

1 Chapter 1: Probabilities and Interference (Mackay Ch 2-3)

An ensemble: x random variable

$$\begin{aligned} A_x &= (a_1, a_2, \dots, a_n) \\ P_x &= (p_1, p_2, \dots, p_n) \\ p(x = a_i) &= p_i \end{aligned}$$

x takes value a_i with probability p_i

$$p \geq 0, \quad \sum_{a_i \in A_x} p(x = a_i) = 1$$

Short hand for $p(x = a_i)$ is $p(a_i)$, $p(x)$

Joint ensemble: X, Y ensembles

$$\begin{aligned} XY &= \text{ordered pairs}(x, y) \quad x \in A_X, y \in A_Y \\ P(x, y) &= \text{joint probability of } x \text{ and } y \end{aligned}$$

Marginal probability: $P(x, y) \rightarrow P(x), P(y)$

$$\begin{aligned} P(x) &= \sum_{y \in A_y} P(x, y) \\ P_x(x = a_i) &= \sum_{b \in A_y} P_{XY}(x = a_i, y = b) \end{aligned}$$

Conditional probability:

$$P(x = a_i | y = b_j) = \frac{P(x = a_i, y = b_j)}{P(y = b_j)}$$

“Probability of $x = a_i$ given that $y = b_j$ (is true)”

Example 1 $XY = 2$ successive letters in english alphabet. P_x and P_y are identical ‘frequency of a letter in english’

$$A_{xy} = \{aa, ab, ac, \dots, zz\}$$

$$P(y|x = 'q')$$

Peak at $y = 'u'$

$$\neq P_Y(y)$$

because x and y are not independent

X, Y “independent” if (and only if) $P(x, y) = P(x)P(y)$

Userful relations: $P(x, y) = P(x|y)P(y) = P(y|x)P(x)$

For any assumption H

$$\forall H : \quad P(x, y|H) = p(x, y|H)p(y|H)$$

‘Sum rule’:

$$P(x|H) = \sum_{y \in A_y} P(x, y|H) = \sum_{y \in A_y} P(x|y, H)P(y|H)$$

2 Lecture 1/18

Last time: Main point $P(y|x) \neq P(y)$

Useful relations: Conditional probability

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

where the joint relation is

$$P(x,y) = P(x|y)P(y) = P(y|x)P(x)$$

this can be rewritten into *Baye's theorem*

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Example 2: Apply Baye's theorem Alex is test for a nast disease.

- Disease status: a (sick or healthy)
- Test outcome: b (positive or negative)

"Test is 95% reliable" or

$$P(+|sick) = 0.95, \quad P(-|healthy) = 0.95$$

Disease is nasty but rare $P(sick) = 0.01$; $P(Healthy) = 0.99$

Test is positive, what is the probability that Alex is sick? $P(sick|+) = ?$

Solution Use Baye's theorem:

$$P(sick|+) = \frac{P(+|sick)P(sick)}{P(+)}$$

where $P(+)$ is the probability of a positive test result. This can be found using the sum rule

$$P(+) = P(+|sick)P(sick) + P(+|healthy)P(healthy)$$

Thus

$$P(sick|+) = \frac{0.95 * 0.01}{0.95 * 0.01 + 0.05 * 0.99} = 0.161$$

It is useful to write the probabilities in a table

	$b = +$	$b = -$	$P(b)$
$a = \text{sick}$	$0.95 * 0.01$	$0.05 * 0.01$	0.01
$a = \text{healthy}$	$0.05 * 0.99$	$0.95 * 0.99$	0.99
$P(a)$	0.161	0.839	1

where columns represent the 95:5 reliable test.

Exclam!

$$P(S|+) \neq P(+|S)$$

A brief philosophical interlude... The 'Bayesian viewpoint':

Probability as degree of beliefs in propositions given assumptions & evidence, or Probability as 'freq of outcomes in repeat random experiments'

Forward and inverse problems

So far we have talked about Cond Prob, Baye's thrm, and an example.

Generative Model: Parameters $\Theta \rightarrow P(D|\Theta) \rightarrow (P)$ outcomes (data) AKA 'forward problem' 'a model' predicts an outcome given parameters. The model is a probability distribution due to all the uncertainties and errors we have in the real world.

The Inverse Problem $P(\Theta|D)$

The inverse problem is the opposite of the forward problem (obviously). Also related to the issues regarding 'inference' and using Baye's theorem.

Example 3: A forward problem

An urn contains K balls, B balls are black, and $K - B$ balls are white. A ball is drawn at N times with replacement.

- $n_B = \#$ of times a black ball is drawn
- $P(n_B)$, average n_B ?, STD?

With

$$f_B = \frac{B}{K}$$

The probability is given by the binomial distribution

$$P(n_B|N, f_B) = \binom{N}{n_B} f_B^{n_B} (1 - f_B)^{N - n_B}$$

The mean is $N * f_B$ and the STD is $\sqrt{N * f_B * (1 - f_B)}$

Example 4: An inverse problem

We have 11 urns, each with 10 balls. u is the number of black balls in each urn and the urns have $u = 0, 1, \dots, 10$ black balls. Alex selects an urn at random and draws N balls at random with replacement. Bob wates Alex, but does not know which urn u was selected. For Bob, what is $P(u|N, n_B)$?

We have the data, but we are trying to infer the parameter u

Solution Use Baye's theorem

$$P(u|N, n_B) = \frac{P(n_B|u)P(u)}{P(n_B)}$$

where $P(n_B|u)$ is the 'forward' part from Ex 2, $P(u) = 1/11$, and $P(n_B)$ is the 'normalization' that makes it a valid prob. distribution:

$$P(n_B) = \sum_{u'} P(n_B|u')P(u')$$

Therefore

$$P(u|N, n_B) \propto \binom{N}{n_B} \left(\frac{u}{10}\right)^{n_B} \left(1 - \frac{u}{10}\right)^{N - n_B}$$

e.g. $n_B = 3, N = 10$

insert figure 1.2

The (0,0) point is impossible because we picked 3 black balls, and the urn $u = 0$ has no black balls. The same is true for the (10,10) point. The most likely point is $u = 3 \dots$

Exclam! This is known as ‘Posterior Probabilty’

- Θ is the parameter
- D is the data
- $P(\Theta)$ is the prior
- $P(D|\Theta)$ is the likelihood: a function of D prob of data given param (sums to 1 over all options for D). As a function of $\Theta \rightarrow$ likelihood of Θ
- $P(\Theta|D)$ is the posterior
- $P(D)$ is the normalization

! **Probability of *data***

! **Likelihood of *parameters***

Role of Prior:

! You can’t do inference without making assumptions

Lecture 1/23/24

Last time:

- Forward $p(\text{data}|\text{param})$
- Inverse $p(\text{param}|\text{data})$

Using Baye's theorem

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} = \frac{\text{likelihood} \cdot \text{prior}}{\text{norm}}$$

Note: You can't do inference w/o working assumptions (prior) priors are subjective. From the inverse problem ex from last week: what is the probability that next ball Alex draws is black?

$$P(B) = \sum P(u)P(B|u)$$

Note: Inference \neq decision/choice of model. Inference is assigning probabilities to hypotheses.

Problem USB Cable frustrations "It takes 3 tries to plug in a USB cable"

During our first try to plug in the cable, we are collecting data. And if its wrong, we 'believe' that the orientation is wrong, thus we flip it believing that the 2nd try is the correct one. But in fact, this is wrong and the 3rd try is the correct one.

How to collect data?

Lecture 1/25/24

3 Chapter 2: Probabilities and Interference (Mackay Ch 2-3)

Example 5: Tossing a coin

- 3 times: H, H, H
- 10 times: H, H, ... H

what is the probability of the next toss being H?

Ex 5.1 Coin with freq of heads f_H is tossed N times and n_H heads. What is the probability of the next toss being H? (Ex 4 but with fixed unknown parameter)

Prior: subjective assumption (e.g. could be uniform) then do inference.

Ex 5.2 N tosses, n_H heads. What is the probability that the coin is biased? (Model Comparison)

Lecture 1/30/24

Last time: Simple inference (within a model) where we solve for $p(data|param)$ and now we move on to model comparison!

Ch 2: Model Comparison Mackay Ch 3 & 28

A coin that is possibly bent has a frequency of heads f_H . For $N = 100$ tosses, $n_H = 90$ heads which is definitely a bent coin (biased).

For the case $N = 100$, $n_H = 55$, we are not sure if the coin is biased or not. The best fit to data is $f_H = 0.55$ we say that it is probably not bent from our intuition.

For the case $N = 10000$, $n_H = 5500$ we believe that the coin is more likely to be ‘bent’

Which model? We know that the fair coin model fits the model less than the bent coin model, but we believe that the fair coin model fits the data better than the bent coin model. From “Occam’s Razor” (simplicity): Accept the simplest explanation that fits the data. We would prefer the simpler fair coin model since it is simpler. This is merely a ad hoc rule of thumb. But Bayesian Calculus naturally implements Occam’s Razor.

Comparing hypothesis H_o (fair coin) and H_1 (bent coin) Warning! We should choose the hypothesis set before we see the data, otherwise it is cheating!

Big Picture Two levels of inference

- Level 1: Hypothesis set H_o with parameter f_H : Inferring $P(p_a) = ?$
- Level 1: Hypothesis set H_o no params: no inference
- Level 2: Hypothesis set H_o, H_1 : Inferring both $P(H_o)$ and $P(H_1)$

2.1 Coin tosses: 1-param model H_1 (L1 inference)

Outcomes: $X = \{a, b\}$ for heads and tails with probabilities p_a and $p_b = 1 - p_a$

Assumption: The prior on p_a is uniform

F Tosses: data = sequence, $s = aaba\dots$ with $F_a = \#$ of a’s and $F_b = \#$ of b’s; $F_a + F_b = F$

The model:

$$P(s|p_a, F, H_1) = p_a^{F_a} (1 - p_a)^{F_b}$$

since the tosses are a specific sequence e.g. aaba... From the definition of H_1

$$p_a \in [0 \dots 1]$$

is equiprobable and the prior tells us that $p(p_a) = 1$

Questions Given a sequence s of F observations, with $\# a = F_a$ and $\# b = F_b$,

1. What is my posterior belief about p_a ? or $P(p_a) = ?$
2. What is the probability that next draw is a ?

As this is an inverse problem, we use Bayes’s theorem

$$P(p_a|s, F, H_1) = \frac{P(s|p_a, F, H_1)P(p_a)|H_1}{P(s|F, H_1)}$$

the bottom takes the full probability of the data no matter the value of p_a and is the normalization

$$= \frac{p_a^{F_a}(1-p_a)^{F_b}(1)}{\int_0^1 p_a^{F_a}(1-p_a)^{F_b} dp_a}$$

where we use the sum rule for the denominator

$$\sum_{p_a} P(s|p_a, F, H_1) P(p_a|H_1)$$

but since it is a continuous variable, we use the integral instead of the sum. The math gives us the gamma function

$$\text{normalization factor} = \frac{F_a! F_b!}{(F_a + F_b)!}$$

Examples $s = aba$ vs $s = bbb$

$$P(p_a|s = aba) \propto p_a^2(1-p_a) \quad \text{vs} \quad P(p_a|s = bbb) \propto (1-p_a)^3$$

The first looks like a parabola and the second looks like a decaying cubic function. In each case, the most probable p_a is $2/3$ and 0 respectively which is shown by the data.

Probability of next toss is a We need to integrate over the prior to get the probability of the next toss being a .

$$P(\text{next} = a) = \int dp_a P(\text{next} = a|p_a) P(p_a|s, F, H_1) = \int dp_a P(p_a|s, F, H_1) p_a = \text{average of } p_a$$

the average of p_a for the first example is $3/5 = 0.6$ and for the second example is $1/5 = 0.2$

Conclusion: We found Probability of s given p_a and H_1 (Data given biased coin model) and the probability of p_a given s, F, H_1 (inference), or forward and inverse probabilities for the biased coin model H_1 .

2.2 Zero-parameter model H_o (Fair coin) & model comparison where $p_a = 1/2$. The forward probability is

$$P(s|H_o) = \frac{1}{2^F}$$

Question: Given a string of F observations, what comparison can we make between the biased coin model and the fair coin model, H_o vs H_1 ?

The Hypothesis space is now $\{H_o, H_1\}$ where only models are under consideration. Using Baye's theorem again

$$P(H_o|s, F) = \frac{P(s|F, H_o) P(H_o)}{P(s|F)}$$

and

$$P(H_1|s, F) = \frac{P(s|F, H_1) P(H_1)}{P(s|F)}$$

where $P(s|F) = \sum_{H \in \{H_o, H_1\}} P(s|F, H) P(H)$. looking at the ratio of the two probabilities

$$\frac{P(H_1|s, F)}{P(H_o|s, F)} = \frac{P(s|F, H_o) P(H_1)}{P(s|F, H_1) P(H_o)}$$

where the first fraction is what the data told us, and the second fraction is what we know before (prior).

Lecture 2/1/24

Last time: We discussed the zero-parameter model H_o (fair coin) and the one-parameter model H_1 (biased coin). We used Baye's theorem to compare the two models to find the ratio of the two probabilities

$$\mathcal{R} = \frac{P(H_1|s, F)}{P(H_o|s, F)} = \frac{P(H_1)}{P(H_o)} \frac{P(s|F, H_o)}{P(s|F, H_1)}$$

where we set no a priori model (prior) preference, so $P(H_1) = P(H_o) = 1/2$. So the ratio is

$$\mathcal{R} = \frac{P(s|F, H_o)}{P(s|F, H_1)} = \frac{\frac{F_a! F_b!}{(F_a + F_b + 1)!}}{\frac{1}{2^F}} = \frac{2^F F_a! F_b!}{(F + 1)!}$$

what does this plot look like? As the number of tosses goes to infinity, this ratio will go to the truth! Simulation is shown by Figure 3.1. where the the bent coin $p_a = 0.9$ probability goes to infinity as well

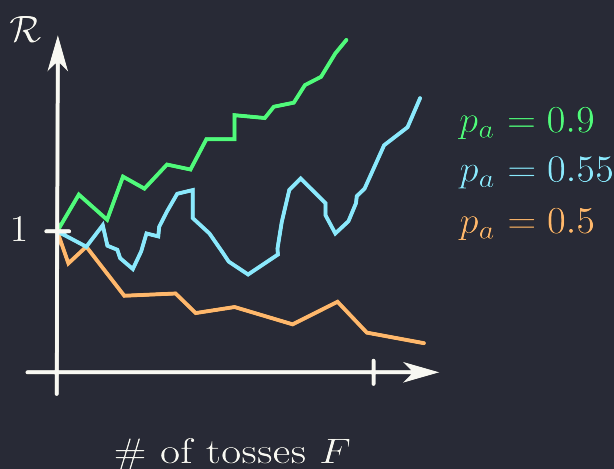


Figure 3.1: Ratio of the two probabilities as a function of the number of tosses

as the slightly biased coin (but at a slower pace) and the fair coin goes to zero. We know this from the probability

$$P(s|F, H_o) = \int_0^1 P(s|p_a, F, H_1) P(p_a|F, H_1) dp_a$$

NOTE: There exists a p_a that fits data better than H_o , but this evidence term includes averaging over p_a Bayes theorem in the context of model comparison

$$\text{bayes} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$$

TAKEHOME: Bayesian model comparison naturally includes Occam's Razor!

2.4 P-values? Why not just use p-values? e.g.

$$F = 250 \quad F_a = 141, F_b = 109$$

Do these data suggest that the coin is biased?

P-value: Probability to get data as extreme or more, assuming the null hypothesis is true.

- Null hypothesis: Coin is fair (H_0)
- Our hypothesis: Coin is biased (H_1)
- mean = $F/2$
- $\sigma = \sqrt{F}/2$
- Our observation: $\frac{F_a - F/2}{\sqrt{F}/2} = 2.02\sigma$
- p-value = $0.0497 < 0.05!!!!$

Google “a small p-value (< 0.05) indicates strong evidence against the null hypothesis so you reject it”

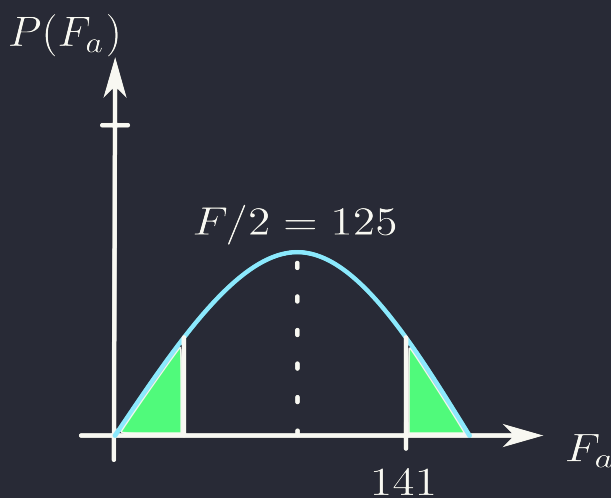


Figure 3.2: Finding p-value based on the Gaussian distribution

From sterling approximation

$$\ln(k!) \approx k \ln(k) - k + \dots$$

With uniform prior on p_a

$$\mathcal{R} = \frac{2^{250} 141! 109!}{251!} = 0.61$$

if anything, there is weak evidence *against* coin being biased.

Non-uniform priors? For a reasonable family of priors, across the entire set of priors, strongest evidence for bias is 2.5 : 1 (From Mackay) This differs from the p-value which is 20 : 1.

4 Chapter 3: Maximum Likelihood *Approximation*

(Ch 22 Mackay)

GOAL: Connect to the stat you may have seen before. Going back to Example 4 (Urns and more urns)

- Unknown u^* selected at random
- 10 draws (with replacement): 3 black

- $P(\text{next draw} = \text{black}) = ?$
- Most likely $u : 3 \rightarrow$ predicts 0.3
- Correct answer: predicts 0.33

but the two numbers are kinda similar...

NOTE: Bayesian model comparison, not model selection, but complete enumeration of hypotheses (integration over hyp space) is computationally expensive (especially in high dimensions)

e.g. Comparing 2 models:

- 1 Gaussian: 2 parameters μ, σ
- 2 Gaussian ($a_1 G_1 + a_2 G_2$): 5 parameters $\mu_1, \sigma_1, \mu_2, \sigma_2, a_1/a_2$

This problem of an increasing number of parameters motivates *Max likelihood (ML) approximation*: instead of enumeration, focus on 1 hypothesis that maximized the likelihood function.

Max Likelihood Estimation (MLE)

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

instead of [assuming prior \rightarrow compute posterior \rightarrow integrate over hyp space] we just [compute the likelihood function \rightarrow maximize it] (MLE).

3.1 A single Gaussian

- Data: $\{x_n\} \quad n = 1, \dots, N$
- model: these observations were sampled from a gaussian with probability

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

where we have 2 parameters μ, σ to determine.

Log likelihood (multiplying likelihoods is hard, adding log likelihoods is easier)

$$\begin{aligned} \ln P(\{x_n\}|\mu, \sigma) &= \sum_{n=1}^N \left(-\ln \sqrt{2\pi\sigma^2} - \frac{(x_n - \mu)^2}{2\sigma^2} \right) \\ &= -N \ln \sqrt{2\pi\sigma^2} - \frac{N}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \end{aligned}$$

Sufficient statistics: Denote

$$\hat{x} \equiv \sum_n \frac{x_n}{N} \quad \text{empirical mean}$$

$$S = \sum_n (x_n - \hat{x})^2 \quad \text{sum of square deviations}$$

These two numbers refer to the sufficient statistics. From these we get the log likelihood

$$\ln P = -N \ln \sqrt{2\pi\sigma^2} - \frac{N(\mu - \hat{x})^2 + S}{2\sigma^2}$$

Thus the max likelihood estimate of μ, σ are

$$\mu_{ML} = \hat{x}$$

$$\sigma_{ML} = \sqrt{\frac{S}{N}} = \sqrt{\frac{\sum_n (x_n - \hat{x})^2}{N}}$$

If σ is known, then $P(\mu)$ is a Gaussian we know that σ/\sqrt{N} is the width of the likelihood (error bars)