

ST1131

Michael Yang

1 dplyr

1.1 filter

Pass in any number of conditions:

```
1 starwars %>% filter(skin_color == "light", eye_color == "brown")
```

is roughly equivalent to

```
1 starwars[starwars$skin_color == "light" & starwars$eye_color == "brown", ]
```

1.2 arrange

Order rows by column, tiebreak by next column

```
1 starwars %>% arrange(height, mass) #sort by height, tiebreak by mass
2
3 starwars %>% arrange(desc(height)) #sort descending
```

1.3 slice

```
1 starwars %>% slice(5:10)
2
3 starwars %>% slice_head(n = 3) # or slice_tail for last few rows
4
5 starwars %>% slice_sample(n = 5) # random sample
6
7 starwars %>%
8   filter(!is.na(height)) %>%
9   slice_max(height, n = 3) # get highest 3 rows by height
```

1.4 select

```
1 # Select columns by name
2 starwars %>% select(hair_color, skin_color, eye_color)
3
4 # Select all columns between hair_color and eye_color (inclusive)
5 starwars %>% select(hair_color:eye_color)
6
7 # Select all columns except those from hair_color to eye_color (inclusive)
8 starwars %>% select(!(hair_color:eye_color))
9
10 # Select all columns ending with color
11 starwars %>% select(ends_with("color"))
12
13 # Select variables (drops all other variables not explicitly mentioned)
14 starwars %>% select(home_world = homeworld)
15
16 # Rename variables
17 starwars %>% rename(home_world = homeworld)
```

Helper functions are `starts_with`, `ends_with`, `matches()`, `contains()`

1.5 mutate

Add new columns

```
1 starwars %>% mutate(height_m = height / 100)
2
3 # Add and select
4 starwars %>%
```

```

5   mutate(height_m = height / 100) %>%
6   select(height_m, height, everything())
7
8 starwars %>%
9   mutate(
10     height_m = height / 100,
11     BMI = mass / (height_m^2)
12   ) %>%
13   select(BMI, everything())
14
15 # Only keep new variables
16 starwars %>%
17   mutate(
18     height_m = height / 100,
19     BMI = mass / (height_m^2),
20     .keep = "none"
21   )

```

1.6 relocate

Reorder columns

```

1 starwars %>% relocate(sex:homeworld, .before = height)

```

1.7 summarise

Collapse data frame to a single row

```

1 starwars %>% summarise(height = mean(height, na.rm = TRUE))
2
3 food_consumption %>%
4   # Filter for rice food category
5   filter(food_category == "rice") %>%
6   # Summarize the mean_co2 and median_co2
7   summarize(mean_co2 = mean(co2_emission),
8             median_co2 = median(co2_emission))

```

1.8 count

```

1 # Usage
2 # count(x, ..., wt = NULL, sort = FALSE, name = NULL)
3
4 starwars %>% count(species, sort = TRUE)
5 starwars %>% count(sex, gender, sort = TRUE)
6 starwars %>% count(birth_decade = round(birth_year, -1))

```

2 ggplot2

```
1 # Usage
2 ggplot(data = NULL, mapping = aes(), environment = parent.frame())
3
4 aes(x = ..., y = ..., colour = ..., size = ...)
```

`ggplot()` is used to construct the initial plot object, and is almost always followed by a plus sign (+) to add components to the plot.

```
1 + geom_histogram(bars = 5)
2 +
```

3 Measures of Spread

```
1 # Variance
2 var(food_consumption$consumption)
3
4 # SD
5 sd(food_consumption$consumption)
6
7 # Quartile
8 quartile(food_consumption$consumption)
9
10 # Quintile
11 quantile(food_consumption$co2_emission, probs = seq(0, 1, 0.2))
12 quantile(food_consumption$co2_emission, probs = c(0, 0.2, 0.4, 0.6, 0.8, 1.0))
13
14 # Decile
15 quantile(food_consumption$co2_emission, probs = seq(0, 1, 0.1))
16
17 # IQR
18 q1 <- quantile(food_consumption$co2_emission, 0.25) # 25 percentile
19 q3 <- quantile(food_consumption$co2_emission, 0.75) # 75 percentile
20 iqr <- q3 - q1
21
22 # Cutoff for outliers
23 lower <- q1 - 1.5 * iqr
24 upper <- q3 + 1.5 * iqr
25
26 # Filter to find outliers
27 food_consumption %>%
28   filter(co2_emission < lower | co2_emission > upper) # note the | or condition instead of ,
```

4 Probability

```
1 # Fixed seed
2 set.seed(5)
3
4 # Random sample
5 sales_count %>%
6   sample_n(1)
7
8 # Sampling without replacement
9 sales_count %>%
10  sample_n(2)
11
12 # Sample with replacement
13 sales_count %>%
14   # Sampling with replacement is independent, without replacement is dependent
15   sample_n(2, replace = TRUE)
16
17 # Count number of deals and probability of each deal
18 amir_deals %>%
19   count(product) %>%
20   mutate(prob = n / sum(n))
21
22 # mean
23 mean(x, trim = 0, na.rm = FALSE) #na.rm means whether or not to remove NA
```

5 Distributions

Law of large sizes: sample mean approaches theoretical mean as sample size increases

5.1 Discrete

```
1 rolls_10 <- die %>%
2   sample_n(10, replace = TRUE) %>%
3   ggplots(aes(n)) +
4     geom_histogram(bins = 6)
5
6 expected_val <- sum(size_distribution$group_size *
7                     size_distribution$probability)
8
9
10 # Calculate probability of picking group of 4 or more
11 size_distribution %>%
12   # Filter for groups of 4 or larger
13   filter(group_size >= 4) %>%
14   # Calculate prob_4_or_more by taking sum of probabilities
15   summarize(prob_4_or_more = sum(probability))
```

5.2 Continuous

```
1 # P(<= 7)
2 punif(7, min = 0, max = 12)
3
4 # P(> 7)
5 punif(7, min = 0, max = 12, lower.tail = FALSE)
6
7 # P(4 <= p <= 7)
8 punif(7, min = 0, max = 12) - punif(4, min = 0, max = 12)
9
10 # Set random seed to 334
11 set.seed(334)
12
13 # Generate 1000 wait times between 0 and 30 mins, save in time column
14 wait_times %>%
15   mutate(time = runif(1000, min = 0, max = 30)) %>%
16   # Create a histogram of simulated times
17   ggplot(aes(time)) +
18     geom_histogram(bins = 30)
19
20 # Simulate n generations of values between 0 to 30
21 runif(n, min = 0, max = 30)
```

5.3 Binomial

```
1 # Simulating trials
2 rbinom(trials, coins, probability of success)
3 rbinom(8, 1, 0.5) # 8 flips of 1 coin with 50% success
4 rbinom(1, 8, 0.5) # 1 flip of 8 coins with 50% success
5 rbinom(10, 3, 0.5) # 10 flips of 3 coins
6
7 # Calculating probability
8 # dbinom(num heads, num trials, prob of heads)
9 dbinom(7, 10, 0.5) # P(heads = 7)
10 pbinom(7, 10, 0.5) # P(heads <= 7)
11 pbinom(7, 10, 0.5, lower.tail = FALSE) # P(heads > 7)
12 1 - pbinom(7, 10, 0.5) # P(heads > 7)
```

```

13
14 # Simulate 52 weeks of 3 deals
15 deals <- rbinom(52, 3, 0.3)
16 # Calculate mean deals won per week
17 mean(deals)
18
19 # Probability of closing 3 out of 3 deals with 30% success
20 dbinom(3, 3, 0.3)
21
22 # Probability of closing <= 1 deal out of 3 deals
23 pbinom(1, 3, 0.3)
24
25 # Probability of closing > 1 deal out of 3 deals
26 pbinom(1, 3, 0.3, lower.tail = FALSE)
27
28 # Expected number won with 30% win rate
29 won_30pct <- 3 * 0.3

```

Expected value = $n \times p$

Each trial must be independent for binomial distribution

5.4 Normal

```

1 # Percentage of values <= 154
2 pnorm(154, mean = 161, sd = 7)
3
4 # values > 154
5 pnorm(154, mean = 161, sd = 7, lower.tail = FALSE)
6
7 # 154 <= value <= 157
8 pnorm(157, mean = 161, sd = 7) - pnorm(154, mean = 161, sd = 7)
9
10 # value that 90% of values <=
11 qnorm(0.9, mean = 161, sd = 7)
12
13 # value that 90% of values >=
14 qnorm(0.9, mean = 161, sd = 7, lower.tail = FALSE)
15
16 # Generate 10 random numbers
17 rnorm(10, mean = 161, sd = 7)
18
19 # Compare performance over 2 quarters
20 pnorm(1000, mean = 5000, sd = 2000, lower.tail = FALSE)
21 pnorm(1000, mean = 5600, sd = 2600, lower.tail = FALSE)

```

5.5 Central Limit Theorem

```

1 die <- c(1, 2, 3, 4, 5, 6)
2 sample_5 <- sample(die, 5, replace = TRUE)
3 sample(die, 5, replace = TRUE) %>% mean()
4
5 # Replicating 10 times of random sample of 5
6 sample_mean <- replicate(10, sample(die, 5, replace = TRUE) %>% mean())
7 sample_sd <- replicate(10, sample(die, 5, replace = TRUE) %>% sd())

```

Central Limit Theorem: sample distribution becomes closer to normal distribution as number of trials increase.
(*samples should be random and independent)

5.6 Poisson

Events appear to happen at a certain rate, but completely at random.

E.g. number of animals adopted from animal shelters per week / Number of people arriving at a restaurant per hour / number of earthquakes per year
Poisson distribution = probability of some # of events occurring over a fixed period of time, e.g. probability of < 20 earthquakes in California per year

λ = average number of events per time interval e.g. average number of adoptions per week. (Also the expected value of the distribution)

Lambda is the distribution's peak.

```
1 # P(= 5), average = 8 (exact value)
2 dpois(5, lambda = 8)
3
4 # P(<= 5), average = 8 (less than or equal to)
5 ppois(5, lambda = 8)
6
7 # P(> 5), avg = 8
8 ppois(5, lambda = 8, lower.tail = FALSE)
9
10 # Sampling from Poisson distribution
11 rpois(10, lambda = 8)
```

CLT still applies to sample means of poisson distribution

5.7 Exponential

Represents probability of certain time passing between Poisson events e.g. probability of ≥ 1 day between adoption, probability of ≥ 10 mins between restaurants, probability of 6-8 months between earthquakes. Also uses lambda (Rate)

Continuous (time)

```
1 # P(wait < 1min)
2 pexp(1, rate = 0.5)
3
4 # P(wait > 4min)
5 pexp(4, rate = 0.5, lower.tail = FALSE)
6
7 # P(1min < wait < 4min)
8 pexp(4, rate = 0.5) - pexp(1, rate = 0.5)
9
10 # Average response time = 2.5 hours from receiving request
11 # P(< 1 hour to respond)
12 pexp(1, rate = 1 / 2.5)
13
14 # P(> 4hrs to respond)
15 pexp(1, rate = 0.4, lower.tail = FALSE)
16
17 # P(3 < wait < 4 hrs to respond)
18 pexp(4, rate = 0.4) - pexp(3, rate = 0.4)
```

Expected time = $1/\lambda$

5.8 T Distribution

Same shape as normal distribution, but different shaped tails (observations more/less likely to fall further than/- closer to mean)

Degrees of freedom (df) affect thickness of tails. Lower df = thicker tails, higher sd. Higher df = closer to normal distribution

T distribution is not skewed.

5.9 Log-normal Distribution

Skewed distributions e.g. length of chess games, adult blood pressures, number of hospitalizations during 2003 SARS outbreak

6 Correlation

6.1 Relationship

positive correlation = same direction, negative correlation = opposite direction

```
1 # Visualizing relationships
2 ggplot(df, aes(x, y)) +
3   geom_point +
4   geom_smooth(method = "lm", se = FALSE) # lm for linear trend line, se = false for no error margin
5
6 # Computing correlation
7 # Usage: corr(vec1, vec2)
8 cor(df$x, df$y)
9
10 cor(df$x, df$y, use = "pairwise.complete.obs") # ignore NA values, otherwise corr returns NA
```

If correlation = 0, there is no relationship and the points will appear to be randomly scattered.

Correlation can be written as r .

If the correlation has a high magnitude, the data points will be clustered closely around a line.

Plot of A against B means A is on the y -axis while B is on the x -axis.

6.2 Caveats of Correlation

Can't account for non-linear relationships e.g. quadratic (U shaped). Correlation shows linear relationships only. Always **visualize** the data.

Log transformation can be used to normalize highly skewed distributions.

```
1 msleep %>%
2   mutate(log_bodywt = log(bodywt)) %>%
3   ggplot(...) + ...
```

Other transformations: square root \sqrt{x} , reciprocal $1 / x$

Correlation does not imply causation. x can be correlated with y but x does not *cause* y .

Spurious correlation: two variables appear to be related but only due to a third, unseen factor, i.e. **confounder** / **lurking variable**.

E.g. holidays are associated with higher retail sales, but it is actually because of more holiday deals (**confounder**)

6.3 Experimental Design

Experiments aim to answer: what is the effect of **treatment** on the **response**? Treatment = explanatory/independent variable Response = dependent variable

E.g. what is the effect of an advertisement on the number of products purchased?

Treatment: advertisement Response: number of products purchased

Controlled experiments - participants are assigned at random to either treatment group (receives treatment) or control group (no treatment).

Groups should be **comparable** so that causation can be inferred. E.g. if average age of control group = 25 while average age of treatment group is 50, age could be a confounder.

Gold standard of experiments:

- ◇ Randomized controlled trial (assign to groups randomly)
- ◇ Placebo - resembles treatment, but has no effect, so participants don't know which group they are in, so causation is due to treatment, not the idea of the treatment (common in clinical trials for drug testing e.g. sugar pills)
- ◇ Double blind trial - person administering treatment also doesn't know if they are administering the actual treatment or placebo

Fewer opportunities for bias = more reliable response

Observational studies - participants assign themselves. Answers questions that are not conducive to a controlled experiment (e.g. can't force people to smoke or have a disease)

However, observational studies only establish association, not causation.

Longitudinal vs Cross-sectional study

Longitudinal: participants followed over a period of time to examine effect of treatment on response

Cross-sectional: data on participants is collected from a single snapshot in time

Longitudinal - more expensive, results take longer. Cross-sectional - cheaper, more convenient