

Homework 5

Solutions

1. (Short questions)

(a) True. In the notes we showed that

$$E(\tilde{a}^2) - E^2(\tilde{a}) = E((\tilde{a} - E(\tilde{a}))^2). \quad (1)$$

For discrete random variables this equals the sum of a nonnegative quantity multiplied by a nonnegative quantity:

$$E((\tilde{a} - E(\tilde{a}))^2) = \sum_{a \in R_{\tilde{a}}} (a - E(\tilde{a}))^2 p_{\tilde{a}}(a), \quad (2)$$

so the result must be nonnegative, which implies $E(\tilde{a}^2) - E^2(\tilde{a}) \geq 0$. For continuous random variables the argument is exactly the same with integrals instead of sums.

(A more straightforward solution): Let \tilde{a} take values $\{-1, 0, 1\}$ with equal probability. $E(\tilde{a}) = 0$ but $E(\tilde{a}^2) > 0$. That is $E^2(\tilde{a})$ can be less than $E(\tilde{a}^2)$

(b) True. Let m denote the median of \tilde{a} , we have

$$F_{\tilde{a}+b}(m+b) = P(\tilde{a} + b \leq m+b) \quad (3)$$

$$= P(\tilde{a} \leq m) \quad (4)$$

$$= \frac{1}{2} \quad (5)$$

so $m+b$ is the median of $\tilde{a} + b$.

(c) True. Since \tilde{a} and \tilde{b} have the same distribution $E(\tilde{a}) = E(\tilde{b})$ (since the expectation operator only depends on the pmf or pdf). By independence $E(\tilde{a}\tilde{b}) = E(\tilde{a})E(\tilde{b}) = E^2(\tilde{a})$.

(d) The probability of the event occurring is $1/n$ so

$$p_{I_i}(0) = 1 - \frac{1}{n}, \quad (6)$$

$$p_{I_i}(1) = \frac{1}{n}. \quad (7)$$

By linearity of expectation

$$E(\text{Number of kids that get their own present}) = E\left(\sum_{i=1}^n I_i\right) \quad (8)$$

$$= \sum_{i=1}^n E(I_i) \quad (9)$$

$$= \sum_{i=1}^n p_{I_i}(1) \quad (10)$$

$$= 1. \quad (11)$$

On average one kid gets their parents' own present.

2. (Paste and Rice)

- (a) The constraints are $100 \leq \max\{\tilde{x}, \tilde{r}\} \leq 300$. The joint pdf is constant over that region. Figure 1 contains the diagram.

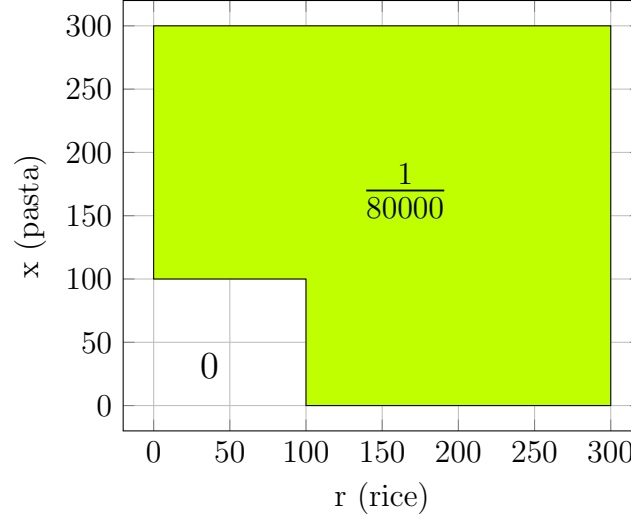


Figure 1: Joint pdf of r and x.

- (b) We compute

$$E(\tilde{r}\tilde{x}) = \int_{x=100}^{300} \int_{r=0}^{300} \frac{rx}{80000} dx dr + \int_{x=0}^{100} \int_{r=100}^{300} \frac{rx}{80000} dx dr \quad (12)$$

$$= \frac{1}{80000} \left(\frac{x^2}{2} \Big|_{100}^{300} \frac{r^2}{2} \Big|_0^{300} + \frac{x^2}{2} \Big|_0^{100} \frac{r^2}{2} \Big|_{100}^{300} \right) \quad (13)$$

$$= 25000. \quad (14)$$

$$E(\tilde{x}) = \int_{x=100}^{300} \int_{r=0}^{300} \frac{x}{80000} dx dr + \int_{x=0}^{100} \int_{r=100}^{300} \frac{x}{80000} dx dr \quad (15)$$

$$= \frac{1}{80000} \left(300 \frac{x^2}{2} \Big|_{100}^{300} + 200 \frac{x^2}{2} \Big|_0^{100} \frac{r^2}{2} \Big|_{100}^{300} \right) \quad (16)$$

$$= 162.5. \quad (17)$$

By symmetry $E(\tilde{r}) = 162.5$. The covariance $\text{Cov}(\tilde{r}, \tilde{x}) = E(\tilde{r}\tilde{x}) - E(\tilde{r})E(\tilde{x}) = -1406.25$ so \tilde{r} and \tilde{x} are negatively correlated.

- (c) The variables are correlated, which implies they cannot be independent.

3. (Law of conditional variance)

- (a) Given $\tilde{a} = a$, \tilde{a} is constant, and its value is x , therefore once x is fixed $\text{Var}(\tilde{b}|\tilde{a} = a)$ is a **number**. It represents the variance of the random variable \tilde{b} given the information

$\tilde{a} = a$. In other words, given two random variables \tilde{a} and \tilde{b} , look at their joint density, condition \tilde{a} to the value x . This slice gives rise to a new distribution of the random variable $\tilde{b}|\tilde{a} = a$, $\text{Var}(\tilde{b}|\tilde{a} = a)$ is precisely the variance of this random variable.

$$\text{Var}(\tilde{b}|\tilde{a} = a) = E((\tilde{b} - E(\tilde{b}|\tilde{a} = a))^2|\tilde{a} = a) \quad (18)$$

- (b) Now we don't keep x fixed but rather we treat it as a variable, for any x we assign a value $\text{Var}(\tilde{b}|\tilde{a} = a)$, this is a function from the range of \tilde{a} to real numbers, which we call by h , so $h(x) = \text{Var}(\tilde{b}|\tilde{a} = a)$. Since this function is defined on a probability space we can regard it as a **random variable** by $h(\tilde{a}) = \text{Var}(\tilde{b}|\tilde{a})$.
- (c) Using iterated expectations we get:

$$E(\text{Var}(\tilde{b}|\tilde{a})) = E(E(\tilde{b}^2|\tilde{a})) - E(E(\tilde{b}|\tilde{a})^2) \quad (19)$$

$$= E(\tilde{b}^2) - E(E(\tilde{b}|\tilde{a})^2) \quad (20)$$

$$\text{Var}(E(\tilde{b}|\tilde{a})) = E(E(\tilde{b}|\tilde{a})^2) - E(E(\tilde{b}|\tilde{a}))^2 \quad (21)$$

$$= E(E(\tilde{b}|\tilde{a})^2) - E(\tilde{b})^2 \quad (22)$$

$$\implies E(\text{Var}(\tilde{b}|\tilde{a})) + \text{Var}(E(\tilde{b}|\tilde{a})) = E(\tilde{b}^2) - E(\tilde{b})^2 = \text{Var}(\tilde{b}) \quad (23)$$

Average of the variance of \tilde{b} given \tilde{a} , plus the variance of the average of \tilde{b} given \tilde{a} .

- (d) Let \tilde{t} be the time at which a runner gets injured. And let \tilde{a} be the random variable that describes the age group: $\tilde{a} = 1$ be the group of runners below 30, and $\tilde{a} = 2$ be the group above 30. So that $\tilde{t}|\{\tilde{a} = 1\} \sim \exp(1)$ and $\tilde{t}|\{\tilde{a} = 2\} \sim \exp(2)$. We are also given that $P(\tilde{a} = 2) = 0.2$. Also note that if $\tilde{a} \sim \exp(\lambda)$ then $E(\tilde{a}^2) = \text{Var}(\tilde{a}) + E(\tilde{a})^2 = 1/\lambda^2 + 1/\lambda^2 = 2/\lambda^2$.

$$E(\tilde{t}) = E(E(\tilde{t}|\tilde{a})) = E(\tilde{t}|\tilde{a} = 1)P(\tilde{a} = 1) + E(\tilde{t}|\tilde{a} = 2)P(\tilde{a} = 2) \quad (24)$$

$$= 1 \cdot 0.8 + \frac{1}{2} \cdot 0.2 = 0.9 \quad (25)$$

$$E(\tilde{t}^2) = E(E(\tilde{t}^2|\tilde{a})) = E(\tilde{t}^2|\tilde{a} = 1)P(\tilde{a} = 1) + E(\tilde{t}^2|\tilde{a} = 2)P(\tilde{a} = 2) \quad (26)$$

$$= 2/1^2 \cdot 0.8 + 2/2^2 \cdot 0.2 = 1.7 \quad (27)$$

$$\implies \text{Var}(\tilde{t}) = 1.7 - (0.9)^2 = 0.89 \quad (28)$$

4. (Water salinity and temperature)

- (a) Please refer to Fig 2. The linear estimate has a negative slope, implying that the salinity decreases with temperature.
- (b) Please refer to Fig 2. According to the conditional means, the salinity decreases with temperature and then increases again.
- (c) The reliability of the conditional mean estimator depends on the number of data points that fell into the bin. In Fig 2 the conditional mean estimates to the left and more reliable than the one on the right, as each bin have more data points.
- (d) I have no idea!

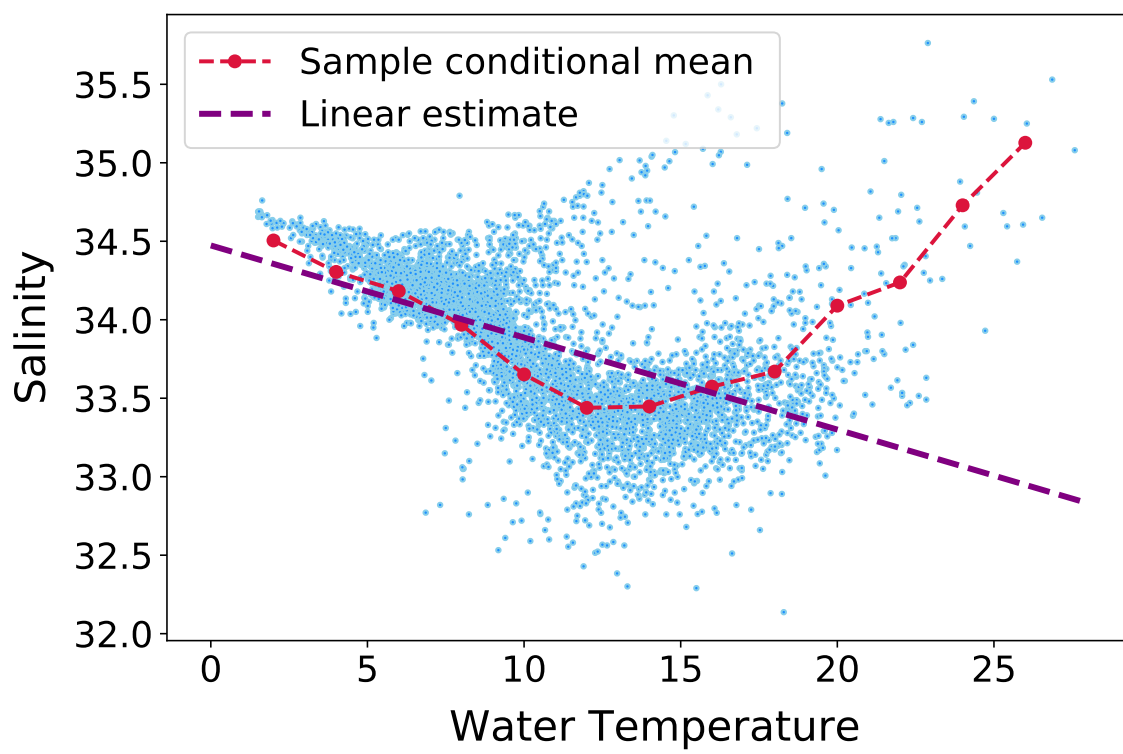


Figure 2