

# Statistical Rethinking: Notes and Selected Exercises

Somnath Sikdar

May 24, 2021

# Contents

<b>2</b>	<b>Small Worlds and Large Worlds</b>	<b>2</b>
<b>3</b>	<b>Generalized Linear Models</b>	<b>4</b>
3.1	Binomial Distributions and Maximum Entropy . . . . .	5
3.2	The Poisson and Related Families of Distributions . . . . .	7
3.2.1	Number of arrivals . . . . .	7

## Chapter 2

# Small Worlds and Large Worlds

**Exercise 2.1.** Suppose there are two species of panda bear. Both are equally common in the wild and live in the same places. They look exactly alike and eat the same food, and there is yet no genetic assay capable of telling them apart. They differ however in their family sizes. Species  $A$  gives birth to twins 10% of the time, otherwise birthing a single infant. Species  $B$  births twins 20% of the time, otherwise birthing singleton infants. Assume these numbers are known with certainty, from many years of field research. Now suppose you are managing a captive panda breeding program. You have a new female panda of unknown species, and she has just given birth to twins. What is the probability that her next birth will also be twins?

**Solution.** We have to estimate  $\Pr\{\text{twins again}|\text{twins before}\}$ . We may write this conditional probability as follows:

$$\begin{aligned}\Pr\{\text{twins again}|\text{twins before}\} &= \Pr\{\text{twins again}|A, \text{twins before}\} \cdot \Pr\{A|\text{twins before}\} + \\ &\quad \Pr\{\text{twins again}|B, \text{twins before}\} \cdot \Pr\{B|\text{twins before}\} \\ &= \Pr\{\text{twins again}|A\} \cdot \Pr\{A|\text{twins before}\} + \\ &\quad \Pr\{\text{twins again}|B\} \cdot \Pr\{B|\text{twins before}\}\end{aligned}\tag{2.1}$$

The last equality follows, since when we condition on any given species, the probability of having twins again is independent of whether twins were born before. Strictly speaking, this is also an assumption but a reasonable one. As such, we do not use the “before” and “again” qualifiers. To evaluate the last expression, we calculate  $\Pr\{\cdot|\text{twins}\}$  for both species  $A$  and  $B$ . Using Bayes’ Theorem, we obtain:

$$\begin{aligned}\Pr\{A|\text{twins}\} &= \frac{\Pr\{A|\text{twins}\} \cdot \Pr\{A\}}{\Pr\{A|\text{twins}\} \cdot \Pr\{A\} + \Pr\{B|\text{twins}\} \cdot \Pr\{B\}} \\ &= \frac{0.10 \times 0.5}{0.5 \times (0.1 + 0.20)} \\ &= \frac{1}{3}.\end{aligned}$$

A similar calculation yields  $\Pr\{B|\text{twins}\} = \frac{2}{3}$ . The final expression in (2.1) evaluates to

$$0.10 \times \frac{1}{3} + 0.20 \times \frac{2}{3} = \frac{1}{6} = 0.17.$$

This probability lies between 0.10 and 0.20 as expected and lies closer to 0.20 than 0.10. This is also clear since if the panda has had twins before, it is more likely to be of species  $B$  than species  $A$ . Hence it is more likely that she will birth twins again.

# Chapter 3

## Generalized Linear Models

This is a re-derivation of the result that the normal distribution with variance  $\sigma^2$  has the largest entropy amongst all distributions defined on  $[-\infty, +\infty]$  with variance  $\sigma^2$ . Let  $p(x)$  be the pdf of the normal distribution with mean  $\mu$  and variance  $\sigma^2$ . Let  $q(x)$  be a pdf with the same variance. Since the entropy of a distribution does not depend on its mean, we may assume that the mean of  $q(x)$  is  $\mu$ .

The entropy  $H(p)$  of the normal distribution is:

$$\begin{aligned} H(p) &= - \int_{-\infty}^{+\infty} p(x) \log p(x) dx \\ &= - \int_{-\infty}^{+\infty} p(x) \log((2\pi\sigma^2)^{-1/2}) dx + \int_{-\infty}^{+\infty} p(x) \frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 dx \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} \left( \frac{x-\mu}{\sigma} \right)^2 \exp \left\{ - \left( \frac{x-\mu}{\sigma} \right)^2 \right\} dx \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma} \int_{-\infty}^{+\infty} \frac{\sigma}{\sqrt{2\pi}} z^2 \exp \left\{ - \frac{z^2}{2} \right\} dz \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} E[Z^2] \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} (\text{Var}(Z) + (E[Z])^2) \\ &= \frac{1}{2} \log(2\pi e \sigma^2). \end{aligned}$$

Now that we have the entropy of the normal, we use the KL-divergence metric to measure the distance of  $q(x)$  from  $p(x)$ .

$$D_{\text{KL}}(q, p) = \int_{-\infty}^{+\infty} q(x) \log \frac{q(x)}{p(x)} dx = -H(q) + H(q, p).$$

At this point, we use the fact that  $D_{\text{KL}}(q, p) \geq 0$  for all distributions  $q$  and  $p$ . This gives us:  $H(q, p) \geq H(q)$ . We do not know what  $H(q)$  is but the expression for  $H(q, p)$  can be

evaluated quite easily.

$$\begin{aligned}
H(q, p) &= - \int_{-\infty}^{+\infty} q(x) \log p(x) dx \\
&= - \int_{-\infty}^{+\infty} q(x) \left[ \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right] dx \\
&= - \log \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{+\infty} q(x) dx + \frac{1}{2\sigma^2} \int_{-\infty}^{+\infty} q(x) (x - \mu)^2 dx \\
&= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \mathbb{E}_q[(x - \mu)^2] \\
&= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sigma^2 \\
&= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \\
&= \frac{1}{2} \log(2\pi e \sigma^2) \\
&= H(p).
\end{aligned}$$

So  $H(q, p) = H(p)$  and combining this with  $H(q, p) \geq H(q)$ , we obtain  $H(p) \geq H(q)$ .

### 3.1 Binomial Distributions and Maximum Entropy

Just as the normal distribution is the one with maximum entropy amongst all distributions that have a constant variance, the binomial is a distribution with the maximum entropy amongst all distributions defined on an experiment with just two outcomes and with a constant mean.

To make this precise, consider binary sequences  $S = \{s_i\}_{i=0}^{2^n-1}$  of length  $n$  sampled from a binomial distribution  $\text{Binom}(n, p)$ . Let  $x_i$  denote the number of 1s in  $s_i$ ; let  $p_i = \mathbf{Pr}\{s_i\} = p^{x_i}(1-p)^{n-x_i}$  and define  $\mathbf{p} = [p_0, \dots, p_{2^n-1}]$ . Let  $\mathbf{q} = [q_0, \dots, q_{2^n-1}]$  be some distribution over  $S$ . We wish to show that  $H(p) \geq H(q)$ . As in the case of the normal distribution, we start with the KL-divergence metric. We have

$$D_{\text{KL}}(\mathbf{q}, \mathbf{p}) = \sum_{i=0}^{2^n-1} q_i \log \frac{q_i}{p_i} = -H(\mathbf{q}) + H(\mathbf{q}, \mathbf{p}) \geq 0,$$

which yields that  $H(\mathbf{q}, \mathbf{p}) \geq H(\mathbf{q})$ .

Now consider the expression  $H(\mathbf{q}, \mathbf{p})$ . We may write:

$$\begin{aligned}
H(\mathbf{q}, \mathbf{p}) &= - \sum_{i=0}^{2^n-1} q_i \log p_i \\
&= - \sum_{i=0}^{2^n-1} q_i \log (p^{x_i} (1-p)^{n-x_i}) \\
&= - \sum_{i=0}^{2^n-1} q_i \left[ x_i \log \frac{p}{1-p} + n \log(1-p) \right] \\
&= - \left( \log \frac{p}{1-p} \right) \sum_{i=0}^{2^n-1} q_i x_i - n \log(1-p) \sum_{i=0}^{2^n-1} q_i \\
&= - \left( \log \frac{p}{1-p} \right) \bar{q} - n \log(1-p).
\end{aligned}$$

Here  $\bar{q}$  is the expected value of the distribution  $\mathbf{q}$ . If we assume that  $\bar{q} = \sum_{i=0}^{2^n-1} p_i x_i$ , then we can retrace the steps in the last derivation and show that  $H(\mathbf{q}, \mathbf{p}) = H(\mathbf{p})$ . Combining this with the inequality  $H(\mathbf{q}, \mathbf{p}) \geq H(\mathbf{q})$ , we obtain that  $H(\mathbf{p}) \geq H(\mathbf{q})$ . That is, any distribution  $\mathbf{q}$  on binary sequences of length  $n$  with the same expected value  $\sum_{i=0}^{2^n-1} p_i x_i$  as  $\mathbf{p}$ , has entropy at most that of  $\mathbf{p}$ .

As a matter of fact,

$$\begin{aligned}
\sum_{i=0}^{2^n-1} p_i x_i &= \sum_{k=0}^n \binom{n}{k} k p^k (1-p)^{n-k} \\
&= \sum_{k=0}^n \frac{n!}{k!(n-k)!} k p^k (1-p)^{n-k} \\
&= \sum_{k=1}^n \frac{n!}{(k-1)!(n-k)!} p^k (1-p)^{n-k} \\
&= np \sum_{k=1}^n \frac{(n-1)!}{(k-1)!(n-k)!} p^{k-1} (1-p)^{n-k} \\
&= np \sum_{j=0}^{n-1} \frac{(n-1)!}{j!(n-1-j)!} p^j (1-p)^{n-1-j} \\
&= np (p + (1-p))^n \\
&= np,
\end{aligned}$$

which matches with intuitive result that if the probability of a 1 is  $p$ , then the expected number of 1s in  $n$  trials is  $np$ .

## 3.2 The Poisson and Related Families of Distributions

The Poisson distribution is closely connected to distributions such as the binomial, exponential and the discrete Gamma (known as the Erlang distribution). To understand this connection, we begin with a description of Poisson process. A Poisson process is an arrival process that is described by a function  $P(k, \tau)$  that gives the probability of  $k$  arrivals in a time interval of length  $\tau$  and which satisfies the following conditions:

1.  $P(k, \tau)$  is the same for all intervals of length  $\tau$ ;
2. the number of arrivals during a particular time interval is independent of the history of arrivals outside this time interval;
3. there exists a constant  $\lambda$  such that:
  - $P(0, \tau) = 1 - \lambda\tau + o(\tau)$
  - $P(1, \tau) = \lambda\tau + o_1(\tau)$
  - $P(k, \tau) = o_k(\tau)$  for  $k \in \{2, 3, \dots\}$

The last condition states that in the regime of small time intervals, the probability of a single arrival in a time interval of length  $\tau$  is proportional to the length of the time interval. The functions  $o(\tau), o_k(\tau)$  satisfy the conditions:

$$\lim_{\tau \rightarrow 0} \frac{o(\tau)}{\tau} = 0, \quad \lim_{\tau \rightarrow 0} \frac{o_k(\tau)}{\tau} \text{ for } k \in \{1, 2, \dots\}.$$

The functions  $o(\tau)$  and  $o_k(\tau)$  can be thought of as the higher-order terms in a Taylor series expansion of function of  $\tau$ .

### 3.2.1 Number of arrivals

In order to compute a closed-form expression for  $P(k, \tau)$ , imagine dividing the time interval  $\tau$  in  $n$  equal disjoint pieces intervals, each of length  $\delta$ . If  $n$  is large enough, then the probability of a single arrival in a time interval of length  $\delta$  is approximately  $\lambda\delta$  and that of two or more arrivals is 0. Since the arrivals in each time window of length  $\delta$  is independent of the arrivals in the other intervals, we may approximate this using a Bernoulli distribution with success probability  $p = \lambda\delta$  and  $n = \tau/\delta$  trials. The point to note is that  $np = \lambda\tau$  is a constant. Thus the probability of  $k$  successes is:

$$\begin{aligned} P(k, \tau) &= \binom{n}{k} p^k (1-p)^{n-k} \\ &= \frac{n!}{(n-k)!k!} \left(\frac{\lambda\tau}{n}\right)^k \left(1 - \frac{\lambda\tau}{n}\right)^{n-k} \\ &= \frac{n(n-1)\cdots(n-k+1)}{n^k} \cdot \frac{1}{k!} \cdot (\lambda\tau)^k \cdot \left(1 - \frac{\lambda\tau}{n}\right)^{-k} \cdot \left(1 - \frac{\lambda\tau}{n}\right)^n \\ &= \frac{n}{n} \left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right) \cdot \frac{1}{k!} \cdot (\lambda\tau)^k \cdot \left(1 - \frac{\lambda\tau}{n}\right)^{-k} \cdot \left(1 - \frac{\lambda\tau}{n}\right)^n. \end{aligned}$$



Fix  $k$  and let  $n \rightarrow \infty$ . In the limit we obtain:

$$P(k, \tau) = \frac{(\lambda\tau)^k}{k!} e^{-\lambda\tau}.$$

Note that  $P(0, \tau) = e^{-\lambda\tau}$  and  $P(1, \tau) = \lambda\tau e^{-\lambda\tau}$ . Recall that the Taylor series for an infinitely differentiable function at  $x = 0$  can be written as:

$$f(0 + \delta) = f(0) + f^{(1)}(0)\delta + \frac{f^{(2)}(0)}{2!}\delta^2 + \dots + \frac{f^{(n)}(0)}{n!}\delta^n + \dots$$

The Taylor series expansion for  $e^{-\lambda\tau}$  is:

$$e^{-(0+\lambda\tau)} = 1 - \lambda\tau + \frac{(\lambda\tau)^2}{2!} - \frac{(\lambda\tau)^3}{3!} + \dots + (-1)^n \frac{(\lambda\tau)^n}{n!} + \dots$$

Using this, we obtain that  $P(0, \tau) = 1 - \lambda\tau + o(\tau)$  and  $P(1, \tau) = \lambda\tau + o_1(\tau)$ , consistent with the small-interval probability specifications.