

Sleeping beauty

Sleeping Beauty is told the following:

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Now suppose that you are sleeping beauty, and you are woken up from your sleep. You know the above, and you know that you are being awoken on Day 1 or on Day 2. What should you think is the chance that the coin flipped on Day 1 came up heads?

The argument for $\frac{1}{2}$ seems straightforward: Sleeping Beauty knows that the coin is fair, and so also knows that there is a $\frac{1}{2}$ chance that it comes up heads on any given throw, and a $\frac{1}{2}$ chance that it comes up tails. She has learned nothing which makes her doubt these probabilities for the Day 1 coin toss; so she should still estimate that there's a $\frac{1}{2}$ chance that the coin came up heads.

This involves some principle of the following sort:

If you estimate that the probability of some particular event occurring are N, then, if you learn nothing new relevant to the determination of the odds of that event occurring, you should stick with your estimate that the probability of its occurrence is N.

This principle --- which sums up the idea that you should only change your view about the probabilities of events in response to new information about the probabilities of those events --- seems almost too obvious to be worth stating.

However, there are powerful arguments for the conclusion that Sleeping Beauty should **shift** the probability she assigns to the coin being heads from $\frac{1}{2}$ before she is put to sleep, to $\frac{1}{3}$ after she is awoken.

Let's let "T1" abbreviate the proposition that the coin came up tails and it is Day 1, and "H1" abbreviate the corresponding proposition about the coin coming up heads and it being Day 1. Then the first argument can be stated as follows:

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Now suppose that you are sleeping beauty, and you are woken up from your sleep. You know the above, and you know that you are being awoken on Day 1 or on Day 2. What should you think is the chance that the coin flipped on Day 1 came up heads?

The argument for $\frac{1}{2}$ seems straightforward: Sleeping Beauty knows that the coin is fair, and so also knows that there is a $\frac{1}{2}$ chance that it comes up heads on any given throw, and a $\frac{1}{2}$ chance that it comes up tails. She has learned nothing which makes her doubt these probabilities for the Day 1 coin toss; so she should still estimate that there's a $\frac{1}{2}$ chance that the coin came up heads.

However, there are powerful arguments for the conclusion that Sleeping Beauty should **shift** the probability she assigns to the coin being heads from $\frac{1}{2}$ before she is put to sleep, to $\frac{1}{3}$ after she is awoken.

Let's let "T1" abbreviate the proposition that the coin came up tails and it is Day 1, and "H1" abbreviate the corresponding proposition about the coin coming up heads and it being Day 1. Then the first argument can be stated as follows:

1. $(P(T1 T1 \text{ or } T2)=P(T2 T1 \text{ or } T2))$	premise
2. $P(T1)=P(T2)$	(1)
3. $(P(H1 H1 \text{ or } T1)=\frac{1}{2})$	premise
4. $P(H1)=P(T1)$	(3)
5. $P(T1)=P(T2)=P(H1)$	(2,4)
C. $P(H1)=\frac{1}{3}$	(5)

This argument requires two assumptions, which are stated in premises (1) and (3). How would you argue for these assumptions, if you were trying to defend the argument?

Can you argue for (3) by changing the case so that the coin flip occurs **after** the Day 1 awakening?

Let's consider a second argument in favor of $\frac{1}{3}$.

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Let's consider a second argument in favor of $\frac{1}{3}$.

We can imagine changing the case so that whether or not the coin comes up heads, Sleeping Beauty is awakened **both** days. In this version, when Sleeping Beauty is awakened, her probability assignments should clearly be as follows:

$$P(T1)=P(T2)=P(H1)=P(H2)=\frac{1}{4}$$

since there's no reason to favor any of the possibilities over the others. But now suppose that we change the case slightly, so that if it is Day 2 and the coin toss was heads, soon after awakening you are told this fact. Suppose now that you are awoken, and that you are not told this. So you can rule out H2 as a possibility. What probability should you assign to the other possibilities? Well, it seems that you have learned only that one of four equiprobable theses is false, so you should maintain the view that

$$P(T1)=P(T2)=P(H1)$$

But then we can infer that

$$P(H1)=\frac{1}{3}.$$

But it seems that our modified case is the same as the original one: we can rule out, given the rules of the game, H2; and the other three possibilities seem equiprobable.

There's also a kind of intuitive argument for the conclusion that $P(H1)=\frac{1}{3}$. Sleeping Beauty would be reasonable to believe that, were this experiment performed over and over again, she would have twice as many tails-awakenings as heads-awakenings. So, given a random awakening, she should think that it is twice as likely that it be a tails-awakening as that it is a heads-awakening. So, she should think that the odds of heads having been thrown on any particular awakening of this sort is $\frac{1}{3}$. (Imagine us forcing Sleeping Beauty to bet on whether the coin came up heads on each awakening over a series of trials of the case. Wouldn't she stand to do much better if she adopted as a hypothesis to guide her betting that $P(H1)=\frac{1}{3}$?)

Nonetheless, it is hard to get rid of the intuition that when asked "What are the odds that a fair coin came up heads?", one should always say: $\frac{1}{2}$. **Especially given that upon being woken, we seem to learn nothing new - after all, we knew that we would be woken on one of the days.**

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Nonetheless, it is hard to get rid of the intuition that when asked “What are the odds that a fair coin came up heads?”, one should always say: $\frac{1}{2}$. **Especially given that upon being woken, we seem to learn nothing new - after all, we knew that we would be woken on one of the days.**

Some indication that this is more than an intuition is given by consideration of the **generalized sleeping beauty problem**.

Imagine that we vary the original case as follows: each time you would be awoken in that scenario, you have a $\frac{1}{100}$ chance of being awoken in the new version. So in this new version, when you are awoken, you do acquire genuinely new information: **you learn that you were awoken at least once**. In this case, how should you estimate the chances of heads versus tails?

It seems quite plausible that you should reason as follows: first, the probability that you will be woken up once, given that heads was flipped, is pretty clearly $\frac{1}{100}$.

Now consider the probability that you will be awoken at least once, given that tails came up: in that case, you get two chances at being woken up, so the probability is higher:

$$1 - \left(\frac{99}{100}\right)^2$$

this is, intuitively, because the chances of you not being woken up on either day is the square of the chances of you not being woken up on one day (i.e., $\frac{99}{100}$), and the chances of you being woken up once is $1 -$ the chances of you not being woken up on either day (since you will definitely either be woken at least once, or not woken either day, and not both).

So now imagine that you are woken up; you know then that you have been woken at least once. What are the odds that you should now assign to the coin having come up heads?

To figure out how to answer this question, an analogy may help: suppose that you think that the odds of ND winning the National Championship next year are $\frac{15}{100}$, and the odds of Purdue doing so are $\frac{1}{100}$. Suppose you are now told that next year the National Championship winner will come from the state of Indiana. (IU and Ball State have 0 chance.) What odds should you now assign to ND winning the National Championship?

A natural thought is: $\frac{15}{16}$, since we have now eliminated 84 of the 100 “possibilities.” Slight more formally, it seems that we should take it to be the odds of ND winning before learning the information about Indiana over the odds of ND winning plus the odds of Purdue winning (again, prior to our information about the NC winner coming from Indiana).

Let’s apply this to the generalized sleeping beauty.

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Some indication that this is more than an intuition is given by consideration of the **generalized sleeping beauty problem**.

Imagine that we vary the original case as follows: each time you would be awoken in that scenario, you have a 1/100 chance of being awoken in the new version. So in this new version, when you are awoken, you do acquire genuinely new information: **you learn that you were awoken at least once**. In this case, how should you estimate the chances of heads versus tails?

It seems quite plausible that you should reason as follows: first, the probability that you will be woken up once, given that heads was flipped, is pretty clearly 1/100.

Now consider the probability that you will be awoken at least once, given that tails came up: in that case, you get two chances at being woken up, so the probability is higher:

$$1 - \left(\frac{99}{100}\right)^2$$

this is, intuitively, because the chances of you not being woken up on either day is the square of the chances of you not being woken up on one day (i.e., 99/100), and the chances of you being woken up once is 1 - the chances of you not being woken up on either day (since you will definitely either be woken at least once, or not woken either day, and not both).

So now imagine that you are woken up; you know then that you have been woken at least once. What are the odds that you should now assign to the coin having come up heads?

Let's apply this to the generalized sleeping beauty.

In this case, you know that the coin came up either heads or tails, and you were woken; so the probability that the coin came up heads is presumably the probability of you being woken given heads divided by the probability of you being woken given heads + the probability of you being woken given tails, i.e.:

$$\frac{\frac{1}{100}}{\frac{1}{100} + 1 - \left(\frac{99}{100}\right)^2} = \frac{.01}{.0299} \approx .334$$

So this says that the probability you should assign to heads = just over 1/3. The interesting thing, though, is that the probability of heads, using this reasoning, steadily increases as the odds of you being woken up on the various occasions increases.

You are going to sleep for three days, during which time you will be woken up either once or twice. On Day 1, a fair coin is tossed. If that coin comes up heads, she will be woken **only** on Day 1, if tails then on Day 1 **and** on Day 2. If she is woken on Day 1, then she will be given a drug to put her back to sleep which also causes her to forget that awakening.

Some indication that this is more than an intuition is given by consideration of the **generalized sleeping beauty problem**.

Imagine that we vary the original case as follows: each time you would be awoken in that scenario, you have a 1/100 chance of being awoken in the new version. So in this new version, when you are awoken, you do acquire genuinely new information: **you learn that you were awoken at least once**. In this case, how should you estimate the chances of heads versus tails?

In this case, you know that the coin came up either heads or tails, and you were woken; so the probability that the coin came up heads is presumably the probability of you being woken given heads divided by the probability of you being woken given heads + the probability of you being woken given tails, i.e.:

$$\frac{\frac{1}{100}}{\frac{1}{100} + 1 - (\frac{99}{100})^2} = \frac{.01}{.0299} \approx .334$$

So this says that the probability you should assign to heads = just over 1/3. The interesting thing, though, is that the probability of heads, using this reasoning, steadily increases as the odds of you being woken up on the various occasions increases.

For example, suppose that there is a 99/100 chance, rather than a 1/100 chance, that you will be woken in the relevant occasions. Then, using the above reasoning, the probability which should be assigned to heads works out as follows:

$$\frac{\frac{99}{100}}{\frac{99}{100} + 1 - (\frac{1}{100})^2} = \frac{.99}{1.9899} \approx .498$$

And, in general, as the odds of being woken on the relevant occasions approaches 1, the probability which should be assigned to heads approaches 1/2 - which is strong evidence that the right response to the initial sleeping beauty problem is 1/2, not 1/3. If this is right, then something must be wrong with the three arguments for 1/3 which we have considered.

Or more, generally, we can say that **either** there is something wrong with what we have said about the generalized sleeping beauty, **or** there is something wrong with the arguments we have presented for 1/3. What you should think about is which of these you think is correct.

If the correct answer is 1/3 - and there is something wrong with the above argument about generalized sleeping beauty - there may be an interesting parallel between the two-envelope paradox and sleeping beauty.

If the correct answer is $\frac{1}{3}$ - and there is something wrong with the above argument about generalized sleeping beauty - there may be an interesting parallel between the two-envelope paradox and sleeping beauty.

Remember that we were discussing the relationship between the “open” and “closed” versions of the paradox, and I suggested that we only get the truly paradoxical results which follow from consideration of the closed versions if we adopt the following principle:

Inference from an unknown

Suppose that you are choosing between two actions, act 1 and act 2. It is always rational to do act 2 if the following is the case: there is truth about the situation which you do not know but which is such that, were you to come to know it, it would be rational for you to do act 2.

I suggested that there is something of a puzzle about how this principle could be false, but that rejecting it seems to give us a decent treatment of the various versions of the two-envelope paradox we discussed. Sleeping beauty - if $\frac{1}{3}$ is the right view - seems to be another case which leads us to reject this principle.

After all, if $\frac{1}{3}$ is the right answer, it seems that we can describe Beauty's situation as follows: before being put to sleep, she knows that she will be awoken; and she knows that upon learning that she has been awoken, she will be rational to judge that the odds of the coin having come up heads is $\frac{1}{3}$. But now, prior to being woken, she is not rational to judge that the odds that the coin will come up heads is $\frac{1}{3}$ - now, obviously, the right answer is $\frac{1}{2}$, despite the fact that she knows that she is about to undergo an experience which will lead her to revise that estimate.

However, sleeping beauty offers a further challenge. Even if we reject inference from an unknown, what is the relevant unknown? One wants to say something like: the proposition that I have been woken **today**. But suppose that this is day 2; is this different than the proposition that I was woken on day 2 (which, after all, I already knew would be the case)? How?