

**Problem 4 (Checking independence):** Most of the techniques used in checking (conditional) independence between random variables are based on hypotheses testing. Formally, we have following two hypotheses:

$$H_0 = \{\text{Random variables are not (conditionally) independent}\}$$

and

$$H_1 = \{\text{Random variables are (conditionally) independent}\}.$$

By assuming the data follow some types of distribution, we can propose several types of metrics and approximate their distributions, in order to claim which hypothesis is true given our observations. The confidence level of our judgement can be measured using following metric, namely  $p$ -value:

$$p = P(H_1 \text{ appears to be true} | \text{data}, H_0 \text{ is actually true}).$$

Obviously, small  $p$ -value implies stronger confidence that  $H_1$  is true. Typically,  $p = 0.05$  or  $p = 0.1$  is used as thresholds to determine whether accept  $H_1$  or  $H_0$ . In this problem, I prefer to use G-test to determine whether two random variables are independent. See <https://en.wikipedia.org/wiki/G-test>. G-test is basically comparing the actual distribution and the product of two marginal distributions by counting each bin. The divergence between two distributions is measured as:

$$G(p, q) = 2 \sum_x p(x) \log\left(\frac{p(x)}{q(x)}\right),$$

where  $p(x)$  and  $q(x)$  are number of observed counts in a bin, instead of a probability mass. Therefore it is different from Kullback-Leibler divergence.

For the implementation I used the code from an existing package: Python `gsq` package. The code I used in this problem can be found in the corresponding attachment.

The conclusion is summarized as following:

**Marginal independence**  $p$ -value of  $d_0$  and  $d_1$  is: 0.0

$p$ -value of  $d_0$  and  $d_2$  is: 0.231658736846

$p$ -value of  $d_0$  and  $d_3$  is: 0.90331010682

$p$ -value of  $d_0$  and  $d_4$  is: 0.0

$p$ -value of  $d_1$  and  $d_2$  is: 0.166504089458

$p$ -value of  $d_1$  and  $d_3$  is: 0.535798622053

$p$ -value of  $d_1$  and  $d_4$  is: 0.0

$p$ -value of  $d_2$  and  $d_3$  is: 0.0253836754723

$p$ -value of  $d_2$  and  $d_4$  is: 0.804517201563

$p$ -value of  $d_3$  and  $d_4$  is: 0.218212152554

Therefore I am confident that  $(d_0, d_1)$ ,  $(d_0, d_4)$ ,  $(d_1, d_4)$  and  $(d_2, d_3)$  are pairs of independent random variables, and the  $p$ -value is the probability that I make such conclusion by mistake given that they are actually dependent.

As for 3-way mutual independence, the only possibility is  $(d_0, d_1, d_4)$  since mutual independence implies pairwise independence. In order to verify that, we only need to check whether  $d_1$  and  $d_4$  are independent conditioned on  $d_0$ . Because if that statement is true, we have the following:

$$P(d_1)P(d_4) = P(d_1|d_0)P(d_4|d_0) = P(d_1, d_4|d_0) = P(d_1, d_4, d_0)/P(d_0),$$

then  $P(d1, d4, d0) = P(d1)P(d4)P(d0)$ . And the  $p$ -value is 0.0 so they are mutually independent.

**Conditional independence** p-value of d1 and d4 conditioned on d0 is: 0.0

p-value of d0 and d1 conditioned on d2 is: 0.0

p-value of d0 and d1 conditioned on d3 is: 0.0

p-value of d0 and d4 conditioned on d1 is: 0.0

p-value of d0 and d4 conditioned on d2 is: 0.0

p-value of d0 and d4 conditioned on d3 is: 0.0

p-value of d1 and d2 conditioned on d3 is: 0.0916463098153

p-value of d1 and d4 conditioned on d0 is: 0.0

p-value of d1 and d4 conditioned on d2 is: 0.0

p-value of d2 and d3 conditioned on d1 is: 0.0551626358929

There is no conditional independence when conditioned on 2 or more other random variables.