

Lecture 17

Instructor: Madhu Sudan

Scribe: Patrick Guo

1 Overview

Today we will conclude that Information = Amortized Complexity. The main theorem we will show (from [1]) is that

Theorem 1 (Information = Amortized Complexity).

$$\frac{1}{n}CC_{\varepsilon,\mu}^n(f) = IC_{\varepsilon,\mu}(f)(1 + o_n(1))$$

Along the way we will also introduce Interactive Correlated Sampling.

Administrative:

- Project presentations May 1, 9am – 4pm, LISE 303
- Project writeups due May 7

2 Preliminaries

Recall the following definition from last lecture:

Definition 2 (Direct product of a function). *Given a function $f : X \times Y \rightarrow R$, its n -fold product is denoted by $f^{\otimes n} : X^n \times Y^n \rightarrow R^n$*

$$f^{\otimes n}(x_1, \dots, x_n, y_1, \dots, y_n) = (f(x_1, y_1), \dots, f(x_n, y_n))$$

We are interested in the communication complexity of $f^{\otimes n}$. Trivially we have $CC(f^{\otimes n}) \leq n \cdot CC(f)$ since we can solve $f^{\otimes n}$ by running the communication for f n times in parallel. $CC(f^{\otimes n})$ is interesting because perhaps the problem is easier when we are asked about n independent copies of the problem – How can solving $f(x_1, y_1)$ help solve $f(x_2, y_2)$, etc.? We will see today that there is reason to believe the iterated version of f is easier.

As a note, when we go from protocols on f on input distribution μ to protocols on $f^{\otimes n}$, we implicitly go from working with μ supported on $X \times Y$ to working with a distribution μ^n on $X^n \times Y^n$, where all n instances of the problem are generated independently of one another. We will just write μ to specify input distribution from now on, and it will be clear when we mean μ^n .

Recall we also redefined error:

Definition 3. $f^{\otimes n}$ is solved by $\bar{\Pi}$ with error ε if for all i ,

$$\Pr[f^{\otimes n}(\bar{x}, \bar{y})_i = \bar{\Pi}(\bar{x}, \bar{y})_i] \geq 1 - \varepsilon$$

This gives us the following definitions for communication and information complexity on the n -fold product:

Definition 4.

$$CC_{\varepsilon,\mu}^n(f) = \min_{\bar{\Pi}} \{CC(\bar{\Pi})\}$$

$$IC_{\varepsilon,\mu}^n(f) = \min_{\bar{\Pi}} \{IC(\bar{\Pi})\}$$

where the minimums are taken over $\bar{\Pi}$ solving $f^{\otimes n}$ on input distribution μ^n with error ε

The motivation behind relaxing our definition of error is that otherwise a protocol erring with probability ε for f when iterated on n independent instances of the problem gives a protocol for $f^{\otimes n}$ erring with probability $1 - (1 - \varepsilon)^n$, which is unideal. The relaxed definition is nice since it means an ε -error protocol on f leads to an ε -error protocol on $f^{\otimes n}$, though it also means an ε -error protocol on $f^{\otimes n}$ in actuality errs more often than ε , but this is fine since we are mainly proving lower bounds.

3 Information = Amortized Complexity

We will make use of the following lemma:

Lemma 5.

$$IC_{\varepsilon,\mu}^n(f) = n \cdot IC_{\varepsilon,\mu}^1(f)$$

Informally, this means that the information leaked by the best protocol for $f^{\otimes n}$ grows linearly in n , and will be used in the proof of Theorem 1 to allow us to compare $CC_{\varepsilon,\mu}^n(f)$ to $IC_{\varepsilon,\mu}^n(f)$, as well as to show that some terms are lower order.

Proof of Lemma. Note that $IC^n(f) \leq n \cdot IC^1(f)$ is intuitively obvious, since we can solve n copies of the problem by running the solution for 1 copy n times in parallel, i.e. by using $\Pi^{\otimes n}$.

Exercise 6. Rigorously show

$$IC_{\mu}(\Pi^{\otimes n}) \leq n \cdot IC_{\mu}(\Pi)$$

by expanding the definition of information complexity and applying the chain rule

Hence we are left with showing the other side of the inequality, $IC_{\varepsilon,\mu}^n(f) \geq n \cdot IC_{\varepsilon,\mu}^1(f)$, or equivalently,

$$IC_{\varepsilon,\mu}^1 \leq 1/n \cdot IC_{\varepsilon,\mu}^n(f)$$

Intuitively, we are trying to extract from a solution for n copies a solution a solution for 1 copy which somehow compresses its information cost. To this end, we do simulation with the same embedding trick we've seen previously in this class for proving the communication complexity of DISJOINTNESS (embed the task of solving 1 instance of the problem into a protocol that solves n)

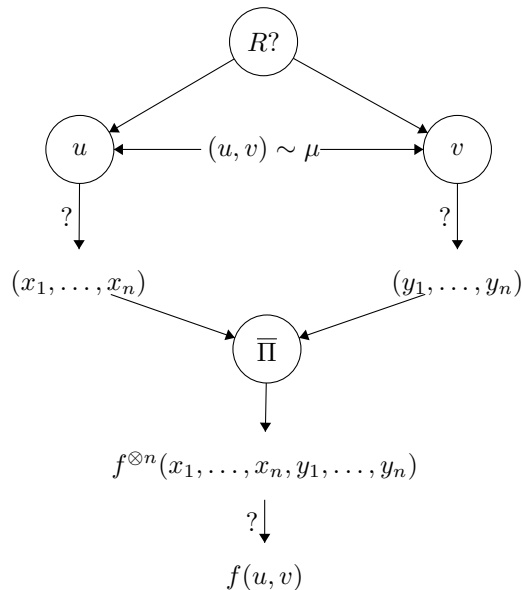


Figure 1: Extracting solution for 1 instance of problem from a solution for n instances through embedding

Like before, we need to determine the following:

- How to embed u, v into (x_1, \dots, x_n) and (y_1, \dots, y_n) , respectively
- How to generate the rest of the inputs' coordinates
- How to extract from $f^{\otimes n}(x_1, \dots, x_n, y_1, \dots, y_n)$ the value $f(u, v)$

Toward the first two points, we use shared randomness R . We simply sample $i \in_{Unif} [n]$ uniformly at random and set $x_i = u, y_i = v$. This is important since we don't know exactly at which coordinate in Π that information is leaked, so a uniform i ensures that we sum up equally over all possibilities to capture the information. Then the extraction is simply projection to the i th coordinate, or $f(u, v) = (f^{\otimes n}(x_1, \dots, x_n, y_1, \dots, y_n))_i$.

Now the question is, how do we generate the remaining coordinates of x, y ? We *could* just sample them iid from their marginal distributions according to μ , but the point is that we want to generate them from some joint, correlated distribution so that we can leverage information being leaked in the protocol for $f^{\otimes n}, \bar{\Pi}$. Thus, we sample with shared randomness $x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n$ from their marginal distributions according to μ . Then, to fill out the rest of the coordinates $x_{i+1}, \dots, x_n, y_1, \dots, y_{i-1}$, the players use private randomness to sample x_{i+1}, \dots, x_n according to their conditional distribution conditioned on the public y_{i+1}, \dots, y_n , and similarly for the other player to sample y_1, \dots, y_{i-1} .

As recap, our protocol Π to solve $f(u, v)$ from a protocol $\bar{\Pi}$ for $f^{\otimes n}$ is

- Using shared randomness, sample i uniformly at random from $[n]$, then sample $x_1, \dots, x_{i-1}, y_{i+1}, y_n$ from their marginal distributions according to μ
- Using private randomness, sample $x_{i+1}, \dots, x_n, y_1, \dots, y_{i-1}$ according to their conditional distributions conditioned on $x_1, \dots, x_{i-1}, y_{i+1}, y_n$
- Communicate according to $\bar{\Pi}$ to solve $f^{\otimes n}(x_1, \dots, x_n, y_1, \dots, y_n)$
- Output $f(u, v) = (\bar{\Pi}(x_1, \dots, x_n, y_1, \dots, y_n))_i$

Note that Π and $\bar{\Pi}$ have the same communication.

Now we compute the information complexity of this protocol. We have

$$IC(\Pi) = I(u; \Pi | v, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n) + I(v; \Pi | u, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n)$$

and we want to show $IC(\Pi) \leq \frac{1}{n} IC(\bar{\Pi})$, so we want to go from $I(u; \Pi | v, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n)$ to $\frac{1}{n} I(x_1, \dots, x_n; \bar{\Pi} | y_1, \dots, y_n)$. Since we are given i , we can rewrite the former term as $I(x_i; \Pi | y_i, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n)$. A step we will need is

Exercise 7. Create a Markov chain to argue through conditional independence that

$$I(x_i; \Pi | y_i, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n) = I(x_i; \Pi | y_i, i, x_1, x_{i-1}, y_1, \dots, y_i, y_{i+1}, \dots, y_n)$$

This is intuitively true since Bob generates y_1, \dots, y_{i-1} given x_1, \dots, x_{i-1} , which are used only to generate the communications of $\bar{\Pi}$, so given Π , y_1, \dots, y_{i-1} give no further information about x_i .

Hence, we can compute that

$$\begin{aligned} I(u; \Pi | v, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n) &= I(x_i; \Pi | y_i, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n) \\ &= I(x_i; \Pi | y_i, i, x_1, x_{i-1}, y_1, \dots, y_i, y_{i+1}, \dots, y_n) \\ &= \frac{1}{n} \sum_{j=1}^n I(x_j, \Pi | x_1, \dots, x_{j-1}, y_1, \dots, y_n) \\ &= \frac{1}{n} I(x_1, \dots, x_n; \bar{\Pi} | y_1, \dots, y_n) \end{aligned}$$

where the second to last equality comes from the fact that i is uniform over $[n]$, and the last equality from chain rule, i.e. $I(A_1, \dots, A_n; B|C) = \sum_{i=1}^n I(A_i; B|C, A_1, \dots, A_{i-1})$

By a completely symmetrical argument we have

$$I(v; \Pi|u, i, x_1, \dots, x_{i-1}, y_{i+1}, \dots, y_n) = \frac{1}{n} I(y_1, \dots, y_n; \Pi|x_1, \dots, x_n)$$

and putting them together gives

$$IC_{\varepsilon, \mu}^n(f) = n \cdot IC_{\varepsilon, \mu}^1(f)$$

as desired. □

Now back to the main theorem: recall that we are heading for Information = Amortized Complexity, or

$$\frac{1}{n} CC_{\varepsilon, \mu}^n(f) = IC_{\varepsilon, \mu}(f)(1 + o_n(1))$$

What this means is that we are trying to compress a communication that has little information but lots of communication by using the fact that we have multiple independent instances of the problem. As an equivalent reformulation using our lemma, we want to show

$$\begin{aligned} CC_{\varepsilon, \mu}^n(f) &\geq IC_{\varepsilon, \mu}^n(f)(1 \pm o(1)) \quad (\text{obvious}) \\ CC_{\varepsilon, \mu}^n(f) &\leq IC_{\varepsilon, \mu}^n(f)(1 \pm o(1)) \end{aligned}$$

The top inequality is straightforward, since for any working protocol, the amount of information conveyed cannot be more than the amount of bits sent in total.

Thus, we want to prove $CC_{\varepsilon, \mu}^n(f) \leq IC_{\varepsilon, \mu}^n(f)(1 \pm o(1))$. Specifically we will show

$$CC_{\varepsilon, \mu}^n(f) \leq IC_{\varepsilon, \mu}^n(f) + O(\sqrt{IC_{\varepsilon, \mu}^n(f) + o(1)})$$

Again by our lemma, since the n -fold information cost is linear in n , we know that $\sqrt{IC_{\varepsilon, \mu}^n(f)}$ is truly a lower order term.

Now, let Π be a k -round protocol for $f(x, y)$ (x being Alice's input, y being Bob's) with communication C and information I . We wish to compress this; specifically, we want to show existence of Π' simulating Π with communication $I + O(k\sqrt{I} + k \log \frac{k}{\varepsilon})$. The idea is to compress each step Π_i of the communication $\Pi = (\Pi_1, \dots, \Pi_k)$. Consider the first communication Π_1 from Alice to Bob. This is entirely a function of x , Alice's input, so the communication is exactly sampled from the distribution $\Pi_1|x$, call this distribution P . Now, to compress this, we want to send just enough information for Bob to reconstruct Π_1 with small probability of error, and for this we need to know how much apriori knowledge Bob had about the distribution of Π_1 . He only knows his input y , but if y is correlated with x , then Bob can have some informed apriori estimate of Π_1 from his distribution $\Pi_1|y$, call this distribution Q . Again as a toy example, suppose $P = Q$. Then with shared randomness Bob can simulate the entire communication on his own without any communication from Alice, so the communication needed (and information of the protocol) is 0. In general, the closer Q is to P , the less information is revealed by Alice's communication to Bob, and Bob can with less communication simulate the communication. This suggests that the amount of communication needed is related to the divergence between what Alice and Bob think Π_1 should be, and this is related to the information conveyed from Alice communicating her actual Π_1 .

Thus, this brings us to the problem of "Interactive Correlated Sampling." We want just enough communication between Alice and Bob such that, if P is supported on Ω , for all $a \in \Omega$ that

$$\Pr[Y = a|X = a] \geq 1 - \varepsilon$$

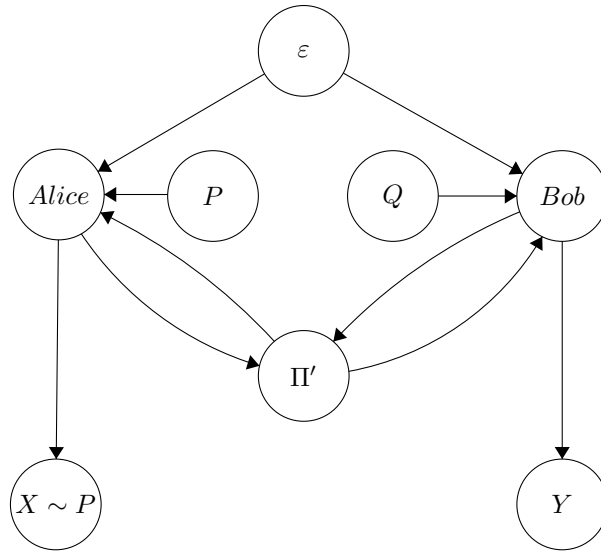


Figure 2: Interactive Correlated Sampling where $\Pr[Y = a|X = a] \geq 1 - \varepsilon$

We will show that there exists protocol Π' that achieves $\Pr[Y = a|X = a] \geq 1 - \varepsilon$ with $CC(\Pi') \leq D(P||Q) + O(\sqrt{D(P||Q)}) + \log \frac{1}{\varepsilon}$. The protocol Π' uses a similar dartboard sampling method as we saw earlier in the class for compressed interactive sampling.

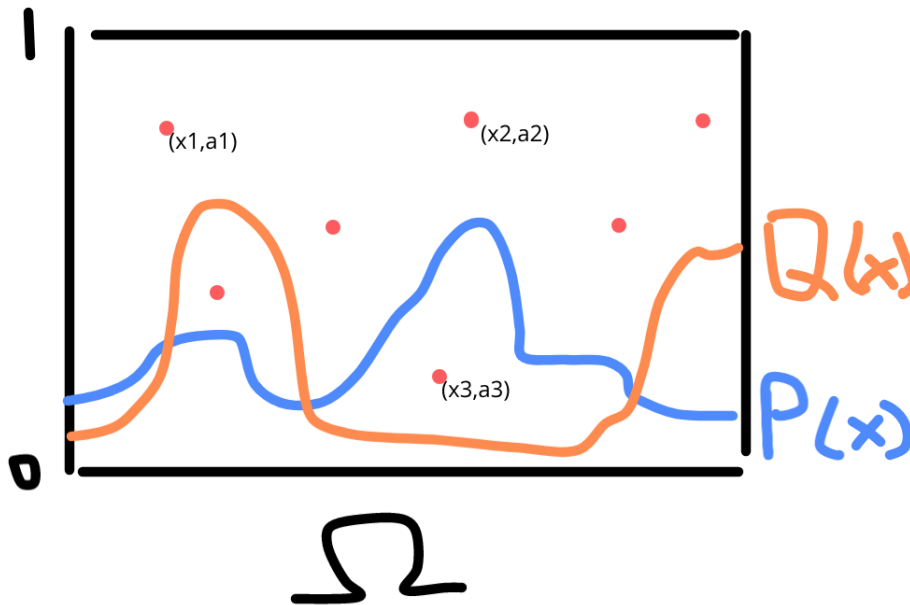


Figure 3: Dartboard sampling method for interactive correlated sampling. The horizontal axis is the discrete axis representing values in the support of P . The vertical axis ranges from 0 to 1, representing the probability of each value. Points (x_i, a_i) are uniformly and independently sampled on the dartboard. (x_3, a_3) is the first point under $P(x)$.

Alice's strategy is simple: she outputs $X = x_i$ where $a_i < P(x_i)$ and i is as small as possible. Note that, since points (x_i, a_i) are uniformly distributed on the board, they are uniformly distributed under the line $P(x)$ which is the PMF of P , so the first such point under the line (and in fact, any point under the line) has x coordinate distributed according to P , so $x_i \sim P$ as desired. Note that with error probability exponentially small in ε we have $i < \frac{|\Omega|}{\varepsilon}$, since the probability of any particular point lying under P is $\frac{1}{|\Omega|}$. Hence, it suffices to sample $\frac{|\Omega|}{\varepsilon}$ points (x, a) .

Now back to our toy example of $P = Q$ for some intuition. In this case, since the points (x_i, a_i) are generated with shared randomness, Bob can simply do the same strategy, since he knows Q which is equal to P , so he takes Y to also be the x coordinate of the first point under Q and no communication is required.

Now suppose we know a constant c such that $P \leq cQ$ (if no such constant exists, divergence is infinity and the result holds trivially). We also remark that the argument works for the weaker condition that $P(x_i) \leq cQ(x_i)$, since our argument only uses the fact that x_i is under the curve cQ . Bob's candidate points for X are then the x -coordinates of all points lying below his line cQ (and we can consider just the points with index $j < \frac{|\Omega|}{\varepsilon}$ since $i < \frac{|\Omega|}{\varepsilon}$ with exponentially small error), and we need enough communication between Alice and Bob to determine exactly which candidate point is correct. To do this, we use shared hash functions $h_j : \Omega \rightarrow \{0, 1\}$. In particular, Alice sends $O(\log c/\varepsilon)$ hash values of $x_i, (h_1(x_i) \dots h_{m \log c/\varepsilon}(x_i))$, Bob does the same hashes for all his candidate points, and then we will have with error probability ε that Bob identifies the correct point x_i .

So how do we determine c ? We do a search, trying 1, 2, 16..., each time giving us some candidate points, and as long as we don't get fooled by a wrong point (which hashing takes care of with the desired probability), once we try large enough c , we are done.

Algorithm 1 Protocol for Interactive Correlated Sampling

```

1: Assume that  $i < |\Omega|/\varepsilon$ 
2: for  $t = 0, 1, 2, 3, \dots$  do
3:   Let  $C_t = 2^{t^2}$ , "Hope that  $P(x_i) < C_t Q(x_i)$ "
4:   Alice sends  $\log C_t/\varepsilon$  bits of hash of  $x_i$  to Bob
5:   for  $j = 1, 2, 3, \dots, \frac{|\Omega|}{\varepsilon}$  do
6:     if  $a_j < C_t Q(x_j)$  then
7:       if Hashes of  $x_j$  agree with Alice's message then
8:         Bob sends message saying done
9:         Break
10:    else
11:      Continue

```

Note that the number of bits sent by the highest order computation is exactly $E_{x_i \sim P} \left[\log \frac{P(x_i)}{Q(x_i)} \right]$ which is just the divergence between P and Q .

Exercise 8. Formalize the argument's computations, i.e. show that the probability Bob identifies the wrong point is indeed bounded by ε , and compute that the amount of communication done is $D(P||Q) + O(\sqrt{D(P||Q)}) + \log \frac{1}{\varepsilon}$

This concludes compression based on divergence, and hence we have Information = Amortized Complexity. This result is particularly nice because it illustrates an operational view of divergence – if we both have individual distributions, but want to jointly sample from you, the amount of communication required (up to some lesser order terms) is equal to the divergence between our distributions.

References

- [1] Braverman, Mark, and Anup Rao. "Information equals amortized communication." IEEE Transactions on Information Theory 60.10 (2014): 6058-6069.