Introduction to Probability

Lecture 11: Estimators (Part II) Mateja Jamnik, <u>Thomas Sauerwald</u>

University of Cambridge, Department of Computer Science and Technology email: {mateja.jamnik,thomas.sauerwald}@cl.cam.ac.uk

Easter 2025



Compiled: May 18, 2025 at 19:41

Mean Squared Error

Estimating Population Size (Second Model)

- Suppose we have a sample of a few serial numbers (IDs) of some product
- We assume IDs are running from 1 to an unknown parameter N (so $N = \theta$)
- Each of the IDs is drawn without replacement from the discrete uniform distribution over {1,2,...,N}
- This is also known as Tank Estimation Problem or (Discrete) Taxi Problem

7, 3, 10, 46, 14



First Estimator Based on Sample Mean

Example 1	
Construct an unbiased estimator T_1 usin	g the sample mean.
	Answer

- Suppose n = 5
- Let the sample be

$$7, 3, 10, \frac{46}{14}, 14$$

The estimator returns:

$$T_1 = 2 \cdot \overline{X}_n - 1 = 2 \cdot \frac{80}{5} - 1 = 31 \quad \textcircled{B}$$
This estimator will often unnecessarily
underestimate the true value *N*.
Illenging exercise: Find a lower bound on **P**[$T_1 < \max(X_1, X_2, \dots, X_n)$]

- Achieving unbiasedness alone is not a good strategy
- Improvement: find an estimator which always returns a value at least max(X₁, X₂,..., X_n)

Cha

Intuition: Constructing an Estimator based on Maximum Sample

- Suppose *n* = 15
- Our samples are:

 $9, 82, 39, 35, 20, 51, 54, 62, 81, 29, {\color{red}84}, 59, 3, 34, 55$



Deriving the Estimator Based on Maximum Sample

Example 2	
Construct an unbiased estimator T_2 using m	$a_X(X_1,\ldots,X_n)$
	AllSwei

Empirical Analysis of the two Estimators





Figure: Histogram of 2000 values for T_1 and T_2 , when N = 1000 and n = 10.

Can we find a quantity that captures the superiority of T_2 over T_1 ?

Intro to Probability

Mean Squared Error

Estimating Population Size (Second Model)

Mean Squared Error

Mean Squared Error Definition —

Let T be an estimator for a parameter θ . The mean squared error of T is

$$\mathsf{MSE}[T] = \mathsf{E}\Big[(T-\theta)^2\Big].$$

• According to this, estimator T_1 better than T_2 if **MSE** $[T_1] <$ **MSE** $[T_2]$.

Bias-Variance Decomposition
The mean squared error can be decomposed into:

$$MSE[T] = (E[T] - \theta)^{2} + V[T]$$

$$= Bias^{2} = Variance$$

• If T_1 and T_2 are both unbiased, T_1 is better than T_2 iff $\mathbf{V}[T_1] < \mathbf{V}[T_2]$.



Bias-Variance Decomposition: Illustration



Source: Edwin Leuven (Point Estimation)

Example 4

It holds that **MSE**
$$[T_1] = \Theta\left(\frac{N^2}{n}\right)$$
, where $T_1 = 2 \cdot \overline{X}_n - 1$.

Answer

Analysis of the MSE for T_2 (non-examinable)



Mean Squared Error

Estimating Population Size (Second Model)

A New Estimation Problem



- Population/ID space S = {1, 2, ..., N}
- We take uniform samples from S without replacement
- Goal: Find estimator for N

Similar idea applies to situations where elements are not labelled before we see them first time (Mark & Recapture Method)

New Model

- Population/ID space of size |S| = N
- We take uniform samples from S with replacement
- Goal: Find estimator for N
- Suppose n = 6, N = 11, S = {3,4,7,8,10,15.83356,20,21,56,81,10000}
- Let the sample be

10, 81, 20, 3, 81, 10000

Let us call this a **collision**

As we do not know S, our only clue are elements that were sampled twice.

Birthday Problem

Birthday Problem: Given a set of k people

- What is the probability of having two with the same birthday (i.e., having at least one collision)?
- What is the expected number of people one needs to ask until the first collision occurs?



Estimation via Collision: The Algorithm

Recall: As we do not know S, our only information are collisions.

FIND-FIRST-COLLISION(S)

- 1: *C* = ∅
- 2: **For** *i* = 1, 2, . . .
- 3: Take next i.i.d. sample X_i from S
- 4: If $X_i \notin C$ then $C \leftarrow C \cup \{X_i\}$
- 5: else return T(i)
- 6: End For

T(i) will be the value of the estimator if algo returns after *i* rounds. (We want *T* unbiased)

• Running Time: The expected time until the algorithm stops is:

= the expected number of samples until a collision...

Same as the birthday problem, but now with |S| = N days... \odot

Expected Running Time (Knuth, Ramanujan)

$$\sqrt{\frac{\pi N}{2}} - \frac{1}{3} + O\left(\frac{1}{\sqrt{N}}\right).$$
Exercise: Prove a bound of $\leq 2 \cdot \sqrt{N}$

Estimation via Collision: Getting the Estimator Unbiased

Example 6		
One can define $T(i)$, $i \in \mathbb{N}$, such that $\mathbf{E}[T] = S $ for any finite, non-empty set S .		
	Answer	