STAT 345/445 Lecture 17

Properties of a random sample

Random samples – Section 5.1 Sum of a random sample – Section 5.2

- Random Samples
- Statistic
 - Sampling distributions

Definition of a random sample

Random sample

Random variables X_1, \ldots, X_n are called a

random sample of size n from the population f(x)

if X_1, \ldots, X_n are

- mutually independent, and
- marginal pmf/pdf of each X_i is f(x)
- Alternative name for a random sample: independent and identically distributed (iid) random variables with pdf or pmf f(x)
 - i.i.d. f(x) = random sample from f(x)

Random samples and statistical inference

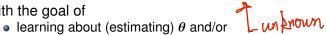
- We view data as observations of random variables
- Usually have more than one observation

$$X_1, X_2, \ldots, X_n$$

- Can often assume that X_1, X_2, \dots, X_n is a random sample
- We model the data by specifying a joint distribution

$$f(x_1, x_2, \ldots, x_n \mid \theta)$$

with the goal of





- predicting observations of X_{n+1}, X_{n+2}, \dots
- Use some summary of the data to do this
 - Need to find the distribution of that summary → sampling distribution

Frample: Say we collected data on lemps on CLB. airport. Say data points are X1=11°C, X2=12°C, X3=12-2°C,... We will assume that \$1, \$2, ... in \n=1) d one realizations of roundom variables. XI, X2,..., Mn if XI, X2,..., Xn ove a roundom sample, from form

then for, can be seen as the (population) distribution of temps at all limes at CLB. Calculate e.g. $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} \chi_i$

Sampling distribution

About random samples

Recall: Random variables X₁,..., X_n are (mutually) independent iif

$$f(x_1,\ldots,x_n)=f_1(x_1)\times\cdots\times f_n(x_n)$$

where $f_i(x_i)$ is the marginal pdf/pmf of X_i

- So, what is the joint pdf/pmf of a *random sample* from f(x)?
 - $f_i(x_i) = f(x_i)$ for all i, so

$$f(x_1, \dots, x_n) = f_1(x_1) \times \dots \times f_n(x_n)$$

= $f(x_1) \times \dots \times f(x_n)$
= $\prod_{i=1}^n f(x_i)$ marginals

Mutually independent

Recall again: Random variables X₁,..., X_n are (mutually) independent iif

$$f(x_1,\ldots,x_n)=f_1(x_1)\times\cdots\times f_n(x_n)$$

where $f_i(x_i)$ is the marginal pdf/pmf of X_i

- ullet any subcollection of X_1, \dots, X_n are also (mutually) independent.
- For example:

$$f(x_1, x_2) = \int \cdots \int f(x_1, \dots, x_n) dx_3 \cdots dx_4$$

$$= \int \cdots \int f_1(x_1) \times \cdots \times f_n(x_n) dx_3 \cdots dx_n$$

$$= f_1(x_1) f_2(x_2) \int f_3(x_3) dx_3 \times \cdots \times \int f_n(x_n) dx_n$$

$$= f_1(x_1) f_2(x_2)$$

More about random samples

- Not all collections of random variables are random samples
 - Need both independence and same (marginal) distributions
- If population is finite and we sample without replacement we don't get a random sample.

Example:

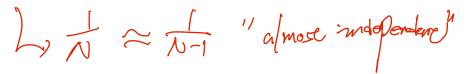
- Draw cards from a standard deck or 52 cards.
- Let X_i = the card we get in draw i, i = 1, ..., 10.
- All X_i have the same (marginal) distribution, but they are not independent since e.g.:

$$P(X_1 = 3\spadesuit) = \frac{1}{52}$$
 but $P(X_1 = 3\spadesuit \mid X_4 = 2\diamondsuit) = \frac{1}{51}$

STAT 345/445 Theoretical Statistics I Lecture 17

A simple random sample

- Sampling without replacement from a finite population is a very common
- A **simple random sample** of size $n, X_1, ..., X_n$ from a finite population of size N comes from a selection procedure were:
 - Any subset of n elements have the same probability of being selected.
- simple random sample ≠ random sample
- If N is huge we have simple random sample \approx random sample



More definitions

You just have to be able

A statistic

Let

- $X_1, ..., X_n$ be a random sample of size n possible
- $T(x_1,...,x_n)$ be a real-valued (or vector-valued) function with domain that includes the sample space of $(X_1,...,X_n)$

then

- The random variable (or random vector) $Y = T(X_1, ..., X_n)$ is called a **statistic**.
- The probability distribution of Y is called the sampling distribution of Y
- In short: A statistic is a function of a random sample.
- Note: Cannot be a function of a parameter.

In general, a statistic is a function of a collection of random variables (does not now to be a random sample).

Commonly seen statistics

sample mean:

$$\overline{X} = \frac{1}{n}(X_1 + \cdots + X_n) = \frac{1}{n}\sum_{i=1}^n X_i$$

sample variance:

$$S^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{i} - \overline{X})^{2}$$

sample standard deviation:

$$S = \sqrt{S^2}$$

The random variables \overline{X} , S^2 and S all have a (sampling) distribution

STAT 345/445 Theoretical Statistics I Lecture 17

Sampling distributions

- A lot of Statistical inference is based on sampling distributions
- Hence our focus on distributions of functions of random variables!
 - A bit easier if we have a random sample
- Example: If the mgf for the population exists:

Let $X_1, ..., X_n$ be a *random sample* of size n from a population with mgf $M_X(t)$. The the mgf of the sample mean is

$$M_{\overline{X}}(t) = (M_X(t/n))^n$$

• Only useful if we recognize the mgf on the right side

$$\overline{X} = \frac{1}{n} X_{1} + \dots + \frac{1}{n} X_{n}$$
 i.e. $b = 0$ and $ai = \frac{1}{n} Y_{i}$

Examples of sampling distributions

- Let $X_1, X_2, ..., X_n$ be a random sample from $N(\mu, \sigma^2)$. What is the distribution of \overline{X} ? • • $\overline{\chi} \sim N(\mu, \mathcal{L})$
- Let $X_1, X_2, ..., X_n$ be a random sample from $\operatorname{Gamma}(\alpha, \beta)$. What is the distribution of \overline{X} ?

From Jost Grouple:
$$\overline{n}$$
 is normal with mean $\sum_{i=1}^{n} \frac{1}{n} M = \frac{1}{n} n \mu = \mu$

and Varionice $\sum_{i=1}^{n} \frac{1}{n^2} L^2 = n \frac{1}{n^2} L^2 = \frac{L^2}{n}$

STAT 345/445 Theoretical Statistics I Lecture 17

 $= \left(\frac{1}{1-ne}\right)^{\times n} = mgf \text{ Gamma}$ $= 1 \text{ for Gramma}(dn, \frac{\beta}{n})$

Sampling distributions: convolution formula

79. if mgfs ove not available,

Convolution formula

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(w) f_Y(z - w) dw$$

proof...

inverse function
$$\chi: W = h.(Z_1, W),$$

$$(\chi, \gamma) \rightarrow (Z_2, W)$$

STAT 345/445

$$\frac{\partial x}{\partial z} = 0, \quad \frac{\partial x}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial z} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

$$= 0, \quad \frac{\partial y}{\partial w} = 1, \quad \frac{\partial y}{\partial w} = 1$$

 $\int_{XY} (X, Y) = \int_{X} (X) \int_{Y} (Y)$ Therefore: $\int_{ZW} (Z, W) = \int_{X} (W) \int_{Y} (Z-W) \cdot [-1]$

$$\int_{\mathbb{R}} f(z, w) dw = \int_{\mathbb{R}} f_{x}(w) f_{y}(z-w) dw$$

 $N > 2$, Just iterate:

$$2_1 = x_1 + x_2$$

 $2_1 = x_1 + x_2 + x_3 = 2_1 + x_3$

Moments of sampling distributions

Lemma

Let $X_1, ..., X_n$ be a random sample of size n from a population and let g(x) be a function such that $\mathbb{E}(g(X_1))$ and $\operatorname{Var}(g(X_1))$ exist. Then

$$\operatorname{E}\left(\sum_{i=1}^{n}g(X_{i})\right)=n\operatorname{E}\left(g(X_{1})\right)$$
 and $\operatorname{Var}\left(\sum_{i=1}^{n}g(X_{i})\right)=n\operatorname{Var}\left(g(X_{1})\right)$

STAT 345/445 Theoretical Statistics I Lecture 17 13/14

$$E\left(\frac{2}{2}g(x_1)\right) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (g(x_1) + g(x_2)) \cdot f(x_1, x_2, \dots x_n) dx_1, \dots dx_n$$

$$= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1, x_2, \dots x_n) dx_1, \dots dx_n + \dots + \dots \int_{-\infty}^{\infty} g(x_n) f(x_1, x_2, \dots x_n) dx_1, dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_1) \cdot f(x_1) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_n$$

$$= \int_{-\infty}^{\infty} g(x_n) f(x_n) dx_1 + \dots + \int_{-\infty}^{\infty} g(x_n) dx$$

For simplicity, www. Now (
$$\sum_{i=1}^{n} x_{i}^{1} = n \angle^{2}$$
)

Your ($\sum_{i=1}^{n} x_{i}^{1}$) = $\mathbb{E}\left(\sum_{i=1}^{n} x_{i}^{1} - \mathbb{E}\left(\sum_{i=1}^{n} x_{i}^{1}\right)^{2}\right)$

= $\mathbb{E}\left(\left(x_{1} - M + x_{2} - M \cdot \cdot \cdot x_{1} - M\right)^{2} + 2(x_{1} - M)(x_{2} - M) + \cdots + 2(x_{n} - M)(x_{n} - M)\right)$

= $\mathbb{E}\left((x_{1} - M)^{2} + \cdots + (x_{1} - M)^{2} + 2(x_{1} - M)(x_{2} - M) + \cdots + 2(x_{n} - M)(x_{n} - M)\right)$

= $\mathbb{E}\left((x_{1} - M)^{2} + \cdots + \mathbb{E}\left((x_{n} - M)^{2}\right) + \mathbb{E}\left((x_{1} - M)(x_{2} - M)\right) + \cdots + 2\mathbb{E}\left((x_{n} - M)(x_{n} - M)\right)$

= $\mathbb{E}\left((x_{1} - M)^{2} + \cdots + \mathbb{E}\left((x_{n} - M)^{2}\right) + \mathbb{E}\left((x_{1} - M)(x_{2} - M)\right) + \cdots + 2\mathbb{E}\left((x_{n} - M)(x_{n} - M)\right)$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= nx^{2}$$

$$= 0 \quad \text{if } x_{1}, \dots x_{n} \quad \text{on}$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= 0 \quad \text{if } x_{1}, \dots x_{n} \quad \text{on}$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + \cdots + x^{2} + 2 \overline{x} Gr(x_{i}, x_{j})$$

$$= x^{2} + x^{$$

Moments of sampling distributions

Theorem

Let $X_1, ..., X_n$ be a random sample of size n from a population with mean μ and variance $\sigma^2 < \infty$. Then

1.
$$E(\overline{X}) = \mu$$
 $E(\overline{h} \stackrel{\mathbf{M}}{\underset{i=1}{\sim}} X_i) = \frac{1}{n} \cdot n\mu = \mu$

2.
$$\operatorname{Var}(\overline{X}) = \frac{\sigma^2}{n} \quad \operatorname{Var}(\frac{1}{2}, \frac{N}{k!}) = \frac{1}{N^2} \cdot N^2 = \frac{L^2}{N}$$

3.
$$E(S^2) = \sigma^2$$
 ... ②

Useful fact: For any numbers x_1, \ldots, x_n we have

$$\sum_{i=1}^{n} (x_i - \overline{x})^2 = \sum_{i=1}^{n} x_i^2 - n\overline{x}^2$$

where
$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

STAT 345/445 Theoretical Statistics I Lecture 17

3:
$$S = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

 $E(S^2) = \frac{1}{n-1} E(\sum_{i=1}^{n} x_i^2 - n\bar{x}^2)$

$$= \frac{1}{n-1} \left(\frac{2}{2} \left[\frac{5(\chi_i^2)}{2} - n \frac{5(\chi_i^2)}{2} \right] - \cdots (x)$$

Recall: Vow(x) = &(x) - E(x) => &(x) = &2 + y 2 $(\frac{1}{2}): -\frac{1}{n-1} \left(\frac{1}{2} (2 + \mu^{2}) - n \left(\frac{2}{n} + \mu^{2} \right) \right)$ $\frac{2}{3} (x) = \sqrt{\alpha r (x^{2}) + \alpha^{2} (x^{2})^{2}}$ $\frac{2}{3} (x^{2})^{2}$

$$= \frac{1}{n-1} \left(n \mathcal{L}^2 + n \mu^2 - \mathcal{L}^2 - n \mu^2 \right) = \frac{1}{n-1} \left(n - n \mathcal{L}^2 \right) = \mathcal{L}^2$$