a good policy at the *end* of the experiment, so we will have to modify these probabilities later to suit our objectives. But first, let's see how TS probabilities are computed.

We begin by estimating the mean and standard deviation of the loss associated with each arm, in

each region. Letting $loss_t := \alpha_t F_t + \lambda C_t$,

$$n_{j,w} = |I_{j,w}|$$
$$\hat{\mu}_{j,w} = \frac{1}{n_{j,w}} \sum_{t \in I_{j,w}} loss_t \tag{29}$$
$$\hat{v}_{j,w} = \frac{1}{n_{j,w}} \sum_{t \in I_{j,w}} (loss_t - \hat{\mu}_{j,w})^2$$

Next, assuming an uninformative prior and a Normal likelihood, we approximate the posterior probability that arm $w$ is optimal in region $j$ as follows. We first draw

$$Z_{j,w,k} \sim \mathcal{N}(\hat{\mu}_{j,w}, \hat{v}_{j,w}/n_{j,w}) \qquad for\, k \in \{1, \cdots, N\},$$

where $N$ is a large number, and then approximate the TS probabilities via

$$\tilde{e}_{m,j}(w) = \frac{1}{N} \sum_{k=1}^{N} \mathbb{I}\{w = \arg\min_{\tilde{w}} Z_{j,\tilde{w},k}\}. \tag{30}$$

Finally, recall from Section 2.1.4 that clients who have a method in mind cannot be assigned to the "side-by-side" (SBS) view, so for those clients we must set the probabilities of arms that assign SBS to zero and renormalize the vector of probabilities so they still sum to one. In other words, in each region we effectively compute two Thompson Sampling probability vectors. For clients who do not have a method in mind, we use

$$\tilde{e}_{t,j}(w) \equiv \tilde{e}_{m,j}(w) \text{ for every } w. \tag{31}$$

For clients who do have a method in mind, we use

$$\tilde{e}_{t,j}(w) = \begin{cases} 0 & \text{if w assigns SBS,} \\ \tilde{e}_{m,j}(w) \Big/ \sum_{w \text{ assigning SEQ}} \tilde{e}_{m,j}(w) & \text{otherwise.} \end{cases} \tag{32}$$

And these TS probabilities are used for all clients $t$ in this batch.

**Exploration sampling probabilities**   As mentioned at the top of the section, assigning arms using Thompson Sampling probabilities is a common heuristic for minimizing mistakes during the experiment. However, our objectives also include finding a good policy at the end of the learning phase. This objective often requires more exploration than Thompson Sampling, so here we modify the TS probabilities following Kasy and Sautmann (2021):

$$e_{t,j}(w) = \frac{\tilde{e}_{t,j}(w)(1 - \tilde{e}_{t,j}(w))}{\sum_w \tilde{e}_{t,j}(w)(1 - \tilde{e}_{t,j}(w))} \tag{33}$$

The transformation (33) has the effect of shrinking positive probabilities towards the uniform probability. Note in particular that no arm has probability larger than one half, although since the transformation is monotonic it is still the case that more promising arms (i.e., arms with high TS probability) are assigned more often.

**Imposing an assignment floor**  Finally, although this may deviate from our objectives, even during the learning phase we impose a lower bound, or "floor," on assignment probabilities. This is done as a guardrail against potential issues during the experiment: having positive probability of assignment for all arms during this phase allows us the possibility of reusing this data if we have need to (using e.g., the techniques described in Hadad et al., 2019).

A floor $\bar{p}_0^{\text{all}}/K_t$ can be imposed as follows. First, for every arm available $w$ that for which $e_{t,j}(w) < \bar{p}_0^{\text{all}}$, set $e_{t,j}(w) \leftarrow \bar{p}_0^{\text{all}}$. Then, decrease the assignment probability of all other arms by setting $e_{t,j}(w) \leftarrow \bar{p}_0^{\text{all}} + C(e_{t,j}(w) - \bar{p}_0^{\text{all}})$, where $C > 0$ is some constant that ensures that sum of all assignment probabilities is one. This process is shown in Figure 13.



Figure 13: Imposing a lower bound on exploration sampling probabilities.

# D    Additional tables

The following tables show statistics computed using the second half of pilot data (i.e., the portion of the data denoted by $S_2^{\text{pilot}}$ in Sections 2.3 and C.5). These are relevant to the computation of probabilities used in the first batch of the "learning" phase.

Table 14 shows the number of observations falling within each data cell defined by region and arm assignment. Tables 15 and 16 show the mean and standard error of negative loss (9) computed in each cell.[20]  Finally, Tables 17 and 18 shows the assignment probabilities in the first batch of the learning phase, for clients who have a method in mind and who do not.

---

[20]The zeros in Table 16 are not a typo. For the cell in row (1000, SBS), column B, there was only one observation. For the cells in the last row, all 8 clients (3 in column C and 5 in column D) ended up choosing no contraceptive method, hence their negative loss is -.25 corresponding to a failure of 25% and no cost. Note that, as Tables 17 and 18 show, this does not preclude these arms from being selected.

|           | A   | B   | C   | D   |
|-----------|-----|-----|-----|-----|
| (0, SEQ)   | 33 | 35 | 40 | 30 |
| (0, SBS)   | 7  | 10 | 9  | 7  |
| (1000, SEQ) | 19 | 24 | 27 | 10 |
| (1000, SBS) | 8  | 1  | 5  | 6  |
| (2000, SEQ) | 16 | 19 | 17 | 12 |
| (2000, SBS) | 6  | 4  | 8  | 8  |
| (4000, SEQ) | 22 | 11 | 22 | 9  |
| (4000, SBS) | 6  | 4  | 3  | 5  |

Table 14: Number of observations in each region and arm.

|           | A       | B        | C        | D       |
|-----------|---------|----------|----------|---------|
| (0, SEQ)   | -0.131  | -0.1293  | -0.09394 | -0.1383 |
| (0, SBS)   | -0.138  | -0.1729  | -0.1848  | -0.25   |
| (1000, SEQ) | -0.1746 | -0.1194  | -0.09678 | -0.1455 |
| (1000, SBS) | -0.1847 | -0.25    | -0.1714  | -0.2149 |
| (2000, SEQ) | -0.1975 | -0.09826 | -0.08041 | -0.25   |
| (2000, SBS) | -0.1572 | -0.1401  | -0.1434  | -0.203  |
| (4000, SEQ) | -0.1773 | -0.1606  | -0.08532 | -0.2223 |
| (4000, SBS) | -0.2917 | -0.1876  | -0.25    | -0.25   |

Table 15: Average negative loss per region and arm, estimated using the second half of pilot data.

|           | A       | B       | C       | D       |
|-----------|---------|---------|---------|---------|
| (0, SEQ)   | 0.01937 | 0.01598 | 0.01251 | 0.01814 |
| (0, SBS)   | 0.0396  | 0.0315  | 0.03262 | 0.0     |
| (1000, SEQ) | 0.03642 | 0.0211  | 0.0178  | 0.03484 |
| (1000, SBS) | 0.04226 | 0.0     | 0.04818 | 0.03508 |
| (2000, SEQ) | 0.03874 | 0.02431 | 0.02353 | 0.0     |
| (2000, SBS) | 0.05043 | 0.06345 | 0.04047 | 0.03122 |
| (4000, SEQ) | 0.03277 | 0.03739 | 0.02487 | 0.02772 |
| (4000, SBS) | 0.04167 | 0.06237 | 0.0     | 0.0     |

Table 16: Standard error of negative loss per region and arm, estimated using the second half of pilot data.

|          | A     | B     | C     | D     |
|----------|-------|-------|-------|-------|
| (0, SEQ) | 0.276 | 0.057 | 0.117 | 0.425 |
| (0, SBS) | 0.273 | 0.025 | 0.025 | 0.025 |
| (1000, SEQ) | 0.070 | 0.172 | 0.136 | 0.417 |
| (1000, SBS) | 0.070 | 0.025 | 0.027 | 0.025 |
| (2000, SEQ) | 0.033 | 0.351 | 0.327 | 0.025 |
| (2000, SBS) | 0.207 | 0.240 | 0.048 | 0.033 |
| (4000, SEQ) | 0.046 | 0.055 | 0.296 | 0.025 |
| (4000, SBS) | 0.025 | 0.075 | 0.025 | 0.025 |

Table 17: Assignment probabilities during first batch of learning phase for clients who do not have a method in mind.

|          | A     | B     | C     | D     |
|----------|-------|-------|-------|-------|
| (0, SEQ) | 0.429 | 0.103 | 0.129 | 0.452 |
| (1000, SEQ) | 0.186 | 0.264 | 0.172 | 0.448 |
| (2000, SEQ) | 0.321 | 0.444 | 0.370 | 0.050 |
| (4000, SEQ) | 0.064 | 0.189 | 0.328 | 0.050 |

Table 18: Assignment probabilities during first batch of learning phase for clients who do have a method in mind.