

---

# Probabilistic Computation for Information Security

---

**Piotr Mardziel**

Department of Computer Science  
University of Maryland, College Park  
College Park, MD 20740  
piotrm@cs.umd.edu

**Kasturi Rangan Raghavan**

Department of Computer Science  
University of California, Los Angeles  
Los Angeles, CA 90024  
kr@cs.ucla.edu

Probabilistic computation is a convenient means of mechanically reasoning about a variety of information security problems. At its core, information security concerns itself with measuring or limiting the knowledge an adversary might attain from interacting with a protected system. Probabilistic inference lets one compute this knowledge explicitly as long as the interaction can be described as a program in a suitable probabilistic language that supports conditioning. Security concerns, however, require soundness guarantees on probabilistic inference which are generally not present in machine learning applications. We summarize some recent work on probabilistic computing for information security and highlight challenging aspects that still need to be addressed.

**Tracking knowledge** Tracking adversary knowledge is a central tool for the applications to follow. We describe this core concept adopting the notation and conventions of Clarkson’s experimental protocol [2].

Bob seeks the answer of a query, the output of a function evaluated on Alice’s data. The data here is secret in that Bob does not already know it. Alice wants to protect her secret, making sure that answering Bob’s query is safe. To do this, Alice will explicitly track what Bob learns about her secret, treating Bob as a rational adversary who will use any information provided to him in order to learn Alice’s secret. Specifying the protection of the secret will be addressed shortly.

Alice’s secret data is a value  $s^*$  which Bob knows is in some set of possibilities, **Secrets**. Furthermore, Bob considers some values more likely than others, a probability distribution  $\delta_1 : \mathbf{Secrets} \rightarrow \mathbb{R}$ . We will refer to this distribution as Bob’s *belief*. Alice then considers the effect of letting Bob learn the output of a potentially non-deterministic query function  $Q$ , written in a probabilistic language. She can do this by computing the conditional distribution  $\delta_2 \stackrel{\text{def}}{=} \delta_1 | (Q(s^*) = o)$  for some potential output value  $o$ . Assuming  $\delta | C$  indeed implements *exact* probabilistic conditioning, this distribution  $\delta_2$  will be the belief Bob would attain, were he to see that the function  $Q$  outputs  $o$  given input  $s^*$ .

If Bob indeed does learn this output, Alice can revise her representation of Bob’s belief to  $\delta_2$  which she can then revise further if Bob learns additional functions over her data. In this manner Alice explicitly tracks what Bob knows about her secret as a probability distribution. The problems of how Alice knew  $\delta_1$ , Bob’s initial belief, in the first place and whether she can actually perform exact inference will need to be addressed.

**Knowledge-based policy enforcement** Alice can use the tracked belief in a variety of ways, the most direct of which is to judge the information security or safety of  $Q$  in order to determine whether she should let Bob get access to the output of the function. To do this, Alice starts by defining a knowledge-based policy  $P : \mathbf{Dist} \rightarrow \{\text{True}, \text{False}\}$ , where **Dist** is a set of all distributions over **Secrets**. For example, Alice might require that Bob never becomes too certain of her value, so she defines  $P(\delta) \stackrel{\text{def}}{=} \forall s \in \mathbf{Secrets} . \delta(s) \leq t$  for some threshold  $t$ . She will then deem the query  $Q$  safe whenever  $P(\delta_2)$  holds. This particular policy is an application of Smith’s vulnerability metric [6], which measures the expected likelihood of Bob guessing a secret value correctly in one try.

Unfortunately for Alice, exact probabilistic conditioning is expensive so she decides has to use an approximation. Instead of computing  $\delta_2$ , she gets  $\delta_2^\sim \stackrel{\text{def}}{=} \delta_1 | \sim (Q(s^*) = o) \approx \delta_1 | (Q(s^*) = o)$ . The

approximation  $\delta_2^\sim$  here need not be a proper probability distribution. She does not want the approximate inference to conclude a query is safe when it really is not so she requires the approximation to be sound relative to her policy:

**Definition.** A probabilistic inference method  $|\sim$  is *sound relative to a policy*  $P$  iff for every belief  $\delta$ , if  $P(\delta|C)$  fails then  $P(\delta|\sim C)$  fails.

For Alice’s example policy, a sound inference method would have to never underestimate the probability of any secret, though it could over-approximate it. In recent work [5], we demonstrate a probabilistic language capable of this approximate, but sound inference. Given the definition of conditional probability,  $\Pr(A|B) = \Pr(A \wedge B) / \Pr(B)$ , the probabilistic interpreter maintains an over-approximation of  $\Pr(A \wedge B)$  and an under-approximation of  $\Pr(B)$ , leading to a sound upper bound for  $\Pr(A|B)$ . The work shows how Alice can enforce her policy without revealing anything about her secret to Bob. Further application of these ideas were proposed to protect information of multiple mutually distrusting parties participating in secure multi-party computation [4].

If Alice deems  $Q$  safe, she can allow Bob access to  $Q(s^*)$  and update her representation of Bob’s belief to  $\delta_2^\sim$ . If she rejects  $Q$ , she will not allow access to  $Q(s^*)$  and keep  $\delta_1$  as the representation of Bob’s unchanged belief. Whether Alice deems  $Q$  safe or not, she can repeat the same process for another program that might come around, starting from either  $\delta_1$  or  $\delta_2^\sim$ , thereby maintaining an explicit representation of Bob’s knowledge while enforcing her knowledge-based policy.

**Obfuscation/Noising** Probabilistic computation is also very convenient for discerning the effects of obfuscation or noising mechanisms, via simple composition with the query function  $Q$ .

If Alice does not want Bob to learn the output of  $Q$  directly, she can add obfuscation to its input, or noising to its output. Let  $O$  be a potentially non-deterministic obfuscation function, written in a probabilistic language, that takes in a value  $s \in \mathbf{Secrets}$  and produces another  $s' \in \mathbf{Secrets}$ . Depending on the structure of  $\mathbf{Secrets}$ , this obfuscation might be in the form of resolution reduction, permutation of locations, or a variety of other mechanisms. Alice can then consider Bob’s knowledge were he to learn the output of an obfuscated function  $Q(O(s^*))$  instead of  $Q(s^*)$  directly, using conditioning,  $\delta_2 = \delta_1 | (Q(O(s^*)) = o)$ .

Alice can also consider adding noise to the output of  $Q$ , using a noising function  $N$ , again written in a probabilistic language. She can then consider Bob’s revised belief  $\delta_2 = \delta_1 | (N(Q(s^*)) = o)$ . She might consider obfuscation and noise at the same time, revising Bob’s belief to  $\delta_2 = \delta_1 | (N(Q(O(s^*))) = o)$ . Bob’s presumed revised belief assumes that he knows  $O$  and  $N$ .

If security is the only issue then a trivial noising function,  $N(s^*) \stackrel{\text{def}}{=} 42$ , ensures that Bob learns absolutely nothing about  $s^*$ . This ignores utility concerns that Bob (and Alice) might be interested in. One approach is to specify Bob’s utility in terms of the accuracy of what Bob can infer about the true output of  $Q$ , having learned some noised output of this function. Depending on the structure of  $Q$ , it is certainly possible to reveal some obfuscated or noised version  $Q'$  that simultaneously satisfy both security and utility objectives.

**Secrets as Probabilistic Programs** It might be appropriate for Alice to express her security concerns in terms of a function of her secret data, instead of the secret data itself. Her data might itself not be all that damaging, and in that setting it does not make sense for Alice to protect against its reconstruction by Bob. Instead Alice can define a “blacklist” function  $B$ , specified as a probabilistic program that operates on her data, whose output she would like to protect from inference by Bob.

In order to check security of  $B$ , and utility of  $Q$ , relative to Bob’s knowledge of her secret, she can evaluate these functions on the probabilistic value representing Bob’s knowledge, written  $B(\delta)$  and  $Q(\delta)$ , respectively. Alice can then attempt to add obfuscation/noise while making sure some agreed-upon utility policy  $U$  will be satisfied on Bob’s view on the output of  $Q$ .

**Definition.** Given belief  $\delta$ , security policy  $P$  and utility policy  $U$ , are satisfied for blacklist function  $B$  and utility function  $Q$  whenever  $P(B(\delta))$  and  $U(Q(\delta))$  hold.

Alice can determine the safety of  $B(\delta)$  and utility of  $Q(\delta)$  in a variety of ways. She could use the policy  $P$  from earlier on  $\delta_B$ , bounding Bob’s chances of guessing the output of  $B$ . She can also use the negation of  $P$ , specifically  $U(\delta) \stackrel{\text{def}}{=} \neg P(\delta) = \exists o. \delta(o) > t$ , to measure some minimum level of utility from  $Q$ , by checking that  $U(Q(\delta))$  holds. Alternate means could involve the ubiquitous

entropy quantity to measure a lower bound on the entropy of  $B(\delta)$  and an upper bound on the entropy of  $Q(\delta)$ .

Having the ability to measure the secrecy and utility of Bob’s belief, Alice can then determine how to construct a function  $Q'(s) = N(Q(O(s)))$  to ensure  $\delta|(Q'(s^*) = o)$  satisfies the security and utility policies (when possible). In recent work [1] we demonstrate how certain sets of features of  $s$  can be safely shared, while preventing Bob from guessing the output of a blacklist function  $B$  defined on same set of features.

We have not considered yet the design of obfuscation and noising functions. One approach is to assume that obfuscation and noising functions are specified as probabilistic programs with free parameters, which then need to be instantiated in such a way that utility and security can both be satisfied. How to do this for such “template” noising/obfuscation functions is an element of our ongoing work.

**Predictive models for blacklist functions** The blacklist (and utility) functions are useful to protect information which is not part of the secret value  $s^*$ , but is somehow correlated with the secret value. For example, Alice’s simple demographical information like age, postal code, and gender might compose the structure of  $s^*$  but Alice might be more interested in protecting her religions preference, even though this is not a field in the system she is securing.

Predictive models or classifiers, however, might exist that do a reasonable job of predicting religion given demographical information. Naturally, Alice could specify one such model as her blacklist function  $B$ . Unfortunately this does not stop Bob from using a different model  $B$  which is as good or better at predicting Alice’s religion. This is especially problematic when there are two disjoint sets of fields in Alice’s secret  $s^*$  that are both good for predicting religion. Alice might obfuscate one, but leave the second set completely untouched.

In general it is impossible for Alice to ensure *all* possible models for religion are defeated by her noising/obfuscation functions. This is simply due to Bob’s potential model  $B(s) = r$  where  $r$  is a constant that just happens to be Alice’s religion. To make this problem reasonable, we restrict ourselves to models that Bob, assumed rational, could produce, given some public body of data available to both Alice and Bob.

**Background knowledge and future work** An assumption on the common background data is useful for both limiting the sorts of blacklist functions  $B$  Alice should be concerned about and for determining what Bob’s initial belief  $\delta$  about Alice’s secret fields can be.

If we assume Bob takes advantage of a common labeled dataset  $D_r$  of tuples  $(s_i, r_i)$  tying demographical information to religion, Bob can learn a variety of predictive models for religion via supervised machine learning techniques. Alice would like to therefore protect herself (by specifying blacklists) against any model that from Bob’s point of view is accurate—i.e. it does well on predicting religion in  $D_r$ . How to do this effectively is a challenging aspect we would like to address in the future.

A related problem is how Alice can construct  $\delta$ , the initial belief Bob has about her secret value. The use of probability distributions to model Bob’s knowledge already presumes that Alice’s secret value  $s^*$  is a sample from  $\delta$ . It is natural, then, to presume a shared data set  $D$  instances  $s_i$  among the population Alice belongs to, is also composed of samples from that distribution.

Alice could thus assume Bob would form an initial belief about her via a statistical model described as a sampling function  $S$  written in a probabilistic language that, given model parameters  $p$ , generates samples from the model. Assuming an initial distribution on  $p$ , labeled  $\delta_p$ , Alice could infer these parameters by conditioning on the data instances from  $D$ , as if they were samples generated by the model, or  $\delta_p|(S(p) = s_1)|(S(p) = s_2)|\dots$  as is demonstrated in recent work of Gordon et al.[3]. Bob’s initial belief about  $s^*$  would then be the distribution  $S(\delta_p)$ .

The problem here is the same as in the problem of specifying blacklist functions; Alice does not know what model Bob could use either to infer the blacklist (e.g. religion) or to form an initial belief. Alice could take into account multiple models, thereby having to track several options for Bob’s belief, but how to construct a sufficient set of models is another challenging element of our future work.

## References

- [1] Supriyo Chakraborty, Kasturi R. Raghavan, Mani B. Srivastava, Chatschik Bisdikian, and Lance M. Kaplan. Balancing value and risk in information sharing through obfuscation. In *Proceedings of the International Conference on Information Fusion (FUSION)*, 2012.
- [2] Michael R. Clarkson, Andrew C. Myers, and Fred B. Schneider. Quantifying information flow with beliefs. *Journal of Computer Security*, 17(5):655–701, 2009.
- [3] Andrew D. Gordon, Mihail Aizatulin, Johannes Borgstroem, Guillaume Claret, Thore Graepel, Aditya Nori, Sriram Rajamani, and Claudio Russo. A model-learner pattern for bayesian reasoning. In *Proceedings of the ACM SIGPLAN Conference on Principles of Programming Languages (POPL)*, 2013.
- [4] Piotr Mardziel, Michael Hicks, Jonathan Katz, and Mudhakar Srivatsa. Knowledge-oriented secure multiparty computation. In *Proceedings of the ACM SIGPLAN Workshop on Programming Languages and Analysis for Security (PLAS)*, 2012.
- [5] Piotr Mardziel, Stephen Magill, Michael Hicks, and Mudhakar Srivatsa. Dynamic enforcement of knowledge-based security policies. In *Proceedings of the IEEE Computer Security Foundations Symposium (CSF)*, 2011.
- [6] Geoffrey Smith. On the foundations of quantitative information flow. In *Proceedings of the Conference on Foundations of Software Science and Computation Structures (FoSSaCS)*, 2009.