# Introduction to probability 2
# Exercise: the Monty Hall problem

Luca <pgl@portamana.org>

## Let's make a deal!

The 'Monty Hall problem', inspired by the TV show *Let's make a deal!* hosted by Monty Hall, was proposed in the *Parade* magazine in 1990 (Lo Bello 1991; I changed the numbers of the doors):

*Suppose you are on a game show and given a choice of three doors. Behind one is a car; behind the others are goats. You pick door No. 1, and the host, who knows what is behind them [and wouldn't open the door with the car], opens No. 2, which has a goat. He then asks if you want to pick No. 3. Should you switch?*



In this exercise we'll try to answer this puzzle using the probability calculus discussed in the lecture notes[1]. First we'll translate the initial information and the question in the language of probability, then we'll find the numerical values of the requested probabilities, and finally we'll check how the answer changes if the initial information is slightly different. Try to solve as many of the queries below, step by step, as your time allows you to. At least skim through the whole exercise to understand the procedure we're following. I'd be happy if you at least tried Query 8, the most important, sometime in the future.

---

[1] https://portamana.org/linko.htm?w=introprob2.pdf

## Query 0:   Intuition

First of all *examine what your intuition tells you the answer should be*, without spending too much time thinking, just as if you were on the game show. Examine which kind of heuristics your intuition uses. If you already know the solution to this puzzle, try to remember what your intuition told you the first time you faced it. Keep your observations in mind for later on.

Now put your intuition aside, close its eyes, and unsheathe your analytic skills. Let the probability calculus be your compass.

## Query 1:   Statements

*Write down the* **statements** *that we have to use in order to analyse and solve the puzzle.*

If the present step is unclear or difficult, take a peek at the answer on the next page. Make sure your answer matches the one given there before proceeding to the next query.

**Answer**

The following *statements* turn out to be enough:

$$k := [\text{The general } k\text{nowledge provided with the puzzle}]$$

$$car1 := \text{"The car is behind door 1"}$$

$$car2 := \text{"The car is behind door 2"}$$

$$car3 := \text{"The car is behind door 3"} \tag{1}$$

$$host2 := \text{"The host opens door 2"}$$

$$you1 := \text{"You initially pick door 1"}$$

We can also consider additional statements, like these:

$$host1 := \text{"The host opens door 1"}$$
$$host3 := \text{"The host opens door 3"} \tag{2}$$

but they aren't used in the solution of the problem.

The symbols I chose ('*car1*', '*host2*', etc.) are of course unimportant. We could for example use '*C1*' instead of '*car1*', and so on.

## Query 2:   Initial probabilities

*Translate the information given in the problem into initial probabilities having definite numerical values.* These probabilities involve the statements found in the previous query.

If the present step is unclear or difficult, take a peek at the answer on the next page. Make sure your answer matches the one given there before proceeding to the next query.

**Answer**

The initial information can be summarized in the following probabilities.

We're assuming that you're initially equally uncertain about where the car is, and that your initial pick of a door doesn't change your uncertainty:

$$P(car_1 \mid k) = P(car_1 \mid you_1\ k) = 1/3,$$
$$P(car_2 \mid k) = P(car_2 \mid you_1\ k) = 1/3, \qquad (3)$$
$$P(car_3 \mid k) = P(car_3 \mid you_1\ k) = 1/3.$$

The host won't open your door and won't open the door with the car. If the car is behind door No. 1 he has a choice between No. 2 and 3. His options are expressed by these probabilities:

$$P(host_1|car_1\ you_1\ k) = 0,$$
$$P(host_2|car_1\ you_1\ k) = 1/2, \qquad P(host_3|car_1\ you_1\ k) = 1/2. \qquad (4a)$$

The 1/2 in the last two probabilities expresses that we don't know which door the host would choose to open, if he has a choice between two.

If the car is behind door No. 2, his options are more limited:

$$P(host_1|car_2\ you_1\ k) = 0, \quad P(host_2|car_2\ you_1\ k) = 0,$$
$$P(host_3|car_2\ you_1\ k) = 1. \qquad (4b)$$

Similarly if the car is behind door No. 3:

$$P(host_1|car_3\ you_1\ k) = 0, \quad P(host_2|car_3\ you_1\ k) = 1,$$
$$P(host_3|car_3\ you_1\ k) = 0. \qquad (4c)$$

No matter what you or the host do, the car is surely behind one of the doors, and it can't be behind more than one door:

$$P(car_1 \lor car_2 \lor car_3 \mid \ldots\ k) = 1,$$
$$P(car_1\ car_2 \mid \ldots\ k) = P(car_1\ car_3 \mid \ldots\ k) = P(car_2\ car_3 \mid \ldots\ k) = 0, \quad (5)$$
$$P(car_1\ car_2\ car_3 \mid \ldots\ k) = 0.$$

## Query 3:   What is the question?

*Translate the question of the problem into probabilities which we want to calculate.*

If the present step is unclear or difficult, take a peek at the answer on the next page. Make sure your answer matches the one given there before proceeding to the next query.

**Answer**

To decide whether we should switch to door No. 3, we must first calculate the probability that the car is behind that door and the probability that it is behind the door we picked, No. 1, *considering all the information we've gathered* – especially the host's choice. We then compare these probabilities:

$$P(car_1 \mid host_2 \; you_1 \; k) = ? \qquad P(car_3 \mid host_2 \; you_1 \; k) = ? \qquad (6)$$

## Query 4:   Solution via Bayes's theorem

*Find the numerical values of the probabilities identified in the previous query. Use Bayes's theorem in the form*

P(*hypothesis1* | *data info*) =

$$\frac{\text{P}(data \mid hypothesis1\ info) \times \text{P}(hypothesis1 \mid info)}{\text{P}(data \mid hyp.1\ info) \times \text{P}(hyp.1 \mid info) + \text{P}(data \mid hyp.2\ info) \times \text{P}(hyp.2 \mid info)}$$

(7)

*and similarly for hypothesis2. What are the 'hypotheses' and the 'data' in this problem?*

If the present step is unclear or difficult, take a peek at the answer on the next page. Make sure your answer matches the one given there before proceeding to the next query.

**Answer**

We have two hypotheses: $car1 :=$ "the car is behind door No. 1", and $car3 :=$ "the car is behind door No. 3". Our data is $host2$: the fact that the host chose door No. 2. Our initial information is the general information of the puzzle, $k$, and the fact that we initially picked No. 1, $you1$. (If you're wondering 'why is $you1$ not part of the data?', read below.)

Bayes's theorem (7) for hypothesis $car1$ in our case becomes

$$P(car1 \mid host2\ you1\ k) =$$

$$\frac{P(host2 \mid car1\ you1\ k) \times P(car1 \mid you1\ k)}{\left[\begin{array}{l} P(host2 \mid car1\ you1\ k) \times P(car1 \mid you1\ k) + \\ \qquad P(host2 \mid car3\ you1\ k) \times P(car3 \mid you1\ k) \end{array}\right]} \tag{8a}$$

and using the numerical values of the initial probabilities (3)–(5) we find

$$P(car1 \mid host2\ you1\ k) = \frac{\frac{1}{2} \times \frac{1}{3}}{\frac{1}{2} \times \frac{1}{3} + 1 \times \frac{1}{3}} = \frac{1}{3}. \tag{8b}$$

Bayes's theorem for hypothesis $car3$ becomes

$$P(car3 \mid host2\ you1\ k) =$$

$$\frac{P(host2 \mid car3\ you1\ k) \times P(car3 \mid you1\ k)}{\left[\begin{array}{l} P(host2 \mid car1\ you1\ k) \times P(car1 \mid you1\ k) + \\ \qquad P(host2 \mid car3\ you1\ k) \times P(car3 \mid you1\ k) \end{array}\right]} \tag{9a}$$

and replacing the numerical values of the initial probabilities we find

$$P(car3 \mid host2\ you1\ k) = \frac{1 \times \frac{1}{3}}{\frac{1}{2} \times \frac{1}{3} + 1 \times \frac{1}{3}} = \frac{2}{3}. \tag{9b}$$

**That is, the car is more likely to be behind door No. 3!** *We should therefore switch.*

Let's see if Bayes's theorem correctly finds also the obvious: it's impossible that the car is behind door No. 2:

$P(car2 \mid host2\ you1\ k) =$

$$\frac{P(host2 \mid car2\ you1\ k) \times P(car2 \mid you1\ k)}{\begin{bmatrix} P(host2 \mid car1\ you1\ k) \times P(car1 \mid you1\ k)\ + \\ P(host2 \mid car2\ you1\ k) \times P(car2 \mid you1\ k)\ + \\ P(host2 \mid car3\ you1\ k) \times P(car3 \mid you1\ k) \end{bmatrix}}$$

(10a)

and substituting the initial probabilities:

$$P(car2 \mid host2\ you1\ k) = \frac{0 \times \frac{1}{3}}{\frac{1}{2} \times \frac{1}{3} + 0 \times \frac{1}{3} + 1 \times \frac{1}{3}} = 0, \qquad (10b)$$

a reassuring result.

In choosing what our 'data' are, you may have asked yourself the following question: Should my initial door pick, *you1*, be included in the 'data'? It's a great question. Your own door choice came as no surprise to you, so it seems more reasonable to consider it as part of the 'other info', as we did above.

***But the important point is this:*** there wouldn't be anything wrong in including your door pick among the 'data'. In fact, *Bayes's theorem would give us exactly the same numerical result even if we took 'host2 you1' as 'data'.* This equality is a result of the self-consistency of the probability calculus.

In this case Bayes's theorem for *car1* takes this form:

$P(car1 \mid host2\ you1\ k) =$

$$\frac{P(host2\ you1 \mid car1\ k) \times P(car1 \mid k)}{\begin{bmatrix} P(host2\ you1 \mid car1\ k) \times P(car1 \mid k)\ + \\ P(host2\ you1 \mid car3\ k) \times P(car3 \mid k) \end{bmatrix}}$$

(11)

and you see that it involves probabilities that we don't have yet, for example $P(host2\ you1 \mid car1\ k)$. We'd need to first calculate these probabilities from our initial ones (3)–(5), using the five probability rules.

It would be a bit of extra work (feel free to do this as an extra query). After this extra calculations you'd find that Bayes's formula (11) above actually simplifies to (8a)! So, a different choice of what's 'data' wouldn't be wrong but would be less convenient, leading to more calculations.

This is a great feature, though: Even if you do your analysis in a slightly roundabout way, the probability calculus will lead you to the correct answer anyway – provided you follow all its rules exactly. Again, this happens because the probability calculus is just an extension of logic.

Finally, note that whether we consider your door pick, *you1*, as part of the data or of the initial information, its specification is essential for solving the problem: if you had picked another door, the host's options would have been different.

## Query 5: Let's educate our intuition

Does the answer you just found agree with your initial intuition? Most people find the correct answer counter-intuitive. I did. If your initial intuition told you differently, try to educate it by examining the results from Bayes's formulae (8)–(10). There's no right or wrong answer. You can see my personal analysis on the next page. People's intuitions often work in different ways, so the only person who can educate your intuition is you.

### *My* **answer**

Where does the final difference between the credibilities of the two hypotheses *car1* and *car3* come from? Bayes's formulae (8) and (9) show that the only difference is in the plausibilities of the host's actions given the two hypotheses and given that you picked door No. 1:

$$\mathrm{P}(host2|car1\ you1\ k) = 1/2, \qquad \mathrm{P}(host2|car3\ you1\ k) = 1. \qquad (12)$$

The host's choice gave us *information* that led to a change in the credibilities of the two hypotheses. The most obvious piece of information given by the host's choice is that the car can't be behind door No. 2: we know that it can't as soon as the host starts to open that door. If the car had been there, the host couldn't have opened that door:

$$\mathrm{P}(host2|car2\ you1\ k) = 0. \qquad (13)$$

But this obvious, large piece of information – so important that it immediately makes the hypothesis *car2* impossible and excluded at the outset – is *not* used in Bayes's formulae (8) and (9). So these formulae are telling us that the host's action contains *additional*, subtler pieces of useful information.

In fact the probabilities (12) tell us that it's more likely that the host opens door No. 2 under the hypothesis *car3* (where it's the only possible action for him) than under the hypothesis *car1* (where he has two possible choices). The observation of *host2* therefore provides slightly stronger evidence for *car3* than for *car1*. Or we can say that *car3* has slightly more 'explanatory power' for *host2* than *car1* does. Since the two hypotheses were initially equally plausible, the new evidence now makes *car1* less plausible. Not impossible, just less plausible. This reasoning is a sort of softened version of the logical impossibility of hypothesis *car2*; it shows that the plausibility calculus is an extension of formal logic (take a look at Hailperin 1984; 1991; 1996).

When I first faced this puzzle I didn't think about these slightly different informational connections between *host2* on one side and *car1*, *car3* on the other. They were eclipsed by the much stronger logical connection between *host2* and *car2*: 'he opens that door – no car there then; that's all there is to it'. So what my intuition learned from the probability calculation is not to exult and stop searching just because a very strong and important piece of information is revealed: there may be

additional crumbs of information hiding around, and together they may lead to far stronger conclusions.[2] This is what happens, for example, with magnetic-resonance imaging: the probability calculus is able to gather and assemble from the signal so many numerous pieces of information, each one invisible to the naked eye, that the final frequency estimate is *orders of magnitude* better than obtained from a Fourier transform. I invite you to read pages 23–24 of Bretthorst (1988) for an insightful discussion of this phenomenon.

The strength of the probability calculus is that it automatically keeps every shred of information into account (unless we too hurriedly skip steps, deluding ourselves to have found everything there was to be found).
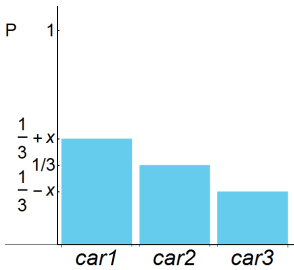
---

[2]As The Wolf concisely puts it in *Pulp Fiction*: 'Well, let's not start sucking each other's dicks quite yet' (https://www.youtube.com/watch?v=7zfbkhj8D0c, 01:00).

* * *

The solution we found depends on our *symmetric* states of ignorance regarding the placement of the car, formulae (3), and regarding the host's decision when he has a choice between two doors, formulae (4a). Let's examine how the credibilities change if we have some additional information.

## Query 6:   Inside information about the car

Let's consider our puzzle from the point of view of a different state of knowledge $k'$. You have a friend who works backstage. They secretly tell you (hey, that's cheating!) that they saw some large object – probably the car – being moved towards the left side (door No. 1). Because of this information you think it's more likely that the car is placed somewhat towards the left (door No. 1) than the right (door No. 3). Let's express this information by subtracting a positive amount $x$ from the initial probability of *car3* and giving it to *car1*:



$$P(car1 \mid k) = P(car1 \mid you1\ k') = \tfrac{1}{3} + x,$$
$$P(car2 \mid k) = P(car2 \mid you1\ k') = \tfrac{1}{3}, \qquad (14)$$
$$P(car3 \mid k) = P(car3 \mid you1\ k') = \tfrac{1}{3} - x,$$

so that we have a linear decrease from door No. 1 to door No. 3. Obviously $x \leqslant 1/3$, because we can't have negative probabilities.

In our previous calculation we saw that the host's opening of No. 2 gives more evidence to *car3* than to *car1*. Now we have initially more evidence for *car1* than for *car3*. It's possible that these two pieces of evidence balance each other out.

*Calculate for which value $x$ it doesn't matter whether you switch or not after the host opens No. 2.*

If the present step is unclear or difficult, take a peek at the answer on the next page. Then proceed to the next query.

**Answer**

The alternative state of knowledge $k'$ differs from $k$ only in the values of the initial probabilities for the car's position: formulae (14) instead of (3). Bayes's formulae (8)–(9) for the final probabilities for the car's position still apply, but with the new values of the initial probabilities. We have

$$\mathrm{P}(car1 \mid host2\ you1\ k') = \frac{\frac{1}{2} \times \left(\frac{1}{3} + x\right)}{\frac{1}{2} \times \left(\frac{1}{3} + x\right) + 1 \times \left(\frac{1}{3} - x\right)} \tag{15}$$

$$\mathrm{P}(car3 \mid host2\ you1\ k') = \frac{1 \times \left(\frac{1}{3} + x\right)}{\frac{1}{2} \times \left(\frac{1}{3} + x\right) + 1 \times \left(\frac{1}{3} - x\right)} \tag{16}$$

The query asks for which value of $x$ it doesn't matter whether we switch or not. This means that the final probabilities for *car1* and *car3*, (8) and (9), are equal. Let's therefore equate the two probabilities above. Note that the two fractions have the same denominator (which is different from zero), so we can just equate the numerators:

$$\frac{1}{2} \times \left(\frac{1}{3} + x\right) = 1 \times \left(\frac{1}{3} + x\right) \quad \implies \quad x = \frac{1}{9}. \tag{17}$$

This means that if the initial credibilities for the car's position are

$$\mathrm{P}(car1 \mid k) = \mathrm{P}(car1 \mid you1\ k') = 4/9,$$
$$\mathrm{P}(car2 \mid k) = \mathrm{P}(car2 \mid you1\ k') = 3/9, \tag{18}$$
$$\mathrm{P}(car3 \mid k) = \mathrm{P}(car3 \mid you1\ k') = 2/9,$$

then, after the host opens door No. 2, we are equally uncertain whether the car is behind No. 1 or No. 2, and so it doesn't matter whether we switch or not.

In general, if $x < 1/9$ we should switch because the credibility of *car3* is higher, and if $x > 1/9$ we should keep No. 1 because the credibility of *car1* is higher.

## Query 7:   Inside information about the host

In the previous query we had inside information about the car's position. Now let's consider yet another state of knowledge $k''$, in which we have inside information about the host instead. Your friend backstage secretly tells you that the host recently had a leg injury and feels some pain when walking. He still wants to present the show, but will limit his walking to a minimum. Since the host initially always stands close to door No. 1, this means that, given the choice to open door No. 2 or No. 3 (this happens when you've picked No. 1 and the car is there too) he will likely choose the closest: No. 2. This knowledge leads us to assign unequal probabilities for *host2* and *host3* conditional on *car1 you1*. Instead of the values (4), let's say that the plausibility that the host opens No. 2 if the car is behind No. 1 is $y$:

$$P(\textit{host1}|\textit{car1 you1 } k) = 0,$$
$$P(\textit{host2}|\textit{car1 you1 } k) = y, \qquad P(\textit{host3}|\textit{car1 you1 } k) = 1 - y. \tag{19}$$

This affects the values of the probabilities for the car after the host opens door No. 2.

*Calculate for which value of $y$ (if any) it doesn't matter whether you switch or not after the host opens No. 2.*

If the present step is unclear or difficult, take a peek at the answer on the next page. Make sure your answer matches the one given there before proceeding to the next – and final! – query.

**Answer**

The alternative state of knowledge $k''$ differs from $k$ only in the values of the probabilities for the host choice: instead of (4a) we have now (19). Bayes's formulae for the final probabilities of *car1* and *car3* still hold, but we now have the values

$$P(car1 \mid host2 \; you1 \; k'') = \frac{y \times \frac{1}{3}}{y \times \frac{1}{3} + 1 \times \frac{1}{3}}, \tag{20}$$

$$P(car3 \mid host2 \; you1 \; k'') = \frac{1 \times \frac{1}{3}}{y \times \frac{1}{3} + 1 \times \frac{1}{3}}. \tag{21}$$

These two probabilities are equal if

$$y + \frac{1}{3} = 1 \times \frac{1}{3} \quad \implies \quad y = 1. \tag{22}$$

This means that it doesn't matter whether we switch only if we're *absolutely certain* that the host would never open door No. 3 when he has the choice between that and No. 2 (he must be in a lot of pain!). Otherwise, it's still best to switch.

## Query 8:    Let's make a deal in your research!

Through the previous queries we've seen that the Monty Hall problem is just another example of calculating the probabilities of some hypotheses ('where's the car?') given some observations or data ('the host chose that specific door') and background information (the rules of the game and that you picked No. 1). The general steps we've taken here would apply identically in a scientific problem. The only difference would be in the number of hypotheses and in the determination of the initial probabilities.

*(a) Consider a problem of hypothesis comparison that you're facing in your research at the moment, or that you faced recently. Simplify it a little.*

*(b) Simplify the number of hypotheses to two or three.*

*(c) Write down the values of initial credibilities of these hypotheses (before you did your experiments); choose values that seem to correctly reflect your initial beliefs.*

*(d) Write down approximate/'toy' values of the probabilities of the result you obtained, conditional on each hypothesis; choose values that seem to correctly reflect the connection between hypothesis and result, but don't overthink too much.*

*(e) Calculate the values of the credibilities of the hypotheses in view of your result, using Bayes's theorem* (7) *with the probabilities you wrote down in the two steps above.*

Obviously you can't trust the result of this analysis in a quantitative (maybe not even qualitative) way, because the probabilities you wrote down in step (d) may not correctly reflect sensitive experimental details: you'd need to analyse that in much more detail, probably using numerical software.

Yet, you have just done a first rough Bayesian analysis of a concrete scientific problem.

## Optional query: the probability calculus

In the lecture[3] we saw that Bayes's theorem is simply a very convenient summary of a sequence of calculations that only involve the five probability rules. If you're curious about the step-by-step calculation, feel free to try it as we did during the lecture in the breast-cancer problem[3] (pdf page 98). Use the probability rules and shortcuts from the slides, and follow the roadmap shown in the next page. Solid red lines: 'and' rule; dashed blue lines: 'or' (∨) rule.

---

[3]https://portamana.org/linko.htm?w=introprob2.pdf

$P(car1|k) =$
$P(car1|you1\ k) = 1/3$

$P(car2|k) =$
$P(car2|you1\ k) = 1/3$

$P(car3|k) =$
$P(car3|you1\ k) = 1/3$

$P(car1|host2\ you1\ k)$

$P(host2\ car1|you1\ k)$   $P(host2\ car2|you1\ k)$   $P(host2\ car3|you1\ k)$

$P((host2\ car1)v(host2\ car2)v(host2\ car3)|you1\ k) =$
$P(host2\ (car1vcar2vcar3)|you1\ k]$

$P(host2|you1\ k)$

$P(car3|host2\ you1\ k)$

$P(car1vcar2vcar3|...\ k) = 1$

$P(host1|car1\ you1\ k) = 0$    $P(host1|car3\ you1\ k) = 0$
$P(host3|car1\ you1\ k) = 1/2$    $P(host3|car2\ you1\ k) = 0$
$P(host2|car1\ you1\ k) = 1/2$    $P(host2|car2\ you1\ k) = 0$

$P(host1|car2\ you1\ k) = 0$    $P(host1|car3\ you1\ k) = 0$
$P(host3|car2\ you1\ k) = 1$    $P(host3|car3\ you1\ k) = 0$
$P(host2|car2\ you1\ k) = 0$    $P(host2|car3\ you1\ k) = 1$

$P(car1\ car2\ car3|...\ k) = 0$

$P(car1\ car2|...\ k) = 0$
$P(car1\ car3|...\ k) = 0$
$P(car2\ car3|...\ k) = 0$

$P(host2\ car1\ car2\ car3)|?1\ k) =$
$P(host2\ car1\ host2\ car2\ host2\ car3)|?1\ k)$

$P(host2\ car1\ car2|you1\ k) =$
$P(host2\ car1\ host2\ car2\ car2|you1\ k)$

$P(host2\ car1\ car3|you1\ k) =$
$P(host2\ car1\ host2\ car3|you1\ k)$

$P(host2\ car2\ car3|you1\ k) =$
$P(host2\ car2\ host2\ car3|you1\ k)$

# Bibliography

Bretthorst, G. L. (1988): *Bayesian Spectrum Analysis and Parameter Estimation*. (Springer, Berlin). https://bayes.wustl.edu/glb/bib.html.

Hailperin, T. (1984): *Probability logic*. Notre Dame J. Formal Logic **25**[3], 198–212.

— (1991): *Probability logic in the twentieth century*. Hist. Philos. Logic **12**[1], 71–110.

— (1996): *Sentential Probability Logic: Origins, Development, Current Status, and Technical Applications*. (Associated University Presses, London).

Lo Bello, A. (1991): *Ask Marilyn: the mathematical controversy in Parade magazine*. Math. Gaz. **75**[473], 275–277.