

## Sampling

DSE 210

- 1 Laws of large numbers
- 2 Basic sampling designs
- 3 Confidence intervals

## Review: expected value

The expected value of a random variable  $X$  is

$$\mathbb{E}(X) = \sum_x x \Pr(X = x).$$

Example: A coin has heads probability  $p$ . Let  $X$  be 1 if heads, 0 if tails.

$$\mathbb{E}(X) = 1 \cdot p + 0 \cdot (1 - p) = p.$$

Linearity properties:

- $\mathbb{E}(aX + b) = a\mathbb{E}(X) + b$  for any random variable  $X$  and any constants  $a, b$ .
- $\mathbb{E}(X_1 + \dots + X_k) = \mathbb{E}(X_1) + \dots + \mathbb{E}(X_k)$  for any random variables  $X_1, X_2, \dots, X_k$ .

Example: Toss  $n$  coins of bias  $p$ , and let  $X$  be the number of heads. What is  $\mathbb{E}(X)$ ?

Let the individual coins be  $X_1, \dots, X_n$ .

$$\mathbb{E}(X) = \mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}(X_1) + \dots + \mathbb{E}(X_n) = np.$$

## Review: variance

$$\text{var}(X) = \mathbb{E}(X - \mu)^2 = \mathbb{E}(X^2) - \mu^2, \text{ where } \mu = \mathbb{E}(X).$$

Toss a coin of bias  $p$ . Let  $X \in \{0, 1\}$  be the outcome.

$$\mathbb{E}(X) = p$$

$$\mathbb{E}(X^2) = p$$

$$\mathbb{E}(X - \mu)^2 = p^2 \cdot (1 - p) + (1 - p)^2 \cdot p = p(1 - p)$$

$$\mathbb{E}(X^2) - \mu^2 = p - p^2 = p(1 - p)$$

This variance is highest when  $p = 1/2$  (fair coin).

The standard deviation of  $X$  is  $\sqrt{\text{var}(X)}$ .

It is the average amount by which  $X$  differs from its mean.

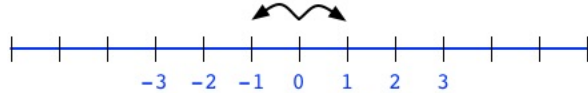
Useful variance rules:

- $\text{var}(X_1 + \dots + X_k) = \text{var}(X_1) + \dots + \text{var}(X_k)$  if  $X_i$ 's independent.
- $\text{var}(aX + b) = a^2 \text{var}(X)$ .

## Variance of a sum

$\text{var}(X_1 + \dots + X_k) = \text{var}(X_1) + \dots + \text{var}(X_k)$  if the  $X_i$  are independent.

Symmetric random walk. A drunken man sets out from a bar. At each time step, he either moves one step to the right or one step to the left, with equal probabilities. Roughly where is he after  $n$  steps?



Let  $X_i \in \{-1, 1\}$  be his  $i$ th step. Then  $\mathbb{E}(X_i) = 0$  and  $\text{var}(X_i) = 1$ .

His position after  $n$  steps is  $X = X_1 + \dots + X_n$ .

$$\mathbb{E}(X) = 0$$

$$\text{var}(X) = n$$

$$\text{stddev}(X) = \sqrt{n}$$

What is the distribution over his possible positions?

Approximately  $N(0, n)$ : Gaussian with mean 0 and std deviation  $\sqrt{n}$ .

## The central limit theorem

Suppose  $X_1, \dots, X_n$  are independent, and that they all come from the same distribution, with mean  $\mu$  and variance  $\sigma^2$ .

Let  $S_n = X_1 + \dots + X_n$ . Then  $S_n$  has mean and variance:

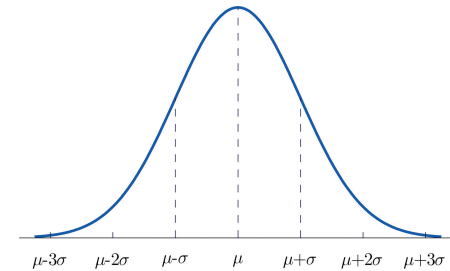
$$\mathbb{E}S_n = n\mu, \quad \text{var}(S_n) = n\sigma^2.$$

**Central limit theorem, very roughly:** For reasonably large  $n$ , the distribution of  $S_n = X_1 + \dots + X_n$  looks like  $N(n\mu, n\sigma^2)$ , the Gaussian with mean  $n\mu$  and variance  $n\sigma^2$ .

Question: What does this imply about the average  $(X_1 + \dots + X_n)/n$ ? What does its distribution look like?

Answer:  $N(\mu, \sigma^2/n)$ .

## The normal distribution



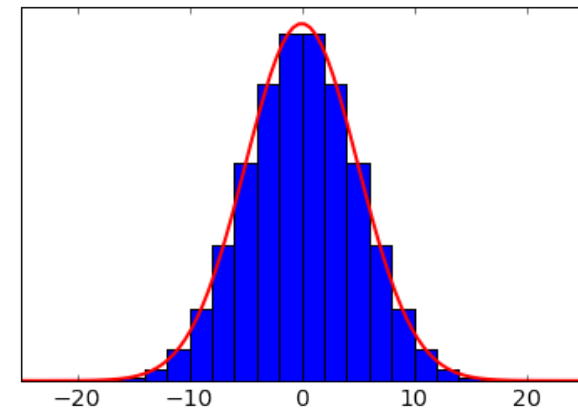
The normal (or *Gaussian*)  $N(\mu, \sigma^2)$  has mean  $\mu$ , variance  $\sigma^2$ , and density function

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

- 68.3% of the distribution lies within one standard deviation of the mean, i.e. in the range  $\mu \pm \sigma$
- 95.4% lies within  $\mu \pm 2\sigma$
- 99.7% lies within  $\mu \pm 3\sigma$

## Symmetric random walk, again

Each  $X_i$  is either 1 or  $-1$ , each with probability  $1/2$ . Therefore,  $X_1 + \dots + X_n$  is distributed like  $N(0, n)$ .



25 steps

## Tosses of a biased coin

A coin of bias (heads probability)  $p$  is tossed  $n$  times.

- What is the distribution of the observed **number** of heads, roughly?

Answer:  $N(np, np(1-p))$

Mean  $np$ , standard deviation on the order of  $\sqrt{n}$ .

- What is the distribution of the observed **fraction** of heads, roughly?

Answer:  $N(p, p(1-p)/n)$ .

Mean  $p$ , standard deviation on the order of  $1/\sqrt{n}$ .

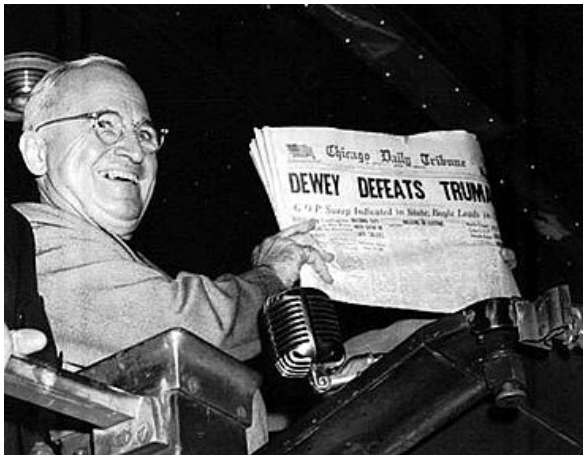
Example: A town has 30,000 registered voters, of whom 12,000 are Democrats. A random sample of 1,000 voters is chosen. How many of them would we expect to be Democrats, roughly?

Answer: The number of Democrats observed will roughly follow a  $N(1000 \times 0.4, 1000 \times 0.4 \times 0.6) = N(400, 240)$  distribution. This has mean 400 and standard deviation  $\approx 15.5$ .

## Outline

- 1 Laws of large numbers
- 2 Basic sampling designs
- 3 Confidence intervals

## Sampling design



In the 1948 Presidential election, the polls all predicted Thomas Dewey as the winner, with at least a five-point margin. But the outcome was quite different.

## Selection bias

The Republican bias in the Gallup Poll, 1936-1948.

Year	Gallup's prediction of Republican vote	Actual Republican vote
1936	44	38
1940	48	45
1944	48	46
1948	50	45

The safest way to sample is **at random**.

## Multistage cluster sampling

Sometimes random sampling is inconvenient, and careful multistage procedures need to be used.

For instance,

### ① Stage 1

- Divide the US into four geographical regions: Northeast, South, Midwest, West.
- Within each region, group together all population centers of similar sizes. E.g. All towns in the northeast with 50-250 thousand people.
- Pick a random sample of these towns.

### ② Stage 2

- Divide each town into wards, and each ward into precincts.
- Select some wards at random from the towns chosen earlier.
- Select some precincts at random from among these wards.
- Then select households at random from these precincts.
- Then select members of the selected households at random, within the designated age ranges.

## Outline

- ① Laws of large numbers
- ② Basic sampling designs
- ③ Confidence intervals

## Sample size versus population size

A certain town in Illinois has the same balance of Democrats and Republicans as the nation at large. We want to determine these fractions using a random sample of 1000 people. Would it be better to choose the 1000 people from the town in Illinois, or from the entire country?

Let the unknown fraction be  $p$ . In both cases, the observed fraction will follow the  $N(p, p(1-p)/1000)$  distribution.

What matters is the sample size, not the overall population size.

## Example: estimating a fraction

A university has 25,000 registered students. In a survey, 400 students were chosen at random, and it turned out that 317 of them were living at home. Estimate the fraction of students living at home.

The observed fraction, out of  $n = 400$  samples, is

$$\hat{p} = \frac{317}{400} \approx 0.79.$$

Give error bars on this estimate.

Let  $p$  be the fraction of students living at home. Then:

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right).$$

Therefore,  $\hat{p}$  has standard deviation  $\sqrt{p(1-p)/n}$ .

But we don't know  $p$ ... so what error bar to use?

In a survey,  $n = 400$  students were chosen at random, and it turned out that 317 of them were living at home.

The observed fraction living at home is  $\hat{p} = 0.79$ . This value  $\hat{p}$  is normally distributed with mean  $p$  and standard deviation  $\sqrt{p(1-p)/n}$ .

Since we don't know the true standard deviation  $\sqrt{p(1-p)}$  of each sample, use the observed standard deviation  $\sqrt{\hat{p}(1-\hat{p})}$ .

$$\text{stddev}(\hat{p}) \approx \sqrt{\frac{0.79 \times 0.21}{400}} \approx 0.02.$$

Using normal approximation gives confidence intervals:

- 68.3% interval:  $0.79 \pm 0.02$
- 95.5% interval:  $0.79 \pm 0.04$
- 99.7% interval:  $0.79 \pm 0.06$

What does a 95% confidence interval mean?

It means that if we were to do this over and over again, the interval would be correct (contain the true value) at least 95% of the time.

## Estimating an average

In a certain town, a random sample is taken of 400 people age 25 and over. The average years of schooling of this sample is 11.6 years, with a standard deviation of 4.1. Find a 95% confidence interval for the average educational level of people 25 and over in this town.

What is the distribution of the observed average?

- Let the true mean educational level be  $\mu$ , with stddev  $\sigma$ .
- We draw  $n$  samples from this distribution, and take the average  $\hat{\mu}$ .
- This  $\hat{\mu}$  has distribution  $N(\mu, \sigma^2/n)$ .

Estimate the standard deviation of  $\hat{\mu}$ .

- Its standard deviation is  $\sigma/\sqrt{n}$ .
- We don't know  $\sigma$ . Instead use the sample standard deviation, 4.1.
- Standard deviation of  $\hat{\mu}$  is roughly  $4.1/\sqrt{400} \approx 0.2$ .

Therefore, 95% confidence interval is  $11.6 \pm 0.4$ .

**And recall: the chance is in the measuring procedure, not in the quantity being estimated.**