

# BASIC STATISTICAL MODELS

ETHAN LEVIEN

## CONTENTS

1. Binomial Distribution	1
2. Uniform distribution and probability density (optional)	2
2.1. Joint density and conditional density	2
2.2. Cumulative density function	2
3. Normal distribution and the central limit theorem	3
4. Transformations of random variables	4
4.1. Standardizing	4
5. Linear regression model	4
5.1. Working with regression	5
5.2. Interpretation of regression parameters	5

## 1. BINOMIAL DISTRIBUTION

A situation that often arises is that we take many, say  $N$ , independent samples from a Bernoulli distribution. Now let  $Y$  be the number of 1s. Symbolically,

(1) 
$$Y = \sum_{i=1}^N y_i, \quad y_i \sim \text{Bernoulli}(q).$$

Then  $Y$  follows **binomial distribution**:

(2) 
$$Y \sim \text{Binomial}(N, q)$$

The binomial distribution has two parameters,  $N$  and  $p$ . Now let's think about the probability distribution. The chance to find any particular configuration of  $k$  ones is  $q^k(1 - q)^{N-k}$  because they are independent. For example

(3) 
$$P(y_1 = 1, y_2 = 0, y_3 = 1) = P(y_1 = 1)P(y_2 = 0)P(y_3 = 1)$$

(4) 
$$= q(1 - q)q = q^2(1 - q).$$

However, there are many configurations with  $k$  ones, in-fact there are

(5) 
$$\binom{N}{k} = \frac{N!}{k!(N - k)!},$$

and therefore

(6) 
$$P(Y = k) = \binom{N}{k} q^k (1 - q)^{N-k}.$$

The binomial distribution has a mean and variance

(7) 
$$\mathbb{E}[Y] = qN \quad \text{var}(Y) = Nq(1 - q).$$

These formulas come from the fact that for sums of independent variables, the variance and expectation sum.

An important feature of the Bernoulli random variables is that the mean grows much faster in  $N$  than the standard deviation. This means that when  $N$  is very large, the deviations from the average will become very small relative to the mean. An important measure of variation relative to the mean is the coefficient of variation

(8) 
$$CV = \frac{\sqrt{\text{var}(Y)}}{\mathbb{E}[Y]}.$$

Binomial samples can be generated in numpy with

```
> y = np.random.binomial(n,p,n_samples)
```

Often we are interested not in  $Y$ , but the fraction  $\phi = Y/N$ . For example, we might be interested in the vote share in an election. You should be able to see that  $\mathbb{E}[\phi] = q$ . What about the variance?

$$(9) \quad \text{var}(\phi) = \text{var}(Y/N) = \frac{1}{N^2} \text{var}(Y) = \frac{q(1-q)}{N}$$

Notice that this will tend towards zero as  $N \rightarrow \infty$ . Meanwhile,  $\mathbb{E}[\phi]$  has no dependence on  $N$ . This is a very important property, as it allows us to determine  $q$  by approximating  $\mathbb{E}[\phi]$  with the sample mean.

**Example 1.** *Coefficient of variation*

**Exercise 1:** *Generating binomial samples*

**Exercise 2:** *Binomial election modeling*

You should recognize that the assumption of independence is very important here. The following example illustrates an instance where this may break down for an election model. It is a bit contrived, but contrived examples, which we sometimes refer to as **toy models**, can be very helpful when it comes to building our intuition.

**Exercise 3:** *More election modeling*

## 2. UNIFORM DISTRIBUTION AND PROBABILITY DENSITY (OPTIONAL)

A uniform random variable, denoted

$$(10) \quad Y \sim \text{Uniform}(a, b)$$

has an equal chance of taking any number in the interval  $[a, b]$  (we assume  $a < b$ ). Let  $L = b - a$ . This is distinct from other distributions we have encountered in that it is a **continuous distribution**, rather than discrete. For the uniform distribution,

$$(11) \quad P(y_1 \leq Y \leq y_2) = \frac{y_2 - y_1}{L}$$

for  $a < y_1 < y_2 < b$ . That is, the chance for  $Y$  to fall in any interval is simply the length of that interval. This insures that the probability of  $Y$  being somewhere in  $[a, b]$  is one:  $P(a \leq Y \leq b) = 1$ . Note that as  $y_2 \rightarrow y_1$ ,  $P(y_1 \leq Y \leq y_2) \rightarrow 0$ . This tells us that the chance for  $Y$  to take any specific value is 0. Indeed, there are simply too many numbers (uncountably many) in any interval to assign positive probability to each. For continuous variables, it is sometimes useful to work with the density,  $f(y)$  (we will use lower case letters for density and uppercase for probability distributions).  $f(y)$  is the probability per unit  $Y$ , meaning that if we look in a small interval

$$(12) \quad f(y)dy = P(y \leq Y \leq y + dy) = \frac{dy}{L}.$$

Thus, for uniform distribution the density is  $1/L$  if  $y \in [a, b]$  and 0 otherwise.

**2.1. Joint density and conditional density.** Conditioning works for probability density just as it does for probability distributions.

**Example 2.** *Conditioning with continuous variables*

**2.2. Cumulative density function.** Sometimes it is useful to characterize a continuous distribution not by the density, but by the **cumulative distribution function (CDF)**, defined as

$$(13) \quad F(y) = P(Y < y).$$

What is the CDF of the uniform distribution? The **median** is the value  $y_m$  for which  $F(y_m) = 1/2$ . What is the median of a Uniform distribution?

To better understand density and CDF, imagine a student says they will arrive at my office between noon and 3. Let  $Y$  represent the time a student arrives, which we will model as a Uniform random variable. Then the density is  $f(y) = 1/3$  which has units 1/hours. We can think of  $f$  as the rate at which the CDF increases – that is, it is the velocity of probability.

# BASIC STATISTICAL MODELS

3

## 3. NORMAL DISTRIBUTION AND THE CENTRAL LIMIT THEOREM

In the previous example, we say that if we take the average of many Bernoulli random variables, we get a histogram that looks a lot like a “bell curve” with a standard deviation was proportional to  $1/\sqrt{n}$ . It turns out this is true when we add up *any* sequence of independent and independent distributed random variables which are not too pathological (actually it is also true for many sequences of random variables which are not independent).

It is useful to define a special random variable which captures the statistical behavior of random sums. We call this a Normal random distribution

(14) 
$$Y \sim \text{Normal}(\mu, \sigma).$$

We can generate Normal random variables in python with

```
> np.random.normal(0,1)
```

The Normal distribution is defined by the Gaussian

(15) 
$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}.$$

This is the classic bell curve shown in Figure 1. The probability distribution for the Normal distribution is defined by the area under this curve. As I discuss in the previous section can think of  $f$  as the probability *per unit of the random variable*, e.g. probability/feet.

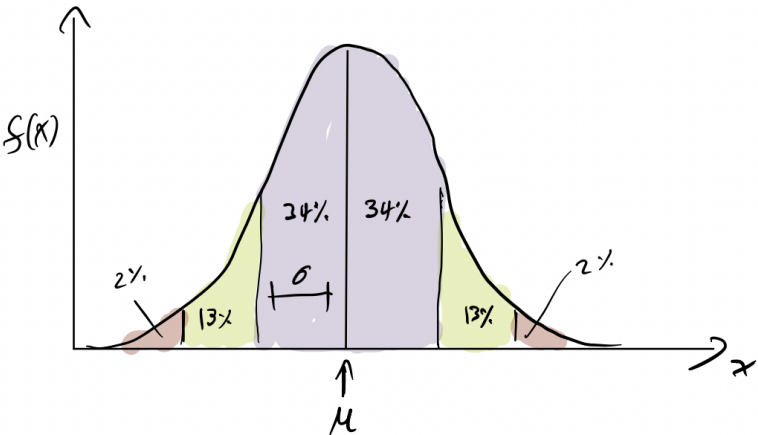


FIGURE 1. Probabilities in the Normal distribution

We use the curve above to calculate probabilities of events in the Normal distribution. For example

(16) 
$$Y \sim \text{Normal}(5, 2)$$

what is (approximately)  $P(Y > 7)$ ? Note that  $5 + 2 = 7$ , so this is asking how likely it is that a Normal variable is greater than 1 standard deviation above the mean. This about  $13.5 + 2 = 15.5\%$

The **central limit theorem** tell us that when  $y_i$  are independent and have a finite mean and variance  $\mu_y$  and  $\sigma_y$  and

(17) 
$$\hat{Y} = \frac{1}{N} \sum_{i=1}^N y_i,$$

then the distribution of  $\hat{Y}$  should be close to that of

(18) 
$$Y \sim \text{Normal}(\mu_y, \sigma_y/\sqrt{N}).$$

**Example 3.** Comparing histograms

**Example 4.** Working with Normal random variables

**Exercise 4:** Hemoglobin levels

## 4. TRANSFORMATIONS OF RANDOM VARIABLES

Now consider

$$(19) \quad X \sim \text{Normal}(\mu_x, \sigma_x)$$

and let

$$(20) \quad Y = aX + b$$

We are just multiplying and shifting everything. Think about what this does to the histogram (and try it in Python). Hopefully you can convince yourself that  $Y$  should also be Normal, but what are the mean and variance? Taking the average of both sides,

$$(21) \quad \mathbb{E}[Y] = a\mu + b$$

and

$$(22) \quad \text{var}(Y) = \text{var}(aX) + \text{var}(b)$$

From the formula for variance, we know  $\text{var}(aX) = a^2\text{var}(X)$ . Also,  $\text{var}(b) = 0$  So

$$(23) \quad Y \sim \text{Normal}(a\mu_x + b, |a|\sigma_x).$$

**4.1. Standardizing.** We can transform any random variable into a so-called standard normal

$$(24) \quad Z \sim \text{Normal}(0, 1).$$

For defined above,

$$(25) \quad Z = \frac{X - \mu_x}{\sigma_x}$$

Then  $a = 1/\sigma_x$  and  $b = -\mu_x/\sigma_x$ . Plugging into Equation (23) yields a standard Normal. **Transforming  $X$  to a standard Normal is equivalent to measuring  $X$  in units of standard deviations.** For example, if we make a histogram of  $X$ , all this transformation does is change the  $X$  axis to units of standard deviations from the mean.

## 5. LINEAR REGRESSION MODEL

We now introduce the concept of regression modeling. A very broad class of models in statistics for the relationship between two variables  $X$  and  $Y$  is a regression model:

$$(26) \quad Y = f(X) + \epsilon$$

where  $f$  is a deterministic function; that is, if we evaluate  $f$  at a particular number, we get something that is not random. The term  $\epsilon$  represents some source of noise **independent of  $X$** , and is typically modeled with a Normal distribution

$$(27) \quad \epsilon \sim \text{Normal}(0, \sigma_\epsilon).$$

In other words, it represents things other than  $X$  which may influence  $Y$ . Regardless of how  $X$  is distributed, for any given values  $X = x$ ,  $Y$  must have a Normal distribution:

$$(28) \quad Y|(X = x) \sim \text{Normal}(f(x), \sigma_\epsilon).$$

Of particular interest (due to its simplicity) is the case

$$(29) \quad f(x) = ax + b$$

which is the subject of this class. That is, we are interested in the model

$$(30) \quad Y = aX + b + \epsilon$$

where

$$(31) \quad \epsilon \sim \text{Normal}(0, \sigma_\epsilon).$$

# BASIC STATISTICAL MODELS

5

**5.1. Working with regression.** In Equation (30),  $X$  could be anything, but let's suppose  $X$  is drawn from a Normal distribution. This gives us the model

$$(32) \quad Y = aX + b + \epsilon$$

$$(33) \quad X \sim \text{Normal}(\mu_x, \sigma_x)$$

(technically this is not a regression model anymore because we specify the distribution of  $X$ ). Now that we have two random variables, we can ask questions like, what is the joint distribution? what are the conditional distributions? What are the Marginal distributions? Using properties of Normal random variables, we get

$$(34) \quad Y \sim \text{Normal}(a\mu + b, |a|\sigma_x + \sigma_\epsilon)$$

This is the distribution of  $Y$  regardless of  $X$ ; that is, it is the distribution we would get if we randomly sampled  $Y$  values ignoring what the value of  $X$  is. What is another name for this? This distribution can be understood visually [DRAW DIAGRAM ON BOARD].

$X$  and  $Y$  are not independent. This can be seen visually. **Note: even though  $X$  and  $Y$  are not independent, we don't say they effect each other. To see why this terminology is problematic note that  $Y$  has no "effect" on  $X$  (it is determined by  $X$ ). However, what is  $X|Y = y$ ?**

**Example 5.** *conditioning with continuous variables*

**5.2. Interpretation of regression parameters.** The interpretation of parameters in regression model is important. Sometimes (the example below) there is a clear meaning in terms of conditional distributions, but suppose we have a model of our time in a race as a function of the temperature:

$$(35) \quad Y = aX + b + \epsilon$$

Let's start with  $a$ : This is the average difference between times for races at temperature which differ by 1. (we expect this should be negative).  $\epsilon$  is the variation around this average. Now what about  $b$ ? The problem is that if we plug in  $X = 0$  we obtain a nonsensical quantity, so really  $b$  does not have physical interpretation. This is common situation with regression models, and for this reason it is sometimes advantages to center  $X$ .

**Exercise 5:** *Kid's test scores*

**Exercise 6:** *The random walk*