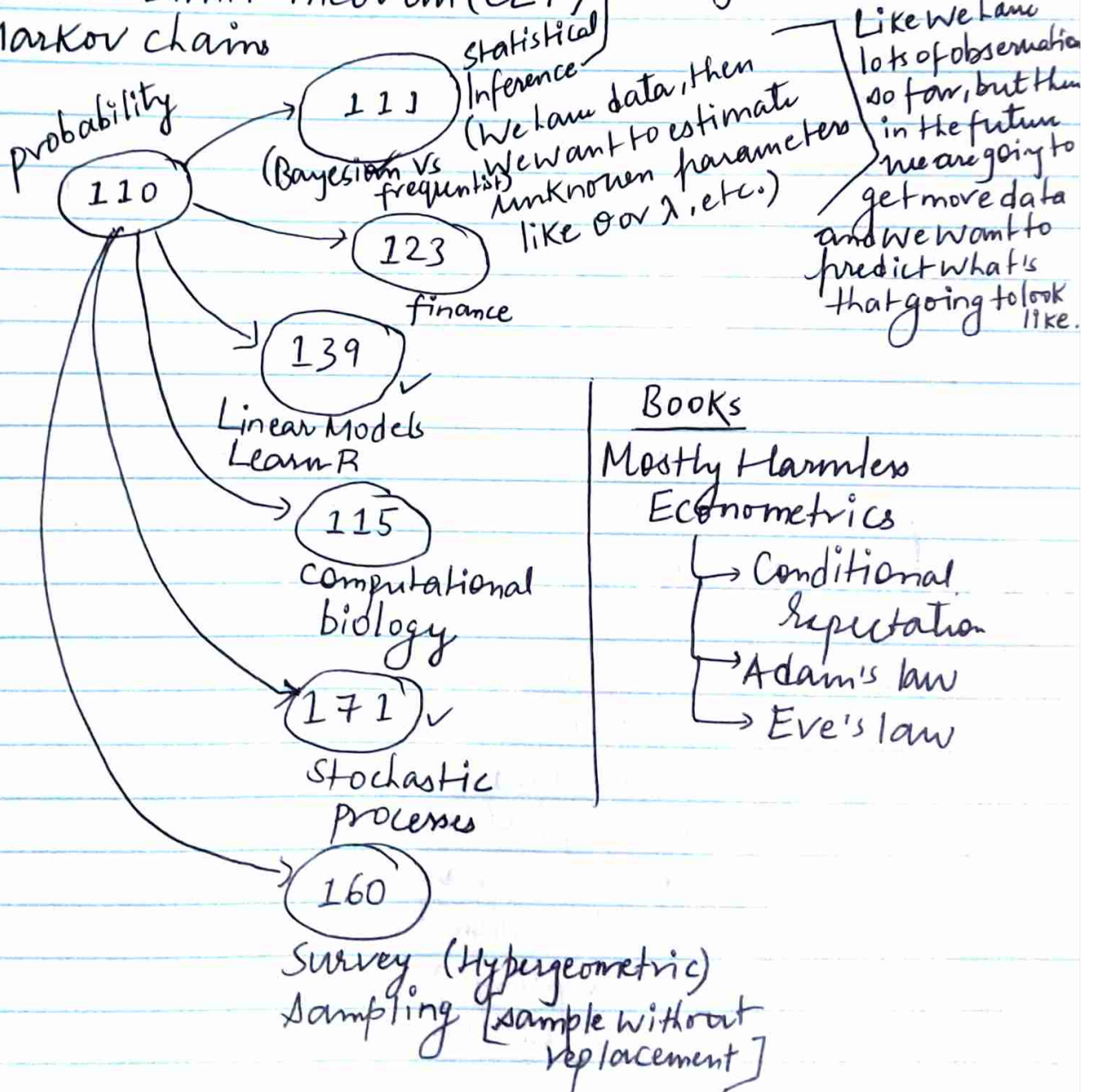# Lecture 34: A look Ahead

1. Conditioning
2. Symmetry
3. Random variables and their distributions
4. Stories
5. Linearity
6. Indicator random variables
7. LOTUS
8. Law of Large Numbers (LLN)
9. Central Limit theorem (CLT)
10. Markov chains

1–4 — What is randomness and uncertainty?

5–7 — Computing Expected values / averages.

8–10 — Long run behaviour

probability

110 → 111  Statistical Inference
(Bayesian vs frequentist)

(We have data, then we want to estimate unknown parameters like $\theta$ or $\lambda$, etc.)

→ 123  finance

→ 139  Linear Models Learn R

→ 115  Computational biology

→ 171  Stochastic processes

→ 160  Survey (Hypergeometric) Sampling [sample without replacement]

Like we have lots of observation so far, but then in the future we are going to get more data and we want to predict what's that going to look like.

Books

Mostly Harmless Econometrics
↳ Conditional Expectation
↳ Adam's law
↳ Eve's law

16

# Sampling from finite population

Let's true values (of something that we are interested in, i.e. for each person we have some variable. It could be their height, their income, their opinion on some question, whatever we are studying) are $(Y_1, Y_2, \ldots, Y_N)$, where $N$ is the size of the population. And let's assume each person has an ID number like social security number, so there's some well defined way to list them out.

→ treated as non-random, fixed (constant).

Sample of size $n$, goal is to estimate the average of $Y_1, Y_2, \ldots, Y_N$, i.e.,
$$\sum_{j=1}^{N} Y_j$$

prob. that person $j$ is in the sample is $P_j$ (Known). $P_j > 0$
So, the simple random sampling would be the case when all the $P_j$'s are equal, i.e., everyone is equally likely to be collected into our sample (But, obviously this may not be true. Some people may be much easier to sample than the others or some are just obscure, hard to reach, etc. So $P_j$'s may not be all equal in practice)

Let $(X_1, Z_1), \ldots, (X_n, Z_n)$ be the sample data. $X_j$ is the Y value we are interested in. Xs are random variables and Ys are fixed. The Xs are random because this is the first we collect into our survey, but that person was randomly chosen with some probabilities, so the value has become random because of the sampling; $Z$ is their ID number of that person, i.e. who did we actually get.

# Unbiased estimator for the total

Then $\sum_{j=1}^{n} \dfrac{X_j}{P_{Z_j}}$ is unbiased.

$Z_j$ is the random id number. So, in the denominator we have random probability.

( This says take each measurement and divide by the prob. that we actually got that person in our survey, then that's unbiased. )

## Proof

$$\sum_{j=1}^{n} \dfrac{X_j}{P_{Z_j}}$$

$$= \sum_{j=1}^{N} \dfrac{I_j Y_j}{P_j}$$

over entire population

, where $I_j$ is the indicator of $j$th person being included.

Expected value: $E\left( \sum_{j=1}^{N} \dfrac{I_j Y_j}{P_j} \right)$

By linearity

$$= \sum_{j=1}^{N} \dfrac{E(I_j) Y_j}{P_j}$$

$$= \sum_{j=1}^{N} \dfrac{P_j}{P_j} Y_j$$

$$= \sum_{j=1}^{N} Y_j$$

The expected value of $I_j$, by defn. and fundamental bridge is the prob. that person $j$ is included in the sample.

$= P_j$

Horwitz-Thomson Estimator

or, Inverse probability waiting

Is it a good estimator?

## Ex Basu's Elephant

# Sampling with Replacement and Sampling without Replacement

When we sample with replacement, the two sample values are independent. Practically, this means that what we get on the first one doesn't affect what we get on the second.
Mathematically, this means that the covariance between the two is 0.

In sampling without replacement, the two sample values aren't independent. Practically, this means that what we got on the for the first one affects what we can get for the second one.
Mathematically, this means that the covariance between the two isn't 0. This complicates the computations.

When we sample without replacement, and get a non-zero covariance, the covariance depends on the population size. If the population is very large, this covariance is very close to 0. In that case, sampling with replacement isn't much different from sampling without replacement. In some discussions, people describe this difference as sampling from an infinite population (sampling with replacement) versus sampling from a finite population (without replacement).