

Learning as conditional inference

The line between “reasoning” and “learning” is unclear in cognition. Just as reasoning can be seen as a form of conditional inference, so can learning: discovering persistent facts about the world (for example, causal processes or causal properties of objects). By saying that we are learning “persistent” facts we are indicating that there is something to infer which we expect to be relevant to many observations over time. Thus, we will formulate learning as inference in a model that (1) has a fixed latent value of interest, the *hypothesis*, and (2) has a sequence of observations, the *data points*.

When thinking about learning as inference, there are several key questions. First, what can be inferred about the hypothesis given a certain subset of the observed data? For example, in most cases, you cannot learn much about the weight of an object based on its color. However, if there is a correlation between weight and color – as is the case in many children’s toys – observing color does allow you to learn about weight.

Second, what is the relationship between the amount of input (how much data we’ve observed) and the knowledge gained? In psychology, this relationship is often characterized with a *learning curve*, representing a belief as a function of amount of data. In general, getting more data allows us to update our beliefs. But some data, in some models, has a much bigger effect. In addition, while knowledge often changes gradually as data is accumulated, it sometimes jumps in non-linear ways; these are usually the most psychologically interesting predictions.

Example: Learning About Coins

As a simple illustration of learning, imagine that a friend pulls a coin out of her pocket and offers it to you to flip. You flip it five times and observe a set of all heads:

[H, H, H, H, H].

Does this seem at all surprising? To most people, flipping five heads in a row is a minor coincidence but nothing to get excited about. But suppose you flip it five more times and continue to observe only heads. Now the data set looks like this:

[H, H, H, H, H, H, H, H, H, H, H]

Most people would find this a highly suspicious coincidence and begin to suspect that perhaps their friend has rigged this coin in some way – maybe it’s a weighted coin that always comes up heads no matter how you flip it. This inference could be stronger or weaker, of course, depending on what you believe about your friend or how she seems to act; did she offer a large bet that you would flip more heads than tails? Now you continue to flip five more times and again observe nothing but heads – so the data set now consists of 15 heads in a row:

在计算层面，这两点难以区分。
学习相比推理可能涉及的更多的是表征上的改变，或者是记忆。

```
[H, H, H, H, H, H, H, H, H, H, H, H, H, H, H]
```

Regardless of your prior beliefs, it is almost impossible to resist the inference that the coin is a trick coin.

This *learning curve* reflects a highly systematic and rational process of conditional inference. For simplicity let's consider only two hypotheses, two possible definitions of `coin`, representing a fair coin and a trick coin that produces heads 95% of the time. A priori, how likely is any coin offered up by a friend to be a trick coin? Of course there is no objective or universal answer to that question, but for the sake of illustration let's assume that the *prior probability* of seeing a trick coin is 1 in a 1000, versus 999 in 1000 for a fair coin.

x

```
var observedData = ['h', 'h', 'h', 'h', 'h']
```

```
var fairPrior = 0.999
```

```
var fairnessPosterior = Infer({method: 'enumerate'}, function() {
```

```
  var fair = flip(fairPrior)
```

```
  var coin = Bernoulli({p: fair ? 0.5 : 0.95})
```

```
  var obsFn = function(datum){observe(coin, datum == 'h')}
```

```
  mapData({data: observedData}, obsFn)
```

```
  return {fair: fair}
```

```
})
```

```
viz(fairnessPosterior)
```

run ▼

Try varying the number of flips and the number of heads observed. You should be able to reproduce the intuitive learning curve described above. Observing 5 heads in a row is not enough to suggest a trick coin, although it does raise the hint of this possibility: its chances are now a few percent, approximately 30 times the baseline chance of 1 in a 1000. After

observing 10 heads in a row, the odds of trick coin and fair coin are now roughly comparable, although fair coin is still a little more likely. After seeing 15 or more heads in a row without any tails, the odds are now strongly in favor of the trick coin.

When exploring learning as a conditional inference, we are particularly interested in **the dynamics of how inferred hypotheses change as a function of amount of data** (often thought of as time the learner spends acquiring data). We can map out the *trajectory* of learning by plotting a summary of the posterior distribution as a function of the amount of observed data. Here we plot the expectation that the coin is fair in the above example:

```
var fairnessPosterior = function(observedData) {  
  
  return Infer({method: 'enumerate'}, function() {  
  
    var fair = flip(0.999)  
  
    var coin = Bernoulli({p: fair ? 0.5 : 0.95})  
  
    var obsFn = function(datum){observe(coin, datum == 'h')}  
  
    mapData({data: observedData}, obsFn)  
  
    return fair  
  
  })  
  
}  
  
  
  
  
  
  
  
  
var trueWeight = 0.9  
  
var fullDataSet = repeat(100, function(){flip(trueWeight)?'h':'t'})  
  
var observedDataSizes = [1,3,6,10,20,30,50,70,100]  
  
var estimates = map(function(N) {  
  
  return expectation(fairnessPosterior(fullDataSet.slice(0,N)))  
  
}, observedDataSizes)
```

```
viz.line(observedDataSizes, estimates)
```

run ▼

Notice that different runs of this program can give quite different trajectories, but always end up in the same place in the long run. This is because **the data set used for learning is different on each run. This is a feature, not a bug: real learners have idiosyncratic experience, even if they are all drawn from the same distribution.** Of course, we are often interested in the average behavior of an ideal learner: we could average this plot over many randomly chosen data sets, simulating many different learners.

Study how this learning curve depends on the choice of `fairPrior`. There is certainly a dependence. If we set `fairPrior` to be 0.5, equal for the two alternative hypotheses, just 5 heads in a row are sufficient to favor the trick coin by a large margin. If `fairPrior` is 99 in 100, 10 heads in a row are sufficient. We have to increase `fairPrior` quite a lot, however, before 15 heads in a row is no longer sufficient evidence for a trick coin: even at `fairPrior` = 0.9999, 15 heads without a single tail still weighs in favor of the trick coin. This is because **the evidence in favor of a trick coin accumulates exponentially as the data set increases in size**; each successive `h` flip increases the evidence by nearly a factor of 2.

bits

从这个角度来说，
learning似乎就是有时
的(dynamic) reasoning

Learning is always about the shift from one state of knowledge to another. The speed of that shift provides a way to diagnose the strength of a learner's initial beliefs. Here, the fact that somewhere between 10 and 15 heads in a row is sufficient to convince most people that the coin is a trick coin suggests that for most people, the a priori probability of encountering a trick coin in this situation is somewhere between 1 in a 100 and 1 in 10,000—a reasonable range. Of course, if you begin with the suspicion that any friend who offers you a coin to flip is liable to have a trick coin in his pocket, then just seeing five heads in a row should already make you very suspicious—as we can see by setting `fairPrior` to a value such as 0.9.

Independent and Exchangeable Sequences

Now that we have illustrated the kinds of questions we are interested in asking of learning models, let's delve into the mathematical structure of models for sequences of observations.

If the observations have *nothing* to do with each other, except that they have the same distribution, they are called *identically, independently distributed* (usually abbreviated to i.i.d.). For instance the values that come from calling `flip` are i.i.d. To verify this, let's first check whether the distribution of two flips in a sequence look the same (are “identical”):

```
var genSequence = function() {return repeat(2, flip)}
```

```
var sequenceDist = Infer({method: 'enumerate'}, genSequence)
```

```
viz.marginals(sequenceDist)
```

run ▼

Now let's check that the first and second flips are independent, by conditioning on the first and seeing that the distribution of the second is unchanged:

```
var genSequence = function() {return repeat(2, flip)}
```

```
var sequenceCondDist = function(firstVal) {
```

```
  return Infer({method: 'enumerate'},
```

```
    function() {
```

```
      var s = genSequence()
```

```
      condition(s[0] == firstVal)
```

```
      return {second: s[1]};
```

```
    })
```

```
  }
```

```
viz(sequenceCondDist(true))
```

```
viz(sequenceCondDist(false))
```

run ▼

It is easy to build other i.i.d. sequences in WebPPL; we simply construct a stochastic thunk (a random function with no arguments) and evaluate it repeatedly. For instance, here is an extremely simple model for the words in a sentence:

```
var words = ['chef', 'omelet', 'soup', 'eat', 'work', 'bake', 'stop']
```

```
var probs = [0.0032, 0.4863, 0.0789, 0.0675, 0.1974, 0.1387, 0.0277]
```

```
var thunk = function() {return categorical({ps: probs, vs: words})};
```

```
repeat(10, thunk)
```

run ▼

In this example the different words are indeed independent: you can show as above (by conditioning) that the first word tells you nothing about the second word. However, constructing sequences in this way it is easy to accidentally create a sequence that is not entirely independent. For instance:

```
var words = ['chef', 'omelet', 'soup', 'eat', 'work', 'bake', 'stop']
```

```
var probs = (flip() ?
```

```
    [0.0032, 0.4863, 0.0789, 0.0675, 0.1974, 0.1387, 0.0277] :
```

```
    [0.3699, 0.1296, 0.0278, 0.4131, 0.0239, 0.0159, 0.0194])
```

```
var thunk = function() {return categorical({ps: probs, vs: words})};
```

```
repeat(10, thunk)
```

run ▼

While the sequence looks very similar, the words are not independent: learning about the first word tells us something about the `probs`, which in turn tells us about the second word. Let's show this in a slightly simpler example:

```
var genSequence = function() {
```

```
    var prob = flip() ? 0.2 : 0.7
```

```
    var thunk = function() {return flip(prob)}
```

```
    return repeat(2, thunk)
```

```
}
```

```
var sequenceCondDist = function(firstVal) {
```

```
    return Infer({method: 'enumerate'},
```

```
    function() {
```

```
var s = genSequence()
```

```
condition(s[0] == firstVal)
```

```
return {second: s[1]}
```

```
});
```

```
};
```

```
viz(sequenceCondDist(true))
```

```
viz(sequenceCondDist(false))
```

run ▼

Conditioning on the first value tells us something about the second. This model is thus not i.i.d., but it does have a slightly weaker property: it is **exchangeable** (https://en.wikipedia.org/wiki/Exchangeable_random_variables), meaning that **the probability of a sequence of values remains the same if permuted into any order**. When modeling learning it is often reasonable that the order of observations doesn't matter—and hence that the distribution is exchangeable.

It turns out that exchangeable sequences can always be modeled in the form used for the last example: **de Finetti's theorem** (https://en.wikipedia.org/wiki/De_Finetti%27s_theorem) says that, under certain technical conditions, **any exchangeable sequence can be represented as follows, for some** `latentPrior` **distribution and observation function** `f`:

exchangeable observations are conditionally independent relative to some latent variable.

```
var latent = sample(latentPrior)
var thunk = function() {return f(latent)}
var sequence = repeat(2,thunk)
```

Example: Polya's urn

A classic example is Polya's urn: Imagine an urn that contains some number of white and black balls. On each step we draw a random ball from the urn, note its color, and return it to the urn along with *another* ball of that color. Here is this model in WebPPL:

```
var urnSeq = function(urn, numsamples) {
```

```
  if(numsamples == 0) {
```

```
    return []
```

```
} else {
```

```
  var ball = uniformDraw(urn)
```

```
  return ball+urnSeq(urn.concat([ball]), numsamples-1)
```

```
}
```

```
}
```

```
var urnDist = Infer({method: 'enumerate'},
```

```
  function(){return urnSeq(['b', 'w'],3)})
```

```
viz(urnDist)
```

run ▼

Polya's urn is an examples of a “rich get richer” dynamic, which has many applications for modeling the real world. Examining the distribution on sequences, it appears that this model is exchangeable—permutations of a sequence all have the same probability (e.g., `bbw`, `bwb`, `wbb` have the same probability; `bww`, `wbw`, `wwb` do too). (Challenge: Can you prove this mathematically?) no, 直觉来说不对, 因为第一个抽到的会影响后面。

Because the distribution is exchangeable, we know that there must be an alternative representation in terms of a latent quantity followed by independent samples. The de Finetti representation of this model is:

```
var urn_deFinetti = function(urn, numsamples) {
```

```
  var numWhite = sum(map(function(b){return b=='w'},urn))
```

```
  var numBlack = urn.length - numWhite
```

```
  var latentPrior = Beta({a: numWhite, b: numBlack}) 假设你知道是beta分布, 但对于一个不知道这个模型的人来说并不知道。
```

```
  var latent = sample(latentPrior)
```

```
  return repeat(numsamples, function() {return flip(latent) ? 'b' : 'w'}).join("")
```



```
}
```

```
var urnDist = Infer({method: 'forward', samples: 10000},
```

```
function(){return urn_deFinetti(['b', 'w'],3)})
```

```
viz(urnDist)
```

run ▼

We sample a shared latent parameter – in this case, a sample from a Beta distribution – generating the sequence samples independently given this parameter. We obtain the same distribution on sequences of draws. (Challenge: show mathematically that these two representations give the same distribution.)

Ideal learners

Recall that we aimed to formulate learning as inference in a model that has **a fixed latent value of interest** and **a sequence of observations**. We now know that this will be possible anytime we are willing to assume the data are exchangeable.

Many Bayesian models of learning are formulated in this way. We often write this in the pattern of Bayes' rule:

```
Infer({...}, function() {  
  var hypothesis = sample(prior)  
  var obsFn = function(datum){...uses hypothesis...}  
  mapData({data: observedData}, obsFn)  
  return hypothesis  
});
```

The `prior` **samples a hypothesis from the *hypothesis space***. This distribution expresses our prior knowledge about how the process we observe is likely to work, before we have **observed any data**. The function `obsFn` **captures the relation between the hypothesis and a single datum**, and will usually contain an `observe` **statement**. (The marginal probability function for `obsFn` is called the *likelihood*. Sometimes `obsFn` itself is colloquially called the likelihood, too.) Here we have used the special operator `mapData` (<https://webppl.readthedocs.io/en/master/functions/arrays.html?highlight=mapData>) whose meaning is the same as `map`. We use `mapData` both to remind ourselves that we are expressing the special pattern of observing a sequence of observations, and because some inference algorithms can use this hint to do better learning.

Overall **this setup of prior, likelihood, and a sequence of observed data** (which implies an exchangeable distribution on data!) describes **an ideal learner**.

Example: Subjective Randomness

What does a random sequence look like? Is 00101 more random than 00000? Is the former a better example of a sequence coming from a fair coin than the latter? Most people say so, but notice that if you flip a fair coin, these two sequences are equally probable. Yet these intuitions about randomness are pervasive and often misunderstood: In 1936 the Zenith corporation attempted to test the hypothesis the people are sensitive to psychic transmissions. During a radio program, a group of psychics would attempt to transmit a randomly drawn sequence of ones and zeros to the listeners. Listeners were asked to write down and then mail in the sequence they perceived. The data thus generated showed no systematic effect of the transmitted sequence—but it did show a strong preference for certain sequences (Goodfellow, 1938 ([https://scholar.google.com/scholar?q="A%20psychological%20interpretation%20of%20the%20results%20of%20the%20Zenith%20experiment](https://scholar.google.com/scholar?q=))). The preferred sequences included 00101, 00110, 01100, and 01101.

Griffiths and Tenenbaum (2001) (<http://web.mit.edu/cocosci/Papers/random.pdf>) suggested that we can explain this bias if people are considering not the probability of the sequence under a fair-coin process, but the probability that the sequence would have come from a fair process as opposed to a non-uniform (trick) process:

```
var isFairDist = function(sequence) {  
  
  return Infer({method: 'enumerate'},  
  
    function () {  
  
      var isFair = flip() 先验是fair的概率50%，不fair的概率也是50%  
  
      var realWeight = isFair ? 0.5 : 0.2  
  
      var coin = Bernoulli({p: realWeight})  
  
      mapData({data: sequence}, function(d){observe(coin, d)})  
  
      return isFair  
  
    })  
  }  
  
  print("00101 is fair?")  
}
```

```
viz(isFairDist([false, false, true, false, true]))
```

```
print("00000 is fair?")
```

```
viz(isFairDist([false, false, false, false, false]))
```

run ▼

This model posits that **when considering randomness people are more concerned with distinguishing a “truly random” generative process from a trick process**. How do these inferences depend on the amount of data? Explore the learning trajectories of this model.

Learning a Continuous Parameter

The previous examples represent perhaps simple cases of learning. Typical learning problems in human cognition or AI are more complex in many ways. For one, learners are almost always **confronted with more than two hypotheses** about the causal structure that might underlie their observations. Indeed, **hypothesis spaces for learning are often infinite**. Countably infinite hypothesis spaces are encountered in models of learning for domains traditionally considered to depend on “discrete” or “symbolic” knowledge; hypothesis spaces of grammars in language acquisition are a canonical example. Hypothesis spaces for learning in domains traditionally considered more “continuous”, such as perception or motor control, are typically uncountable and parametrized by one or more continuous dimensions. In causal learning, both discrete and continuous hypothesis spaces typically arise. (In statistics, making conditional inferences over continuous hypothesis spaces given data is often called *parameter estimation*.)

We can explore a basic case of learning with continuous hypothesis spaces by slightly enriching our coin flipping example. Suppose instead of simply flipping a coin to determine which of two coin weights to use, we can choose *any* coin weight between 0 and 1. The following program computes conditional inferences about the weight of a coin drawn from a *prior distribution* described by the `Uniform` function, conditioned on a set of observed flips.

```
var observedData = ['h', 'h', 'h', 'h', 'h']
```

```
var weightPosterior = Infer({method: 'rejection', samples: 1000}, function() {
```

```
  var coinWeight = sample(Uniform({a: 0, b: 1}))
```

```
  var coin = Bernoulli({p: coinWeight})
```

```
  var obsFn = function(datum){observe(coin, datum == 'h')}  
}
```

```
mapData({data: observedData}, obsFn)
```

```
return coinWeight
```

```
})
```

```
viz(weightPosterior)
```

run ▼

Experiment with different data sets, varying both the number of flips and the relative proportion of heads and tails. How does the shape of the conditional distribution change? The location of its peak reflects a reasonable “best guess” about the underlying coin weight. It will be roughly equal to the proportion of heads observed, reflecting the fact that our prior knowledge is basically uninformative; a priori, any value of `coinWeight` is equally likely. The spread of the conditional distribution reflects a notion of confidence in our beliefs about the coin weight. The distribution becomes more sharply peaked as we observe more data, because each flip, as an independent sample of the process we are learning about, provides additional evidence of the process’s unknown parameters.

We can again look at the learning trajectory in this example:

```
var weightPosterior = function(observedData){
```

```
  return Infer({method: 'MCMC', samples: 1000}, function() {
```

```
    var coinWeight = sample(Uniform({a: 0, b: 1}))
```

```
    var coin = Bernoulli({p: coinWeight})
```

```
    var obsFn = function(datum){observe(coin, datum=='h')}
```

```
    mapData({data: observedData}, obsFn)
```

```
    return coinWeight
```

```
  })
```

```
}
```

```

var fullDataSet = repeat(100, function(){return 'h'})

var observedDataSizes = [0,1,2,4,8,16,25,30,50,70,100]

var estimates = map(function(N) {

  return expectation(weightPosterior(fullDataSet.slice(0,N)))

}, observedDataSizes)

viz.line(observedDataSizes, estimates)

```

run ▼

It is easy to see that this model doesn't really capture our intuitions about coins, or at least not in everyday scenarios. Imagine that you have just received a quarter in change from a store – or even better, taken it from a nicely wrapped-up roll of quarters that you have just picked up from a bank. Your prior expectation at this point is that the coin is almost surely fair. If you flip it 10 times and get 7 heads out of 10, you'll think nothing of it; that could easily happen with a fair coin and there is no reason to suspect the weight of this particular coin is anything other than 0.5. But running the above query with uniform prior beliefs on the coin weight, you'll guess the weight in this case is around 0.7. Our hypothesis generating function needs to be able to draw `coinWeight` not from a uniform distribution, but from some other function that can encode various expectations about how likely the coin is to be fair, skewed towards heads or tails, and so on.

One option is the **Beta distribution**. The Beta distribution takes parameters `a` and `b`, which describe the prior toward `true` and `false`. (When `a` and `b` are integers they can be thought of as *prior observations*.)

```

var pseudoCounts = {a: 10, b: 10};

var weightPosterior = function(observedData){

  return Infer({method: 'MCMC', burn:1000, samples: 1000}, function() {

    var coinWeight = beta(pseudoCounts)

    var coin = Bernoulli({p: coinWeight})

    var obsFn = function(datum){observe(coin, datum=='h')}
  })
}

```

```
mapData({data: observedData}, obsFn)
```

```
return coinWeight
```

```
})
```

```
}
```

```
var fullDataSet = repeat(100, function(){return 'h'});
```

```
var observedDataSizes = [0,1,2,4,6,8,10,20,30,40,50,70,100];
```

```
var estimates = map(function(N) {
```

```
  return expectation(weightPosterior(fullDataSet.slice(0,N)))
```

```
}, observedDataSizes);
```

```
viz.line(observedDataSizes, estimates);
```

run ▼

We are getting closer, in that learning is far more conservative. In fact, it is too conservative: after getting heads 100 times in a row, most humans will conclude the coin can *only* come up heads. The model, in contrast, still expects the coin to come up tails around 10% of the time.

We can of course decrease our priors [a](#) and [b](#) to get faster learning, but then we will just go back to our earlier problem. We would like instead to encode in our prior the idea that fair coins (probability 0.5) are much more likely than even moderately unfair coins.

A More Structured Hypothesis Space

The following model explicitly builds in the prior belief that fair coins are likely, and that all unfair coins are equally likely as each other:

```
var weightPosterior = function(observedData){
```

```
  return Infer({method: 'MCMC', burn:1000, samples: 10000}, function() {
```

```
    var isFair = flip(0.999)
```

```
    var realWeight = isFair ? 0.5 : uniform({a:0, b:1})
```

```

var coin = Bernoulli({p: realWeight})

var obsFn = function(datum){observe(coin, datum=='h')}

mapData({data: observedData}, obsFn)

return realWeight

})

}

var fullDataSet = repeat(50, function(){return 'h'});

var observedDataSizes = [0,1,2,4,6,8,10,12,15,20,25,30,40,50];

var estimates = map(function(N) {

return expectation(weightPosterior(fullDataSet.slice(0,N)))

}, observedDataSizes);

viz.line(observedDataSizes, estimates);

```

run ▼

This model stubbornly believes the coin is fair until around 10 successive heads have been observed. After that, it rapidly concludes that the coin can only come up heads. The shape of this learning trajectory is much closer to what we would expect for humans. This model is a simple example of **a hierarchical prior** which we explore in detail in a later chapter.

Example: Estimating Causal Power

Modeling beliefs about coins makes for clear examples, but it's obviously not a very important cognitive problem. However, many important cognitive problems have a remarkably similar structure.

For instance, a common problem for cognition is **causal learning**: from observed evidence about the co-occurrence of events, attempt to infer the causal structure relating them. An especially simple case that has been studied by psychologists is **elemental causal induction**: causal learning when there are only **two events**, **a potential cause C** and **a potential effect E**. Cheng and colleagues Cheng (1997) (<https://pdfs.semanticscholar.org/ac40/c59cc950959978c42fb0618b1458a93975a3.pdf>) have suggested assuming that C and background effects can both cause E, with a noisy-or interaction. Causal learning then

becomes an example of parameter learning, where the parameter is the “causal power” of C to cause E:

```
var observedData = [{C:true, E:true}, {C:true, E:true}, {C:false, E:false}, {C:true, E:true}]

var causalPowerPost = Infer({method: 'MCMC', samples: 10000}, function() {

  // Causal power of C to cause E

  var cp = uniform(0, 1)

  // Background probability of E

  var b = uniform(0, 1)

  var obsFn = function(datum) {

    // The noisy causal relation to get E given C

    var E = (datum.C && flip(cp)) || flip(b)

    condition( E == datum.E)

  }

  mapData({data: observedData}, obsFn)

  return {causal_power: cp}

});
```



```
viz(causalPowerPost);
```

run ▼

Experiment with this model: when does it conclude that a causal relation is likely (high `cp`)? Does this match your intuitions? What role does the background rate `b` play? What happens if you change the functional relationship in `obsFn`?

Reading & Discussion: Readings (</readings/learning-as-conditional-inference.html>)

Test your knowledge: Exercises (</exercises/learning-as-conditional-inference.html>)

Next chapter: 10. Learning with a language of thought (</chapters/lot-learning.html>)