

UPPSALA UNIVERSITET

FÖRELÄSSNINGSANTECKNINGAR

Sannolikhetssteori 2

Rami Abou Zahra

Inlämningsdatum
September 7, 2023

CONTENTS

1. Repetition	2
1.1. Probability measure	2
1.2. Conditional probability theory	3
1.3. Independent events	3
1.4. Random variables	4
2. Transformations	6
2.1. Pre-knowledge	6
3. Multivariate random variables	6
3.1. Multivariate case	7
3.2. Conditional Probabilities	9
3.3. Conditional expected values and variances	10

1. REPETITION

Anmärkning: Det rekommenderas starkt att läsa igenom anteckningarna från Sannolikhetsteori 1

Definition/Sats 1.1: Random trial

An event is not certain, it usually has a probability associated with it. Taking that "risk" to see what the outcome is, is called a random trial.

Examples of random trials include throwing dice, picking cards, number of people who pass a road

Different possibilities (outcomes).

In the example of the dice, the outcomes are 1-6

Definition/Sats 1.2: Events

An event is something that happens (or does not happen) when you the random trial

You can have an event based on one outcome, or multiple.

Example (one outcome): The dice is 3 after a throw

Example (several outcomes): The card is 7 or lower (1,2,3,4,5,6, all the different colours)

Example (0 outcomes): The card shows both spades and hearts at the same time (impossible)

1.1. Probability measure.

Related to the probability that an event occurs.

Definition/Sats 1.3: Probability measure

A *probability measure* is a function which satisfies Kolmogorov's axioms and for each event gives a number $\in [0, 1]$

The number is called the *probability* of the event. Usually denoted $P = P(A)$ where $A \subset \Omega$

Definition/Sats 1.4: Kolmogorov's axioms

Let $P : 2^\Omega \rightarrow \mathbb{R}$. P is called a *probability measure* if it satisfies the following

- $P(A) \geq 0 \quad \forall A \in 2^\Omega$
- $P(\Omega) = 1$
- $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i) \quad (A_i \text{ disjoint})$

1.2. Conditional probability theory.

One can think of conditional probability in the term of Venn Diagrams in order to create intuition. Usually, if one event has happened it will affect other events and it is of interest to take this into consideration when calculating the probability of events.

The probability of an event A occurring given that the event B has occurred is denoted by

$$P(A|B)$$

Example:

Let A be the event that a person has 2 daughters, let B be the event that a person has 0 daughters, and C be the event that he has at least 1 daughter.

The probability $P(A|B)$ is of course 0, since given that he has 0 daughters, the probability is 0 for him to have 2 at the same time as he has 0

$P(B|C)$ is also 0, using similar argument as above

$P(A|C) = \frac{P(A \cap C)}{P(C)}$ We cannot say much here, other than that the probability is strictly positive since we already have one child

Definition/Sats 1.5: Bayes theorem

Let F_1, \dots, F_n be disjoint events $\in \Omega$ with $P(F_i) > 0$, and $P(\bigcup F_i) = 1$.

$$P(F_j|E) = \frac{P(E|F_j) \cdot P(F_j)}{\sum_{i=1}^n P(E|F_i)P(F_i)}$$

Example:

Suppose we have 3 different cards. The first card is red on both sides (RR), the second card is black on both sides (BB), and the third card is black and red (RB)

We draw a card at random of these three cards such that we only see one side of the card. Now suppose the side we see is red, what is the probability that the other side is black?

We are interested in the event $P(RB|R)$:

$$\begin{aligned} \frac{P(RB \cap R)}{P(R)} &= \underbrace{\quad}_{\text{Bayes}} = \frac{P(R|RB)P(RB)}{P(R|RR)P(RR) + P(R|RB)P(RB) + P(R|BB)P(BB)} \\ &= \frac{(1/2)(1/3)}{1 \cdot (1/3) + (1/2) \cdot (1/3) + (0) \cdot (1/3)} = \frac{1}{3} \end{aligned}$$

1.3. Independent events.

Definition/Sats 1.6: Independent events

If $P(A|B) = P(A)$, then A and B are independent

Example:

Let A be the event that 2 parents get a daughter, and B be the event that the neighbors child ate an ice cream yesterday.

Since these events do not affect each other, they are independent.

Example:

Let A be the event that the first throw of a dice yields 6, and let B be the event that the second throw is 3. Then A and B are independent since the first throw does not affect the second throw:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$$

Anmärkning:

There is an equivalent definition for independence through the following:

$$P(A \cap B) = P(A)P(B)$$

Anmärkning:

Independence is a symmetric relationship.

1.4. Random variables.

Definition/Sats 1.7: Random variable

A *random variable* is a function that for each outcome associates a number with it.

An example is a person's age, or the value of a card drawn. If the outcome is random, the number is also random.

Each random variable has a distribution function associated with it, and is defined as $F(X) = P(X \leq x)$

Anmärkning:

$$\begin{aligned} \lim_{X \rightarrow -\infty} F(X) &= 0 \\ \lim_{X \rightarrow \infty} F(X) &= 1 \end{aligned}$$

If $X_1 < X_2 \Rightarrow F(X_1) \leq F(X_2)$

We also have that F is right-continuous, meaning

$$\lim_{X \rightarrow a^+} F(X) = F(a)$$

There are 2 types of random variables that will be covered in this course, discrete and continuous (there are also absolutely continuous random variables, but they will not be covered)

1.4.1. Discrete random variables.

Definition/Sats 1.8: Discrete random variables

Consists of a finite or countable infinite set of numbers with probabilities:

- $P(X = x_i) = P(x_i) > 0$
- $P(X = \bigcup_{i=1}^{\infty} x_i) = 1$

Anmärkning:

If we have an uncountable infinite set of possibilities, the probability would be 0. Here is where continuous variables come to play

1.4.2. Continuous random variables.

For a continuous random variable, $F(x)$ is differentiable so that there exists a function f such that:

$$F(x) = \int_{-\infty}^x f(t)dt$$

From this comes 2 important things we can derive (both from the discrete and continuous case), namely expected value and the variance

Definition/Sats 1.9: Expected value

For discrete random variables, it is defined as

$$E(X) = \sum xF(x_i)$$

For the continuous case:

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx$$

Definition/Sats 1.10: Variance

$Var(X) = E(X - E(X))^2$ for both discrete and continuous random variables

An equivalent definition is $E(X^2) - (E(X))^2$

Anmärkning:

$E(X^2)$ is called the second moment

Anmärkning:

If the variance is small, we know that the random variable does not fluctuate a lot from the expected value.

If X is a random variable (r.v) with density function f and g is a function, then we can define a new random variable $g(X) = Y$

Y is a random variable with density function \hat{f} .

Then:

$$E(Y) = \int y\hat{f}(y)dy = E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx$$

And for the discrete random variable we have:

$$E(g(X)) = \sum g(x_i)p(x_i)$$

Anmärkning:

It is often better to use the definition of the density function for X rather than Y

Another remark worth noting is that $E(X^2)$ is a special case of $\int g(x)f(x)dx$

Example:

This example considers a distribution with no expected value (∞), and therefore it has no variance.

$P(X = k) = \frac{1}{k} - \frac{1}{k+1}$, this fulfills Kolmogorovs axioms, and

$$E(X) = \sum_{k=1}^{\infty} \frac{k}{k} - \frac{k}{k+1} = \sum_{k=1}^{\infty} \left(1 - \frac{k}{k+1}\right) = \sum_{k=1}^{\infty} \frac{1}{k+1} = \infty$$

2. TRANSFORMATIONS

2.1. Pre-knowledge.

Let X_i be independent random variables with the same mean value (expected value) μ and variance σ^2

Let $S_n = \sum_{i=1}^n X_i \xrightarrow{\text{Law of large numbers}} \frac{S_n}{n} \rightarrow \mu$ (convergence in probability).

Of course, this is assuming some sort of equal distribution.

The notation for convergence in probability is denoted by $Y_n \xrightarrow{P} a$. This follows from Markov's inequality. It is strongly suggested to look in the notes from the first course here.

From the law of large numbers, $S_n \approx n\mu$. But this does not take into account some errors that may take place in the $\frac{S_n}{n}$ side, as this does not affect the convergence to μ .

This is treated with the *Central Limit Theorem* (CLT), which says that $S_n \sim N(n\mu, n\sigma^2)$

This is equivalent to say that:

$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \approx N(0, 1)$$

When n is large.

Anmärkning:

Here we talked about *convergence in distribution*

3. MULTIVARIATE RANDOM VARIABLES

It is strongly suggested to recall the random n -dimensional vector from Probability theory 1.

It is interesting to look at that the distribution of the vector, but we are often interested in a function $g(X)$

Example:

Starting with 1-dimension, and then working our way up.

Let $Y = g(X)$. Suppose g is strictly increasing of X (larger values of $X \rightarrow g(X)$ is larger).

Then

$$Y \leq y \in \mathbb{R} \Leftrightarrow g(X) \leq y \in \mathbb{R} \Leftrightarrow X \leq g^{-1}(y) = h(y)$$

If we look at the distribution function (which gives us everything, the probability the everything):

$$F_Y(y) = P(X \leq g^{-1}(y)) = P(X \leq h(y)) = F_X(h(y))$$

Anmärkning:

The inverse function g^{-1} is denoted by h

From the chain rule for derivatives we get:

$$f_Y(y) = f_X(h(y)) \cdot h'(y)$$

We have gotten some information about X from Y . We can of course do the same for strictly decreasing functions:

$$\begin{aligned} Y \leq y &\Leftrightarrow X \geq h(y) \\ \Rightarrow F_Y(y) &= P(X \geq h(y)) = 1 - F_X(h(y)) \end{aligned}$$

This is good since we know that density functions are always positive. Taking the derivative gives us:

$$f_Y(y) = -f_X(h(y)) \cdot h'(y)$$

Anmärkning:

We assume g is differentiable with an inverse.

We showed that $f_Y(y) = f_X(h(y)) \cdot |h'(y)|$

3.1. Multivariate case.

Think of two n -dimensional space. One for $X = (X_1, X_2, \dots, X_n)$, and one for $Y = (Y_1, Y_2, \dots, Y_n)$

Suppose g is a bijective function such that g and g^{-1} are differentiable and let $Y = g(X) = (g_1(X), g_2(X), \dots, g_n(X))$ where the component $Y_i = g_i(X) = (X_1, X_2, \dots, X_n)$

We have $X = g^{-1}(Y) = h(Y)$ (same as in 1-dimensional case)

Definition/Sats 3.11: Transformation theorem

The density of Y is given by

$$f_Y(y_1, y_2, \dots, y_n) = f_X(h_1(y), h_2(y), \dots, h_n(y)) \cdot |J|$$

Where J is the Jacobian matrix

$$J = \left| \frac{d(x)}{d(y)} \right| = \begin{vmatrix} \frac{dx_1}{dy_1} & \dots & \frac{dx_1}{dy_n} \\ \vdots & \ddots & \vdots \\ \frac{dx_n}{dy_1} & \dots & \frac{dx_n}{dy_n} \end{vmatrix}$$

Anmärkning:

Transformation theorem corresponds to multivariate analysis change of variables

Bevis 3.1: Sketch of Transformation theorem

Let y_0 be a point in the Y -space. Choose an ε -ball C around y_0 . Then we can assume that f_Y is constant in C .

The probability that our random vector Y will happen in this region is given by

$$\Delta C \cdot (f_Y(y_0) - \varepsilon) \leq P(Y \in C) \leq \Delta C \cdot (f_Y(y_0) + \varepsilon)$$

Anmärkning: ΔC is the volume/area of C

In the X -space, there then is a region which consists of all x whose $g(x)$ belongs to C . Since g is bijective \Rightarrow injective, we have that $Y \in C \Leftrightarrow X \in D = g^{-1}(C)$

This means that these probabilities are the same

$$\left| f_Y(y_0) \cdot \Delta C - \int_D f_X(x) dx \right| \leq \Delta C \cdot \varepsilon$$

As C decreases, ΔC decreases as well as ε

Since g is a nice function, D will also decrease

We let $x_0 = g^{-1}(y_0) = h(y_0)$. We can replace the integral by $f_X(x_0) \cdot \Delta D$ and obtain

$$f_Y(y_0) \cdot \Delta C \approx f_X(x_0) \cdot \Delta D \Leftrightarrow f_Y(y_0) \approx f_X(x_0) \frac{\Delta D}{\Delta C}$$

We get equality when $C \rightarrow 0$ (choosing a smaller and smaller ε)

Recall the functional determinant (Jacobian) of the matrix $\frac{\Delta x}{\Delta y} = \left| \frac{d(x)}{d(y)} \right|$ (relative volume change)

Thus, we get $f_Y(y_0) = f_X(h(y_0)) ||J_n(x_0, y_0)||$

Since this is true for all y , we can take away the index y_0 , and we get:

$$f_Y(y) = f_X(h(y)) \cdot |J|$$

□

Example (1-dim case):

Suppose $g(X) = aX + b$. From this it is easy to see what the inverse function $h = g^{-1}$ is, namely $h(y) = \frac{y-b}{a}$.

The Jacobian is just $\frac{1}{a}$. By the transformation theorem we get

$$f_Y(y) = f_X\left(\frac{y-b}{a}\right) \cdot \left|\frac{1}{a}\right|$$

The main thing is that using the density function of X we can get the density function for Y .

Example 2.4:

X, Y are independent normally distributed random variables $N(0, 1)$, show that $X + Y$ and $X - Y$ are independent and determine their distribution function.

In order to solve this we do a variable substitution. Let $U = X + Y$ and $V = X - Y$.

Notice that $\frac{U+V}{2} = X$ and $\frac{U-V}{2} = Y$

We have our function $g(x, y) = (u, v) = (x + y, x - y)$

We have our inverse $g^{-1}(u, v) = (x, y) = \left(\frac{u+v}{2}, \frac{u-v}{2}\right)$

We can now use the transformation theorem:

$$f_{U,V}(u, v) = f_{X,Y}\left(\frac{u+v}{2}, \frac{u-v}{2}\right) \cdot |J|$$

The Jacobian can be found:

$$J = \begin{vmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{vmatrix} = \frac{-1}{4} - \frac{1}{4} = \frac{-1}{2}$$

By the transformation theorem:

$$f_{U,V}(u, v) = f_{X,Y}\left(\frac{u+v}{2}, \frac{u-v}{2}\right) \cdot \frac{1}{2} \stackrel{\text{indep.}}{=} f_X\left(\frac{u+v}{2}\right) f_Y\left(\frac{u-v}{2}\right) \cdot \frac{1}{2}$$

We know their density functions since they are normally distributed with $N(0, 1)$:

$$= \frac{1}{\sqrt{2\pi}} e^{-(1/2)\left(\frac{u+v}{2}\right)^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-(1/2)\left(\frac{u-v}{2}\right)^2} \cdot \frac{1}{2}$$

After simplification we get that $f_{U,V}$ is a product of one function of U and one function of V . This means that U, V are independent since we get them as a product of two different functions.

Example:

Recall the convolution formula from Probability theory 1; if X, Y are independent, we have

$$f_{X+Y}(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

We can show this from the transformation theorem.

By the independance of X, Y , we have $f_{X,Y}(x, y) = f_X(x) f_Y(y) \forall x, y$

Let $Z = X + Y$. Then $g(X, Y) = (X + Y, X)$

The inverse is given by $Y = Z - X$. We can now use the transformation theorem:

$$f_{Z,Y}(z, x) = f_{X,Y}(h_1(z, x), h_2(z, x)) \cdot |J|$$

The Jacobian is given by

$$\begin{vmatrix} 1 & -1 \\ 1 & 0 \end{vmatrix} = 1$$

Since X, Y are independent, we get:

$$f_{Z,X}(z, x) = f_Y(z-x) f_X(x)$$

The marginal density is given by integrating away x :

$$f_Z(z) = \int_{-\infty}^{\infty} f_{Z,X}(z,x)dx = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x)$$

3.2. Conditional Probabilities.

It is suggested to do some examples from Probability theory 1.

Let X, Y be some random variables with joint discrete distribution (**TODO:** Def).

We look at the conditional distribution:

$$p_{Y|X=x}(y) = P(Y = y|X = x) = \frac{P_{X,Y}(x,y)}{P_X(x)}$$

If we look at the conditional probability distribution function, we have:

$$F_{Y|X=x}(y) = \sum_{z \leq y} P_{Y|X=x}(z)$$

If we now look at if X, Y have a joint continuous distribution, we have something similar, with a couple of things changed where we use the density function instead (since some probabilities are 0)

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

Does this make sense for the continuous case? Well, suppose that $f_{Y|X=x_0}(y) = f_{X,Y}(x_0, y)$, would this be a natural definition?

The reason we cannot use this definition is because the probability that $X = x_0$ is very small because of continuous. But now we have this given, we have assumed this has happened, so $f_{X,Y}(x_0, y)$ could very well be too small compared to $f_{Y|X=x_0}(y)$ where we already know that $X = x_0$ happens.

If we instead put $f_{Y|X=x_0}(y) = K f_{X,Y}(x_0, y)$ some constant to compensate, then we need to check if the properties for density functions are preserved:

$$\int_{-\infty}^{\infty} f_{Y|X=x_0}(y)dy = 1 \Leftrightarrow \text{proper density function}$$

However, with the K in front, we get:

$$\begin{aligned} \int_{-\infty}^{\infty} f_{Y|X=x_0}(y)dy &= 1 = \int_{-\infty}^{\infty} K f_{X,Y}(x_0, y)dy \\ &= K f_X(x_0) = 1 \\ \Rightarrow K &= \frac{1}{f_X(x_0)} \end{aligned}$$

This is true for all $X = x_0$, so the formula we have seems to be correct!

Just as in the discrete case, we have:

$$F_{Y|X=x}(y) = \int_{-\infty}^y f_{Y|X=x}(z)dz$$

Now we have all the tools to start define conditional expectations and conditional variances.

3.3. Conditional expected values and variances.

There are some natural definitions

Definition/Sats 3.12: Conditional Expected Value

In the continuous case we have:

$$\mathbb{E}(Y|X = x) = \int_{-\infty}^{\infty} y f_{Y|X=x}(y) dy$$

In the discrete case we have:

$$\mathbb{E}(Y|X = x) = \sum_{-\infty}^{\infty} P_{Y|X=x}(y)$$

Notice that $\mathbb{E}(Y|X = x)$ is a function of x .

We now look at only the random variable $\mathbb{E}(Y|X) = g(X)$

Sometimes it is easier to look at the conditional expected value.

Definition/Sats 3.13

$\mathbb{E}(Y|X)$ has the same expected value as $\mathbb{E}(Y)$:

$$\mathbb{E}(\mathbb{E}(Y|X)) = \mathbb{E}(g(X)) = \mathbb{E}(Y)$$

In the discrete case, we see a variant of the law of total probabilities:

$$\mathbb{E}(Y) = \mathbb{E}(Y|X) = \sum_x \mathbb{E}(Y|X = x) P(X = x)$$

Rules for calculations, as in the unconditional case, can be found in the book (**TODO:** *Ins*)

$$\mathbb{E}(f(X)Y|X) = f(X)\mathbb{E}(Y|X)$$

Another natural rule is if X, Y are independent:

$$\mathbb{E}(Y|X) = \mathbb{E}(Y)$$

Definition/Sats 3.14: Conditional variance

$$v(x) = \text{Var}(Y|X = x) = \mathbb{E}((Y - \mathbb{E}(Y|X = x))^2 | X = x)$$

$$V(X) = \text{Var}(Y|X) = \mathbb{E}((Y - \mathbb{E}(Y|X))^2 | X)$$

Definition/Sats 3.15

We define $e(X) = \mathbb{E}(Y|X)$ and $V(X) = \text{Var}(Y|X)$ and assume $g(X)$ is some function on X . Then we have:

$$\mathbb{E}((Y - g(X))^2) = \mathbb{E}(V(X)) + \mathbb{E}((e(X) - g(X))^2)$$

Bevis 3.2

$$\begin{aligned}
\mathbb{E}((Y - g(X))^2) &= \mathbb{E}((Y - e(X) + e(X) - g(X))^2) \\
&= \mathbb{E}((Y - e(X))^2) + 2\mathbb{E}(Y - e(X))\mathbb{E}(e(X) - g(X)) + \mathbb{E}((e(X) - g(X))^2) \\
&= \mathbb{E}(\mathbb{E}((Y - e(X))^2|X)) + 2\mathbb{E}(\mathbb{E}(Y - e(X))\mathbb{E}(e(X) - g(X))|X) + \mathbb{E}((e(X) - g(X))^2) \\
&= \mathbb{E}(v(X)) + 2(e(X) - g(X))\mathbb{E}(Y - e(X)|X) + \mathbb{E}((e(X) - g(X))^2) \\
&= \mathbb{E}(v(X)) + 2(e(X) - g(X)) \underbrace{(e(X) - e(X))}_{=0} + \mathbb{E}((e(X) - g(X))^2)
\end{aligned}$$

In the middle term we used:

$$\mathbb{E}(\mathbb{E}(Y - e(X)|X)) = \mathbb{E}(e(X) - e(X)|X) = e(X) - e(X) = 0$$

□

There is a nice corollary that follows from this:

Definition/Sats 3.16

$$Var(Y) = \mathbb{E}(Var(Y|X)) + Var(\mathbb{E}(Y|X))$$

The proof of the corollary follows from Proof 3.2 by choosing $g(X) = \mathbb{E}(Y)$. Then the last theorem gives the result directly:

$$\mathbb{E}((Y - e(X))^2) = \mathbb{E}(v(X)) + \mathbb{E}((e(X) - e(Y))^2)$$

But $e(X)$ has the same expected value as Y , as in $\mathbb{E}(Y) = \mathbb{E}(e(X))$, so we get:

$$\mathbb{E}(v(X)) + \mathbb{E}((e(X) - e(Y))^2) = \mathbb{E}(v(X)) + Var(\mathbb{E}(Y|X))$$

Recall that random variables do not only have values attached to them, but also parameters. For example $X \sim N(\mu, \sigma^2)$ where μ, σ^2 are parameters.

For example, if $X \sim N(\mu, \sigma^2)$ and $X = x$ is an outcome of this random variable, then both $x, \mu \in \mathbb{R}$ while X is a random variable.

Sometimes we want to think of the parameters as random variables.

Example:

Suppose we go to a hospital to take a blood-test and count the number of red blood cells in the sample, then we will get some value which also partly depends on some randomness.

First we look at some individual (even if I go to the hospital several times, each time will give different results).

This seemingly random variation from the same person can be explained by the Poisson-distribution with some parameter m .

This means that if X is an observed value, $X \sim Po(m)$. The value m can be different across people.

We can think that we do a random trial in 2 steps:

- Choose a random individual to take the blood test from
- Count the amount of red blood cells from that individual

Then we can let X be the observed value for that person, and we have $X|M = m \sim Po(m)$ with M having some distribution (does not need to have Poisson distribution)

Example:

By the law of total probability, we can look at:

$$P(A) = \int_{-\infty}^{\infty} P(A|M = x)f_M(x)dx$$

Suppose now that $M \sim \text{Exp}(1)$. We then get:

$$\begin{aligned} P(X = k) &= \int_{-\infty}^{\infty} P(X = k|M = x)f_M(x)dx \\ &= \int_0^{\infty} e^{-x} \frac{x^k}{k!} e^{-x} dx = \frac{1}{k!} \int_0^{\infty} x^k e^{-2x} dx \\ &= \frac{1}{k!} \int_0^{\infty} \frac{y^k}{2^k} e^{-y} \frac{1}{2} dy = \frac{1}{k! + 2^{k+1}} \int_0^{\infty} y^k e^{-y} dy \\ &\Rightarrow k! \Rightarrow P(X = k) = \frac{1}{2^{k+1}} \end{aligned}$$

So X has a geometric distribution. One can also get this from the Γ distribution since it is very similar:

$$\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx$$

When $z = k$ then $\Gamma(z) = (k-1)!$

In the last example, we do not need to calculate the unconditional distribution of X if we just want to know the expected value/variance of X using the formulas that we proved above.

Recall that we can write $\mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|M))$. We said that $X|M \sim \text{Po}(m)$, and we know that a Poisson random variable has $\mathbb{E}(X|M) = M$, and we get:

$$\mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|M)) = \mathbb{E}(M)$$

Also remember that $M \sim \text{Exp}(1)$, so $\mathbb{E}(M) = 1$, and we have:

$$\mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|M)) = \mathbb{E}(M) = 1$$

By the corollary, we can also find the variance:

$$\text{Var}(X) = \mathbb{E}(\text{Var}(X|M)) + \underbrace{\text{Var}(\mathbb{E}(X)|M)}_{\text{Var}(X|M) = M} = \mathbb{E}(M) + \underbrace{\text{Var}(M)}_{= m = 1} = 1 + 1 = 2$$

Example:

Suppose we are in a coffee shop, and every customer has a choice between coffee (with probability p) or tea (with probability $1 - p$).

Suppose the number of customers during lunchtime is $\sim \text{Po}(\lambda)$ distributed. We want to count the number of coffees ordered in total (or rather, find the distribution of the number of coffees).

We proceed by letting X be number of coffees ordered, and let N be the number of customers. We actually know because of how we assumed it, we know the amount of customers, we have a binomial distribution:

$$X|N = n \sim \text{Bin}(n, p)$$

From our example we also know that $N \sim \text{Po}(\lambda)$, we have $P(X = k|N = n) = \binom{n}{k} p^k q^{n-k}$.

Now we want to count $P(X = k)$ without $N = n$. We can use the law of large probabilities:

$$\begin{aligned} P(X = k) &= \sum_{n=0}^{\infty} P(X = k|N = n)P(N = n) \\ &= \sum_{n=0}^{\infty} \binom{n}{k} p^k q^{n-k} e^{-\lambda} \frac{\lambda^n}{n!} \end{aligned}$$