

InfoCollection Examples

Information collection problems are incredibly common in application. We give examples, classifying them according to the type of the information collection decisions involved. The difficulty of the problem depends in part on the size of the decision space, but also on the complexity of the underlying statistical model.

Stopping decisions

In this first set of problems, the information collection decision is a stopping time that adaptively chooses how much information to collect.

1. **Early warning systems:** When using radar during wartime, given that we observe a suspicious signal, we would like to decide between raising an alarm or waiting and collecting more information through observation over a longer period of time. See [21].
2. **Statistical process control:** When operating an industrial chemical process and periodically observing the quality of samples pulled from the production line, we would like to stop the process if poor sample quality indicates a production problem. We would like to stop the process as soon as possible when there is a problem, while waiting long enough to control the probability of false alarm under normal fluctuations. See [22].
3. **Group sequential clinical trials:** We would like to adaptively choose the length of time to run a clinical trial. Adaptively choosing the duration allows us to stop the trial early and immediately make the treatment available to the general public if the treatment being tested performs very well on the initial set of patients, saving lives. If the initial results are equivocal we can run the trial longer instead, and if the initial results are very poor we can stop the trial early and declare the drug a failure, saving the cost of a longer trial. See [8].
4. **Drug recall:** While monitoring the incidence of side effects from a drug approved for use in the population, upon observing a few suspicious events we would like to decide whether to recall the drug. By waiting, we collect more information about the incidence of side effects, allowing

us to know with greater accuracy whether the incidents observed were due to chance or indicate a real problem, but we also risk more danger to the population.

5. **Technology adoption:** When operating a business, we would like to decide whether to adopt a new technology (such as a new computer system or manufacturing technology) now, or to wait and gather more information about how well it works by observing the experience of other early-adopters.

Decisions from a finite set

In the following information collection problems, the information collection decision at each time period comes from a finite set.

1. **Screening of shipping containers:** We would like to screen shipping containers entering a port for dangerous nuclear materials. We have a number of examination techniques at our disposal, including various types of passive imaging (e.g., observation of gamma rays) and active interrogation (e.g., with a neutron beam), as well as manual inspection of containers. Each examination technique has its own costs and unique detection abilities. We would like to design a protocol to guide our application of these techniques to maximize the probability of detecting dangerous containers while minimizing cost and delay of containers traveling through the port. See [10].
2. **Outsourcing:** When making the decision of whether or not to outsource a project to another firm, we would like to account for the fact that, by keeping the project in-house, employees in our firm will learn and develop expertise that will be useful in future projects. See [1].
3. **Manufacturing optimization:** We would like to choose the best assembly line configuration from among several related choices to maximize manufacturing efficiency. We have the ability to collect information about the quality of any given configuration from a discrete event simulation of the manufacturing operation, or from using the configuration in practice.
4. **Drug screening:** Upon having identified a target pathway important in the function of a disease, e.g., a protein interaction, we would like

to screen a large number of unrelated compounds (usually between 10^3 and 10^6) using automated high-throughput in vitro experiments in order to find a relatively small number of compounds with the ability to disrupt this pathway.

5. **Collaborative filtering:** We would like to use collaborative filtering to suggest books, movies, or music to users according to their underlying but unknown preferences. After viewing an item, the user reports his or her enjoyment of that item with some probability. We may suggest an item because it possesses a high probability of being enjoyed by the user, but we may also suggest an item to learn more about the user's preferences.
6. **Emergency vehicle routing:** We would like route ambulances from their origin to patients and then to the hospital in a minimal amount of time. Before an ambulance leaves, or while it is in the early part of its route, we have the option of asking a traffic helicopter to observe traffic flows along links in the transportation network to choose a route with the smallest possible travel time. See [15].

Continuous univariate decisions

In the following information collection problems, the information collection decision is chosen from a one-dimensional continuous or discrete space. Because the space has only a single dimension, it can generally be discretized if it was continuous and then truncated to approximate the problem with a finite decision set of manageable size.

1. **Product pricing:** We would like to dynamically vary the price of a product to maximize revenue while learning about demand at different prices. Similarly, we might like to vary the price in test and target markets in order to learn about customer demand and better choose a final market price at which to sell the product.
2. **Venture capital:** We learn about the profitability of startup companies with various characteristics by observing how other companies have fared in the past. We may learn from startup companies funded by other firms, but this may not be possible if our firm provides a large

fraction of the funding in our sector, or if the information about companies funded by other firms is delayed in becoming available. At each point in time we must decide how much money to allocate to each company requesting funding, in order to maximize both short-term profit and our information about the marketplace, which will increase future profits.

3. **Sequential planning:** We would like to discover whether a new experimental technique for growing samples in a biology lab produces better results than a standard technique. We have the option of growing samples under the new and standard techniques in batches. There is a large cost paid per batch due to the time required to grow the samples, and a smaller cost paid per sample included due to the experimentalist's time preparing it. The decision at the completion of each batch is whether to continue with another batch, and if so, how many samples to include. See [18, 17].
4. **Supply chain management:** We would like to make inventory decisions in a supply chain, for which the distribution of customer demand is unknown and changing over time. Our observations of demand come only through the amount of product we sell, and demand lost due to stockout is unobserved. By stocking more product, we collect more information about demand because stockouts are less frequent, but we may pay more in holding costs or unused product. See [9, 12].
5. **Fishery regulation:** We would like to set limits on fishing low enough to maintain the fish population, without unnecessarily reducing profits in the fishing industry. The current fish population is only partially known, as are the dynamics governing it, and one of the main methods with which to learn about them is through the returns from the catch. Each year we set limits on the catch, keeping in mind that we will have less information later about the fish population if we set the limits low now. See [20].
6. **Drug dosing:** In clinical practice, we would like to escalate the dosage of a drug to a therapeutic level as quickly as possible without causing undue side effects. Similarly, in a clinical trial, we would like to find the maximum tolerated dose of a chemotherapy drug as quickly as possible, while minimizing toxic effects to participants. See [5, 2].

Continuous multivariate decisions

In the following information collection problems, the information collection decision is chosen from a multidimensional continuous space. The dimension of the space is often large enough that discretization is computationally infeasible.

1. **Oil exploration:** We would like to discover the best place at which to drill a commercial oil well by drilling a sequence of exploratory test wells. We would like to find a good location with as few exploratory wells as possible. See [3].
2. **Chemical process optimization:** We would like to choose the inputs to an industrial chemical process to maximize the quality of the output. For example, we might be choosing the temperature, pressure, and mixture of gases to use when etching silicon wafers. See [11].
3. **Visual search in computer vision:** We would like to design a computer vision algorithm that adaptively chooses where to point a video camera in order to quickly and effectively find a particular object. See [14].
4. **Physical search:** We would like to direct one or more search teams in order to find a lost person, e.g., a hiker, or a lost object, e.g., a sunken submarine. See [19, 4].
5. **Emergency response:** Following the release of radioactive material in a major city, we would like to identify contaminated areas using an efficient sequential method for testing contamination at points within the city. At each point in time our decision contains the areas to which our inspectors will travel and investigate next.
6. **X-ray crystallography:** We would like to crystallize a protein in order to interrogate its 3-dimensional structure using X-ray crystallography. Many proteins crystallize only under a delicate and unique set of experimental conditions. Given a particular protein, we would like to search for appropriate experimental conditions, trying these conditions sequentially until we find one that works.
7. **Model calibration:** We have a time-consuming computational model for which we would like to find parameters that best fit observed data.

We would like to search through the space of input parameters to the computational model, running the model on parameters of our choosing and learning how well their outputs fit data. See [6].

Decisions with combinatorial structure

In the following information collection problems, each information collection decision has combinatorial structure, often because it is drawn from the set of subsets of some larger finite set.

1. **Drug combination therapy:** Given a particular patient and his or her medical characteristics, a doctor would like to decide what drug or combination of drugs to prescribe. See [16] for an example motivated by diabetes treatment.
2. **Drug design:** We would like to search among chemically similar derivatives of a molecule to find the one that is best at binding to a target protein. These molecules are specified by the chemical functional groups chosen from a larger set and placed at fixed substituent locations. At each point in time, we decide which molecule to test, and then collect information about its effectiveness.
3. **Product feature selection:** We would like to choose the best set of features to include in a product, e.g., a cell phone, to maximize profitability. To determine consumer preferences for features we may conduct a sequence of focus groups in which we ask participants to judge the quality of feature subsets of our choosing.
4. **Assortment planning:** We would like to choose the best composition of products in a product line to simultaneously satisfy demand from several market segments. We may learn about demand from focus groups or from purchases of products currently offered in the product line.
5. **Pricing of credit default swaps:** We have a mathematical model of credit default swap prices parameterized by a clustering of firms into categories [13], under which prices of liquidly traded assets are predicted by taking an expectation via many replications of a Monte Carlo simulation. In order to price more exotic credit derivatives, we

would like to calibrate the model by finding a clustering of firms that accurately replicates market prices.

6. **Research and development portfolio selection:** We would like to allocate research funding among the set of proposed research projects to maximize the expected value to society of the resulting research. See [7].
7. **Network travel times:** We would like to find the path with the shortest expected travel time through a transportation or communications network by measuring travel times along paths of our choosing. We might have the ability to observe individual travel times along links in the path, or perhaps only the total travel time.

References

- [1] E.G. Anderson Jr and G.G. Parker. The Effect of Learning on the Make/Buy Decision. *Production and Operations Management*, 11(3):313–339, 2002.
- [2] J. Babb, A. Rogatko, and S. Zacks. Cancer phase i clinical trials: efficient dose escalation with overdose control. *Stat Med*, 17(10):1103–1120, May 1998.
- [3] J.E. Bickel and J.E. Smith. Optimal sequential exploration: A binary learning model. *Decision Analysis*, 14(15):16, 2006.
- [4] G. Chudnovsky. *Search theory: some recent developments*. CRC Press, 1988.
- [5] B. H. Eichhorn and S. Zacks. Sequential search of an optimal dosage. I. *J. Amer. Statist. Assoc.*, 68:594–598, 1973.
- [6] P. Frazier and W.B. Powell. Simulation model calibration with correlated knowledge-gradients. in review, 2009.
- [7] L. Hannah, W. B. Powell, and J. Stewart. One-stage r&d portfolio optimization with an application to solid oxide fuel cells. in review, 2009.

- [8] C. Jennison and B.W. Turnbull. *Group sequential methods with applications to clinical trials*. CRC Press, 2000.
- [9] M.A. Lariviere and E.L. Porteus. Stalking Information: Bayesian Inventory Management with Unobserved Lost Sales. *Management Science*, 45(3):346–363, 1999.
- [10] D. Madigan, S. Mittal, and F. Roberts. Sequential decision making algorithms for port of entry inspection: Overcoming computational challenges. In *IEEE International Conference on Intelligence and Security Informatics (ISI-2007)*, pages 1–7, 2007.
- [11] R.H. Myers and D.C. Montgomery. *Response Surface Methodology: Process and Product Optimization Using Designed Experiments*. John Wiley & Sons, New York, 2002.
- [12] D. Negoescu, P. Frazier, and W. B. Powell. Optimal learning policies for the newsvendor problem with censored demand and unobservable lost sales. in preparation, 2008.
- [13] E. Papageorgiou and R. Sircar. Multiscale intensity models and name grouping for valuation of multi-name credit derivatives. *Applied Mathematical Finance*, 2009. to appear.
- [14] L.W. Renninger, P. Verghese, and J. Coughlan. Where to look next? eye movements reduce local uncertainty. *Journal of Vision*, 7:1–17, 2007.
- [15] I. Ryzhov and W. B. Powell. Information collection on a graph. in review, 2009.
- [16] I. Ryzhov, W. B. Powell, and P. Frazier. The knowledge-gradient algorithm for a general class of online learning problems. in review, 2008.
- [17] C. Schmegner and M.I. Baron. Principles of optimal sequential planning. *Sequential Analysis*, 23(1):11–32, 2004.
- [18] N. Schmitz. *Optimal sequentially planned decision procedures*. Springer-Verlag Berlin, 1993.
- [19] L.D. Stone. *Theory of Optimal Search*. Academic Press, 1975.

- [20] D. Tomberlin. An approach to managing fisheries when weak and strong stocks mix. In *Proceedings of 2008 International Institute of Fisheries Economics and Trade*, 2008.
- [21] A. Wald and J. Wolfowitz. Optimum Character of the Sequential Probability Ratio Test. *The Annals of Mathematical Statistics*, 19(3):326–339, 1948.
- [22] G.B. Wetherill and D.W. Brown. *Statistical process control: theory and practice*. Chapman & Hall, 1991.