

## Test 2 Checklist

1. probability of an event “E”

(a)  $0 \leq P(E) \leq 1$

(b) the complement of  $E$  (“not E”) is denoted  $\bar{E}$

$$P(\bar{E}) = 1 - P(E)$$

2. types of probability:

(a) empirical probability: relative frequency of historical results, experiments, simulations

(b) subjective probability: educated guess based on the knowledge/belief of the observer

(c) theoretical probability: mathematical model

3. probability distribution for a random variable  $x$  ( $x=L1$ ,  $P(x)=L2$ )

(a) sample space  $\Omega$  is the set of possible outcomes

(b) for empirical distributions, the sum of frequency is the sample size  $n$

(c) PDF = probability = proportion = relative frequency

(d) use CDF to get percentiles

(e)  $\sum P(x) = 1$

(f) “expected value” is the mean

4. randomness, simulations, law of large numbers

5. equally likely outcomes  $P(E) = \frac{\# E}{\# \Omega}$

6. joint probability

(a) “and” - use multiplication:  $P(A \text{ and } B) = P(A)P(B|A)$

(b)  $A$  and  $B$  are independent if  $P(A \text{ and } B) = P(A)P(B)$  (the events don’t affect each other)

(c) multiple events: multiply probabilities of each step, given the previous ones

(d) with / without replacement

(e) “or” - use addition: if an event could happen different ways, add the probabilities  
e.g. probability of 2 heads in three tosses =  $P(\text{HHT}) + P(\text{HTH}) + P(\text{THH})$

(f) don’t double-count, use a Venn diagram if events are not disjoint

(g) “at least one” easier to work with if you translate to “not zero”

7. binomial distribution  $X \sim BI(n, p)$

(a)  $x$  is the sum of repeated, independent, and identical binary trials

(b) theoretical mean  $\mu = np$ , st dev  $\sigma = \sqrt{np(1-p)}$

(c)  $\text{binompdf}(n, p, x)$  gives the probability of exactly  $x$

(d)  $\text{binomcdf}(n, p, x)$  gives the probability of less than or equal to  $x$

(e) be able to recognize binomial distribution in applications

8. normal distribution  $X \sim N(\mu, \sigma)$

(a) bell shaped with thin tails

- (b) 68, 95, 99.7 “empirical rule”, and connection with z-scores
- (c) draw a picture when you work these problems
- (d)  $\text{normcdf}(a, b, \mu, \sigma)$  gives the probability that  $x$  is between  $a$  and  $b$
- (e)  $\text{invnorm}(A, \mu, \sigma)$  gives the  $x$  value with area  $A$  to the left of it

9. sampling distribution

- (a) empirical/sample statistics  $(\bar{x}, s)$  versus theoretical/population statistics  $(\mu, \sigma)$ .
- (b)  $\bar{x}$  (the group average) is itself a random variable, and

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

this is called the “standard error” of  $\bar{x}$

- (c) Central Limit Theorem: if the individuals are independent, then as  $n$  gets larger,

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

- (d) use the CLT when you are interested in the group average or total

10. Linear regression and correlation

- (a) scatterplots
- (b) use LinRegTtest to find and work with the regression line  $y = a + bx$  that best fits the data.
- (c) the slope of the line is the coefficient of the  $x$ .
- (d) You can use the regression line to estimate  $y$  for a given  $x$ :
  - interpolation if  $x$  is within known points
  - extrapolation if it lies outside known points
- (e) the sample correlation  $r$  is between  $-1$  and  $1$
- (f)  $|r|$  measures how tightly the cloud of data points cluster around the regression line.
- (g)  $r$  has the same sign as the slope
- (h) Be able to explain positive/negative and weak/strong correlation intuitively.
- (i) Know the difference between correlation and causation.

11. models

- (a) be able to select an appropriate model for the given problem or data
- (b) discrete (e.g. binomial) vs continuous (e.g. normal) distributions
- (c) area under the graph represents probability
- (d) for a continuous random variable and a number  $c$ ,  $P(x = c) = 0$ ; because a line contains zero area
- (e) translate words into equations or inequalities, and then find probability using PDF or CDF
- (f) be able to work two step problems
- (g) be able to solve for parameters (e.g.  $n, p, \mu, \sigma$ ) to satisfy a given requirement, possibly using guess-check
- (h) remember that many models assume that observations are independent