

# Exploration librairie swirl : Common Distributions

---

> swirl()

Welcome to swirl! Please sign in. If you've been here before, use the same name as you did then. If you are new, call yourself something unique.

What shall I call you? **jlbellier**

Please choose a course, or type 0 to exit swirl.

- 1: Statistical Inference
- 2: Take me to the swirl course repository!

Selection: **1**

Please choose a lesson, or type 0 to return to course menu.

- |                           |                 |                           |
|---------------------------|-----------------|---------------------------|
| 1: Introduction           | 2: Probability1 | 3: Probability2           |
| 4: ConditionalProbability | 5: Expectations | 6: Variance               |
| 7: CommonDistros          | 8: Asymptotics  | 9: T Confidence Intervals |
| 10: Hypothesis Testing    | 11: P Values    | 12: Power                 |
| 13: Multiple Testing      | 14: Resampling  |                           |

Selection: **7**

Attempting to load lesson dependencies...

Package 'ggplot2' loaded correctly!

0%

Common Distributions. (Slides for this and other Data Science courses may be found at github <https://github.com/DataScienceSpecialization/courses/>. If you care to use them, they must be downloaded as a zip file and viewed locally. This lesson corresponds to 06\_Statistical\_Inference/06\_CommonDistros.)

==

2%

Given the title of this lesson, what do you think it will cover?

- 1: Common Bistros
- 2: Common Distributions
- 3: I haven't a clue
- 4: Rare Distributions

Selection: **2**

**That's a job well done!**

====

5%

The first distribution we'll examine is the Bernoulli which is associated with experiments which have only 2 possible outcomes. These are also called (by people in the know) binary trials.

=====

7%

It might surprise you to learn that you've probably had experience with Bernoulli trials. Which of the following would be a Bernoulli trial?

- 1: Tossing a die
- 2: Spinning a roulette wheel
- 3: Drawing a card from a deck
- 4: Flipping a coin

Selection: 4

**Nice work!**

=====

9%

For simplicity, we usually say that Bernoulli random variables take only the values 1 and 0. Suppose we also specify that the probability that the Bernoulli outcome of 1 is  $p$ . Which of the following represents the probability of a 0 outcome?

- 1:  $p$
- 2:  $p(1-p)$
- 3:  $1-p$
- 4:  $p^2$

Selection: 3

**You nailed it! Good job!**

=====

11%

If the probability of a 1 is  $p$  and the probability of a 0 is  $1-p$  which of the following represents the PMF of a Bernoulli distribution? Recall that the PMF is the function representing the probability that  $X=x$ .

- 1:  $p^{(1-x)} * (1-p)^{(1-x)}$
- 2:  $p * (1-p)$
- 3:  $x * (1-x)$
- 4:  $p^x * (1-p)^{(1-x)}$

Selection: 4

**You're the best!**

=====

14%

Recall the definition of the expectation of a random variable. Suppose we have a Bernoulli random variable and, as before, the probability it equals 1 (a success) is  $p$  and probability it equals 0 (a failure) is  $1-p$ . What is its mean?

1:  $p^2$

2:  $1-p$

3:  $p(1-p)$

4:  $p$

Selection: 4

**That's correct!**

=====

**16%**

Given the same Bernoulli random variable above, which of the following represents  $E(X^2)$

1:  $p$

2:  $p(1-p)$

3:  $p^2$

4:  $1-p$

5:  $(1-p)^2$

Selection: 1

**You nailed it! Good job!**

=====

**18%**

Use the answers of the last two questions to find the variance of the Bernoulli random variable. Recall  $\text{Var} = E(X^2) - (E(X))^2$

1:  $p(1-p)$

2:  $p(p-1)$

3:  $p^2(1-p)^2$

4:  $p^2-p$

Selection: 1

**You are really on a roll!**

=====

**20%**

Binomial random variables are obtained as the sum of iid Bernoulli trials. Specifically, let  $X_1, \dots, X_n$  be iid Bernoulli( $p$ ) random variables; then  $X = X_1 + X_2 + \dots + X_n$  is a binomial random variable. Binomial random variables represent the number of successes,  $k$ , out of  $n$  independent Bernoulli trials. Each of the trials has probability  $p$ .

=====

**23%**

The PMF of a binomial random variable  $X$  is the function representing the probability that  $X=x$ . In other words, that there are  $x$  successes out of  $n$  independent trials. Which of the following represents the PMF of a binomial distribution? Here  $x$ , the number of successes, goes from 0 to  $n$ , the number of trials, and  $\text{choose}(n,x)$  represents the binomial coefficient 'n choose x' which is the number of ways  $x$  successes out of  $n$  trials can occur regardless of order.

1:  $\text{choose}(n,x) * p^{(n-x)} * (1-p)^x$

2:  $\text{choose}(n,x) * p^x * (1-p)^{(n-x)}$

```
3: choose(n,x) * p^x*(1-p)^(1-x)
4: p^x
```

Selection: 2

**You nailed it! Good job!**

=====

25%

Suppose we were going to flip a biased coin 5 times. The probability of tossing a head is .8 and a tail .2. What is the probability that you'll toss at least 3 heads.

```
> choose(5,3)*0.8^3*0.2^2 + choose(5,4)*0.8^4*0.2+ choose(5,5)*0.8^5
[1] 0.94208
```

**You are really on a roll!**

=====

27%

Now you can verify your answer with the R function pbinom. The quantile is 2, the size is 5, the prob is .8 and the lower.tail is FALSE. Try it now.

```
> pbinom(q=2,size=5,prob=0.8,lower.tail = FALSE)
[1] 0.94208
```

**Perseverance, that's the answer.**

=====

30%

Another very common distribution is the normal or Gaussian. It has a complicated density function involving its mean  $\mu$  and variance  $\sigma^2$ . The key fact of the density formula is that when plotted, it forms a bell shaped curve, symmetric about its mean  $\mu$ . The variance  $\sigma^2$  corresponds to the width of the bell, the higher the variance, the fatter the bell. We denote a normally distributed random variable  $X$  as  $X \sim N(\mu, \sigma^2)$ .

=====

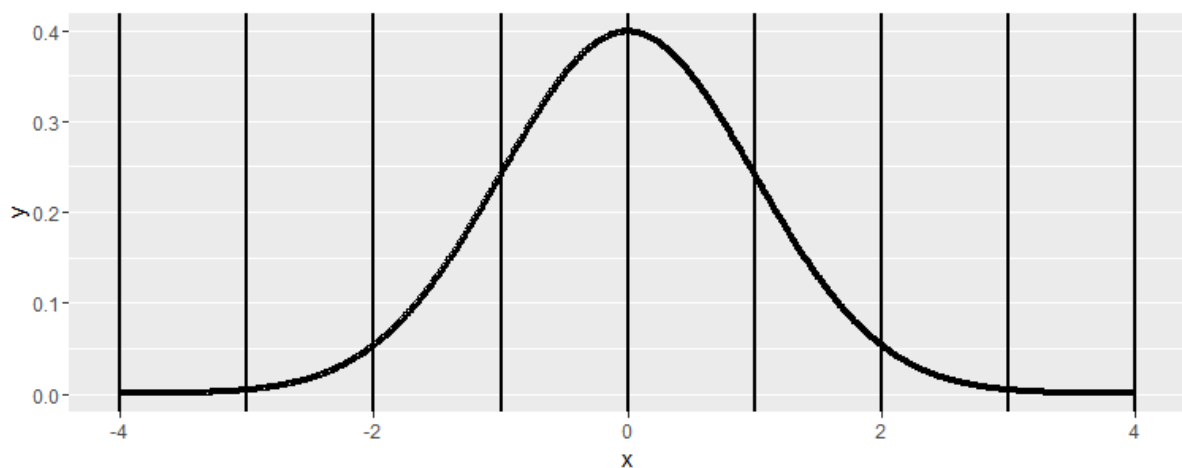
32%

When  $\mu = 0$  and  $\sigma = 1$  the resulting distribution is called the standard normal distribution and it is often labeled  $Z$ .

=====

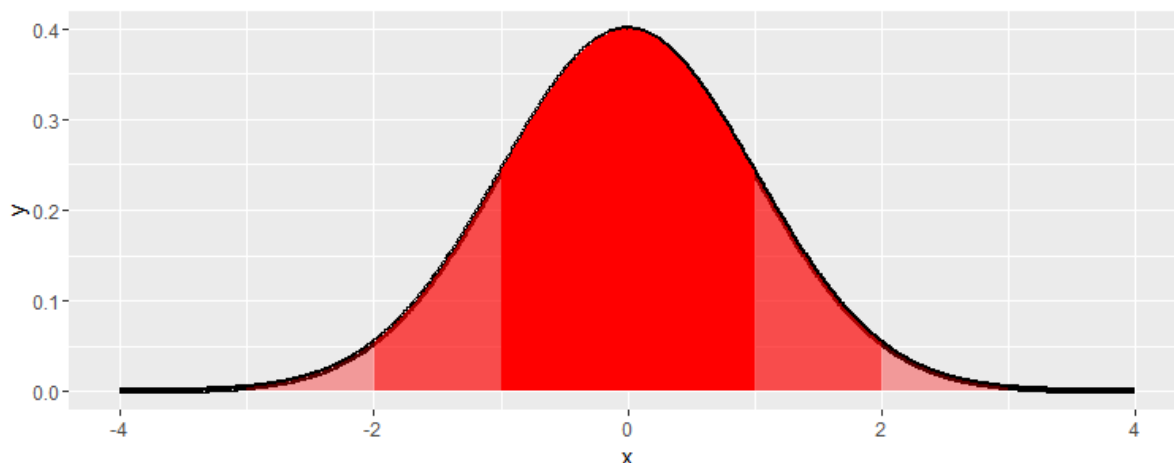
34%

Here's a picture of the density function of a standard normal distribution. It's centered at its mean 0 and the vertical lines (at the integer points of the x-axis) indicate the standard deviations.



36%

Approximately 68%, 95% and 99% of the normal density lie within 1, 2 and 3 standard deviations from the mean, respectively. These are shown in the three shaded areas of the figure. For example, the darkest portion (between -1 and 1) represents 68% of the area.



39%

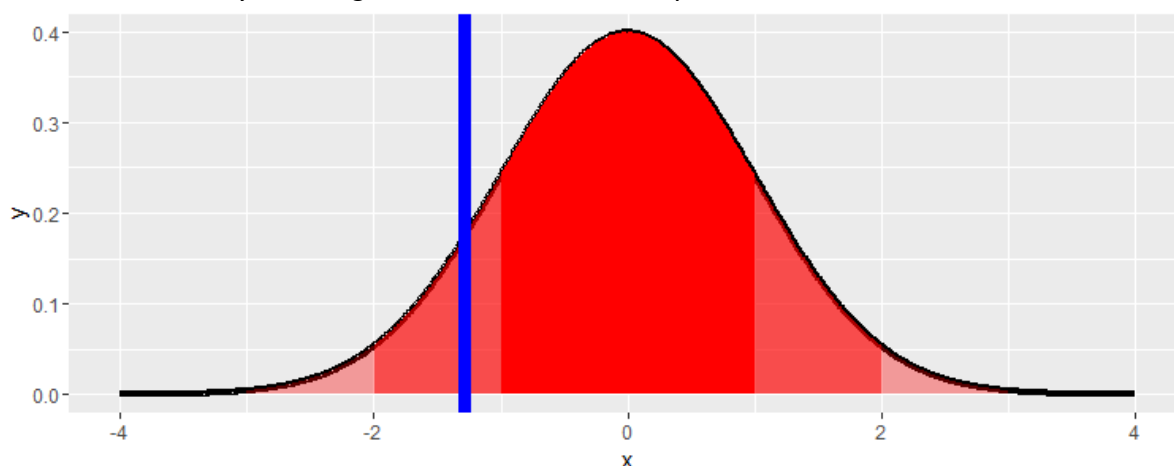
The R function `qnorm(prob)` returns the value of  $x$  (quantile) for which the area under the standard normal distribution to the left of  $x$  equals the parameter `prob`. (Recall that the entire area under the curve is 1.) Use `qnorm` now to find the 10th percentile of the standard normal. Remember the argument `prob` must be between 0 and 1. You don't have to specify any of the other parameters since the default is the standard normal.

```
> qnorm(0.1)
[1] -1.281552
```

Nice work!

41%

We'll see this now by drawing the vertical line at the quantile -1.281552.



43%

Which of the following would you expect to be the 1st percentile?

- 1: 2.33
- 2: -1.28

3: -1.0  
4: 0  
5: -2.33  
Selection: 5

**Keep up the great work!**

===== 45%  
By looking at the picture can you say what the 50th percentile is?

> 0  
[1] 0

**That's correct!**

===== 48%  
We can use the symmetry of the bell curve to determine other quantiles. Given that 2.5% of the area under the curve falls to the left of  $x = -1.96$ , what is the 97.5 percentile for the standard normal?

1: -1.28  
2: 1.96  
3: 2.33  
4: 2  
Selection: 2

**Excellent work!**

===== 50%  
Here are two useful facts concerning normal distributions. If  $X$  is a normal random variable with mean  $\mu$  and variance  $\sigma^2$ , i.e.,  $X \sim N(\mu, \sigma^2)$ ,

===== 52%  
then the random variable  $Z$  defined as  $Z = (X - \mu)/\sigma$  is normally distributed with mean 0 and variance 1, i.e.,  $Z \sim N(0, 1)$ . ( $Z$  is standard normal.)

===== 55%  
The converse is also true. If  $Z$  is standard normal, i.e.,  $Z \sim N(0, 1)$ , then the random variable  $X$  defined as  $X = \mu + \sigma * Z$  is normally distributed with mean  $\mu$  and variance  $\sigma^2$ , i.e.,  $X \sim N(\mu, \sigma^2)$

===== 57%  
These formulae allow you to easily compute quantiles (and thus percentiles) for ANY normally distributed variable if you know its mean and variance. We'll show how to find the 97.5th percentile of a normal distribution with mean 3 and variance 4.

===== 59%  
Again, we can use R's `qnorm` function and simply specify the mean and standard deviation (the square root of the variance). Do this now. Find the 97.5th percentile of a normal distribution with mean 3 and standard deviation 2.

> `qnorm(0.975, mean=3, sd=2)`  
[1] 6.919928

**Keep working like that and you'll get there!**

===== 61%

Let's check it using the formula above,  $X = \mu + \sigma \cdot Z$ . Here we'll use the 97.5th percentile for the standard normal as the value  $Z$  in the formula. Recall that we previously calculated this to be 1.96. Let's multiply this by the standard deviation of the given normal distribution (2) and add in its mean (3) to see if we get a result close to the one `qnorm` gave us.

```
> 1.96*2+3  
[1] 6.92
```

**You're the best!**

===== 64%

Suppose you have a normal distribution with mean 1020 and standard deviation of 50 and you want to compute the probability that the associated random variable  $X > 1200$ . The easiest way to do this is to use R's `pnorm` function in which you specify the quantile (1200), the mean (1020) and standard deviation (50). You also must specify that the `lower.tail` is `FALSE` since we're asking for a probability that the random variable is greater than our quantile. Do this now.

```
> pnorm(q=1200,mean = 1020,sd=50,lower.tail=FALSE)  
[1] 0.0001591086
```

**Perseverance, that's the answer.**

===== 66%

Alternatively, we could use the formula above to transform the given distribution to a standard normal. We compute the number of standard deviations the specified number (1200) is from the mean with  $Z = (X - \mu)/\sigma$ . This is our new quantile. We can then use the standard normal distribution and the default values of `pnorm`. Remember to specify that `lower.tail` is `FALSE`. Do this now.

```
> pnorm(q=3.6,lower.tail=FALSE)  
[1] 0.0001591086
```

**Keep trying! Or, type `info()` for more options.**

Type `pnorm((1200-1020)/50,lower.tail=FALSE)` at the R prompt.

```
> pnorm((1200-1020)/50,lower.tail=FALSE)  
[1] 0.0001591086
```

**That's the answer I was looking for.**

===== 68%

For practice, using the same distribution, find the 75% percentile. Use `qnorm` and specify the probability (.75), the mean (1020) and standard deviation (50). Since we want to include the left part of the curve we can use the default `lower.tail=TRUE`.

```
> qnorm(0.75,mean=1020,sd=50,lower.tail = TRUE)
[1] 1053.724
```

**Great job!**

===== 70%

Note that R functions `pnorm` and `qnorm` are inverses. What would you expect `pnorm(qnorm(.53))` to return?

```
> 0.53
[1] 0.53
```

**You are quite good my friend!**

===== 73%

How about `qnorm(pnorm(.53))`?

```
> 0.53
[1] 0.53
```

**You are quite good my friend!**

===== 75%

Now let's talk about our last common distribution, the Poisson. This is, as Wikipedia tells us, "a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event."

===== 77%

In other words, the Poisson distribution models counts or number of event in some interval of time. From Wikipedia, "Any variable that is Poisson distributed only takes on integer values."

===== 80%

The PMF of the Poisson distribution has one parameter,  $\lambda$ . As with the other distributions the PMF calculates the probability that the Poisson distributed random variable  $X$  takes the value  $x$ . Specifically,  $P(X=x)=(\lambda^x e^{-\lambda})/x!$ . Here  $x$  ranges from 0 to infinity.

===== 82%

**The mean and variance of the Poisson distribution are both  $\lambda$ .**

=====84%

Poisson random variables are used to model rates such as the rate of hard drive failures. We write  $X \sim \text{Poisson}(\lambda * t)$  where  $\lambda$  is the expected count per unit of time and  $t$  is the total monitoring time.

=====86%

For example, suppose the number of people that show up at a bus stop is Poisson with a mean of 2.5 per hour, and we want to know the probability that at most 3 people show up in a 4 hour period. We use the R function `ppois` which returns a probability that the random variable is less than or equal to 3. We only need to specify the quantile (3) and the mean (2.5 \* 4). We can use the default parameters, `lower.tail=TRUE` and `log.p=FALSE`. Try it now.



```
> ppois(q=3,lambda =2.5*4,lower.tail=TRUE,log.p=FALSE)
[1] 0.01033605
```

**Your dedication is inspiring!**

=====89%

Finally, the Poisson distribution approximates the binomial distribution in certain cases. Recall that the binomial distribution is the discrete distribution of the number of successes,  $k$ , out of  $n$  independent binary trials, each with probability  $p$ . If  $n$  is large and  $p$  is small then the Poisson distribution with  $\lambda$  equal to  $n \cdot p$  is a good approximation to the binomial distribution.

=====91%

To see this, use the R function `pbinom` to estimate the probability that you'll see at most 5 successes out of 1000 trials each of which has probability .01. As before, you can use the default parameter values (`lower.tail=TRUE` and `log.p=FALSE`) and just specify the quantile, size, and probability.

```
> pbinom(q=5,size=1000,prob=0.01,lower.tail =TRUE, log.p=FALSE)
[1] 0.06613951
```

**Great job!**

=====93%

Now use the function `ppois` with quantile equal to 5 and  $\lambda$  equal to  $n \cdot p$  to see if you get a similar result.

```
> ppois(q=5,lambda=1000*0.01,lower.tail = TRUE,log.p = FALSE)
[1] 0.06708596
```

**You are quite good my friend!**

=====95%

See how they're close? Pretty cool, right? This worked because  $n$  was large (1000) and  $p$  was small (.01).

===== 98%

**Congrats! You've concluded this uncommon lesson on common distributions.**

=====100%