

# Statistics Cheat Sheet

## Population

The entire group one desires information about

## Sample

A subset of the population taken because the entire population is usually too large to analyze  
Its characteristics are taken to be representative of the population

## Mean

Also called the arithmetic mean or average

The sum of all the values in the sample divided by the number of values in the sample/population

$\mu$  is the mean of the population;  $\bar{x}$  is the mean of the sample

## Median

The value separating the higher half of a sample/population from the lower half

Found by arranging all the values from lowest to highest and taking the middle one (or the mean of the middle two if there are an even number of values)

## Variance

Measures dispersion around the mean

Determined by averaging the squared differences of all the values from the mean

Variance of a population is  $\sigma^2$

Can be calculated by subtracting the square of the mean from the average of the squared scores:

$$\sigma^2 = \frac{\sum (x - \mu)^2}{n}$$

Variance of a sample is  $s^2$ ; note the  $n-1$

$$\sigma^2 = \frac{\sum x^2}{n} - \mu^2$$

Can be calculated by:

$$s^2 = \frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}$$

## Standard Deviation

Square root of the variance

Also measures dispersion around the mean but in the same units as the values (instead of square units with variance)

$\sigma$  is the standard deviation of the population and  $s$  is the standard deviation of the sample

## Standard Error

An estimate of the standard deviation of the sampling distribution—the set of all samples of size  $n$  that can be taken from a population

Reflects the extent to which a statistic changes from sample to sample

For a mean,  $\frac{s}{\sqrt{n}}$

For the difference between two means,

Assuming equal variances  $\sqrt{s^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$ ; unequal variances  $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

## T-test

### One-Sample

Tests whether the mean of a normally distributed population is different from a specified value

Null Hypothesis ( $H_0$ ): states that the population mean is equal to some value ( $\mu_0$ )

Alternative Hypothesis ( $H_a$ ): states that the mean does not equal/is greater than/is less than  $\mu_0$

t-statistic: standardizes the difference between  $\bar{x}$  and  $\mu_0$

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}} \quad \text{Degrees of freedom (df)} = n-1$$

Read the table of t-distribution critical values for the p-value (probability that the sample mean was obtained by chance given  $\mu_0$  is the population mean) using the calculated t-statistic and degrees of freedom.

$H_a: \mu > \mu_0 \rightarrow$  the t-statistic is likely positive; read table as given

$H_a: \mu < \mu_0 \rightarrow$  the t-statistic is likely negative; the t-distribution is symmetrical so read the probability as if the t-statistic were positive

Note: if the t-statistic is of the 'wrong' sign, the p-value is 1 minus the  $p$  given in the chart

$H_a: \mu \neq \mu_0 \rightarrow$  read the p-value as if the t-statistic were positive and double it (to consider both less than and greater than)

If the p-value is less than the predetermined value for significance (called  $\alpha$  and is usually 0.05), reject the null hypothesis and accept the alternative hypothesis.

#### Example:

You are experiencing hair loss and skin discoloration and think it might be because of selenium toxicity. You decide to measure the selenium levels in your tap water once a day for one week. Your results are given below. The EPA maximum contaminant level for safe drinking water is 0.05 mg/L. Does the selenium level in your tap water exceed the legal limit (assume  $\alpha=0.05$ )?

Day	Selenium mg/L
1	0.051
2	0.0505
3	0.049
4	0.0516
5	0.052
6	0.0508
7	0.0506

$$H_0: \mu = 0.05; H_a: \mu > 0.05$$

Calculate the mean and standard deviation of your sample:

$$\bar{x} = 0.0508$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1} = \frac{(0.051 - 0.0508)^2 + (0.0505 - 0.0508)^2 + etc...}{6} = 9.15 \times 10^{-7}$$

$$s = \sqrt{s^2} = 9.56 \times 10^{-4}$$

$$\text{The t-statistic is: } t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}} = \frac{0.0508 - 0.05}{\frac{9.56 \times 10^{-4}}{\sqrt{7}}} = 2.17 \text{ and the degrees of freedom are } n-1 = 7-1 = 6$$

Looking at the t-distribution of critical values table, 2.17 with 6 degrees of freedom is between  $p=0.05$  and  $p=0.025$ . This means that the p-value is less than 0.05, so you can reject  $H_0$  and conclude that the selenium level in your tap water exceeds the legal limit.

## T-test

### Two-Sample

Tests whether the means of two populations are significantly different from one another

#### Paired

Each value of one group corresponds directly to a value in the other group; ie: before and after values after drug treatment for each individual patient

Subtract the two values for each individual to get one set of values (the differences) and use  $\mu_0 = 0$  to perform a one-sample t-test

#### Unpaired

The two populations are independent

$H_0$ : states that the means of the two populations are equal ( $\mu_1 = \mu_2$ )

$H_a$ : states that the means of the two populations are unequal or one is greater than the other ( $\mu_1 \neq \mu_2, \mu_1 > \mu_2, \mu_1 < \mu_2$ )

t-statistic:

$$\text{assuming equal variances: } t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad \text{assuming unequal variances: } t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)}}$$

degrees of freedom =  $(n_1-1)+(n_2-1)$

Read the table of t-distribution critical values for the p-value using the calculated t-statistic and degrees of freedom. Remember to keep the sign of the t-statistic clear (order of subtracting the sample means) and to double the p-value for an  $H_a$  of  $\mu_1 \neq \mu_2$ .

**Example:**

Consider the lifespan of 18 rats. 12 were fed a restricted calorie diet and lived an average of 700 days (standard deviation=21 days). The other 6 had unrestricted access to food and lived an average of 668 days (standard deviation=30 days). Does a restricted calorie diet increase the lifespan of rats (assume  $\alpha=0.05$ )?

$$\mu_1=700, s_1=21, n_1=12; \mu_2=668, s_2=30, n_2=6$$

$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 > \mu_2 \text{ (because we are only asking if a restricted calorie diet increases lifespan)}$$

We cannot assume that the variances of the two populations are equal because the different diets could also affect the variability in lifespan.

$$\text{The t-statistic is: } t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)}} = \frac{700 - 668}{\sqrt{\frac{21^2}{12} + \frac{30^2}{6}}} = 2.342$$

$$\text{Degrees of freedom} = (n_1-1)+(n_2-1) = (12-1)+(6-1)=16$$

From the t-distribution table, the p-value falls between 0.01 and 0.02, so we do reject  $H_0$ . The restricted calorie diet does increase the lifespan of rats.

## Chi-Square Test

### For Goodness of Fit

Checks whether or not an observed pattern of data fits some given distribution

$$H_0: \text{the observed pattern fits the given distribution}$$

$$H_a: \text{the observed pattern does not fit the given distribution}$$

$$\text{The chi-square statistic is: } \chi^2 = \sum \frac{(O - E)^2}{E} \quad (O \text{ is the observed value and } E \text{ is the expected value})$$

Degrees of freedom = number of categories in the distribution – 1

Get the p-value from the table of  $\chi^2$  critical values using the calculated  $\chi^2$  and df values. If the p-value is less than  $\alpha$ , the observed data does not fit the expected distribution. If  $p>\alpha$ , the data likely fits the expected distribution

**Example 1:**

You breed puffskeins and would like to determine the pattern of inheritance for coat color and purring ability.

Puffskeins come in either pink or purple and can either purr or hiss. You breed a purebred, pink purring male with a purebred, purple hissing female. All individuals of the  $F_1$  generation are pink and purring. The  $F_2$  offspring are shown below. Do the alleles for coat color and purring ability assort independently (assume  $\alpha=0.05$ )?

Pink and Purring	Pink and Hissing	Purple and Purring	Purple and Hissing
143	60	55	18

Independent assortment means a phenotypic ratio of 9:3:3:1, so:

$$H_0: \text{the observed distribution of } F_2 \text{ offspring fits a 9:3:3:1 distribution}$$

$$H_a: \text{the observed distribution of } F_2 \text{ offspring does not fit a 9:3:3:1 distribution}$$

The expected values are:

Pink and Purring	Pink and Hissing	Purple and Purring	Purple and Hissing
155.25	51.75	51.75	17.25

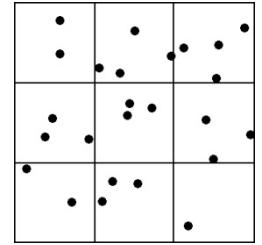
$$\chi^2 = \sum \frac{(O - E)^2}{E} = \frac{(143 - 155.25)^2}{155.25} + \frac{(60 - 51.75)^2}{51.75} + \frac{(55 - 51.75)^2}{51.75} + \frac{(18 - 17.25)^2}{17.25} = 2.519$$

df=4-1=3

From the table of  $\chi^2$  critical values, the p-value is greater than 0.25, so the alleles for coat color and purring ability do assort independently in puffskeins.

### Example 2:

You are studying the pattern of dispersion of king penguins and the diagram on the right represents an area you sampled. Each dot is a penguin. Do the penguins display a uniform distribution (assume  $\alpha=0.05$ )?



$H_0$ : there is a uniform distribution of penguins

$H_a$ : there is not a uniform distribution of penguins

There are a total of 25 penguins, so if there is a uniform distribution, there should be 2.778 penguins per square. The actual observed values are 2, 4, 4, 3, 3, 3, 2, 3, 1, so the  $\chi^2$  statistic is:

$$\chi^2 = \sum \frac{(O - E)^2}{E} = \frac{(1 - 2.778)^2}{2.778} + 2\left(\frac{(2 - 2.778)^2}{2.778}\right) + 4\left(\frac{(3 - 2.778)^2}{2.778}\right) + 2\left(\frac{(4 - 2.778)^2}{2.778}\right) + \left(\frac{(3 - 2.778)^2}{2.778}\right) = 2.72$$

df=9-1=8

From the table of  $\chi^2$  critical values, the p-value is greater than 0.25, so we do not reject  $H_0$ . The penguins do display a uniform distribution.

## Chi-Square Test

### For Independence

Checks whether two categorical variables are related or not (independence)

$H_0$ : the two variables are independent

$H_a$ : the two variables are not independent

Does not make any assumptions about an expected distribution

The observed values ( $\#_1, \#_2, \#_3$ , and  $\#_4$ ) are usually presented as a table. Each row is a category of variable 1 and each column is a category of variable 2.

		Variable 1		Totals
		Category x	Category y	
Variable 2	Category a	$\#_1$	$\#_2$	$\#_1 + \#_2$
	Category b	$\#_3$	$\#_4$	$\#_3 + \#_4$
Totals		$\#_1 + \#_3$	$\#_2 + \#_4$	$\#_1 + \#_2 + \#_3 + \#_4$

The proportion of category  $x$  of variable 1 is the number of individuals in category  $x$  divided by the total number of individuals  $\left(\frac{\#_1 + \#_3}{\#_1 + \#_2 + \#_3 + \#_4}\right)$ . Assuming independence, the expected number of individuals that fall within category  $x$  of variable 2 is the proportion of category  $x$  multiplied by the number of individuals in category  $a$   $\left(\frac{\#_1 + \#_3}{\#_1 + \#_2 + \#_3 + \#_4}\right)(\#_1 + \#_2)$ . Thus, the expected value is:

$$E = \frac{(\#_1 + \#_3)(\#_1 + \#_2)}{\#_1 + \#_2 + \#_3 + \#_4} = \frac{(row\ total)(column\ total)}{grand\ total}$$

Degrees of freedom =  $(r-1)(c-1)$  where  $r$  is the number of rows and  $c$  is the number of columns

The chi-square statistic is still  $\chi^2 = \sum \frac{(O - E)^2}{E}$

Read the p-values from the table of  $\chi^2$  critical values.

### Example:

Given the data below, is there a relationship between fitness level and smoking habits (assume  $\alpha=0.05$ )?

		Fitness Level				
		Low	Medium-Low	Medium-High	High	
Never smoked		113	113	110	159	495
Former smokers		119	135	172	190	616
1 to 9 cigarettes daily		77	91	86	65	319
$\geq 10$ cigarettes daily		181	152	124	73	530
		490	491	492	487	1960

$H_0$ : fitness level and smoking habits are independent

$H_a$ : fitness level and smoking habits are not independent

First, we calculate the expected counts. For the first cell, the expected count is:

$$E = \frac{(row\ total)(column\ total)}{grand\ total} = \frac{(495)(490)}{1960} = 123.75$$

	Fitness Level			
	Low	Medium-Low	Medium-High	High
Never smoked	123.75	124	124.26	122.99
Former smokers	154	154.31	154.63	153.06
1 to 9 cigarettes daily	79.75	79.91	80.08	79.26
$\geq 10$ cigarettes daily	132.5	132.77	133.04	131.69

$$\chi^2 = \sum \frac{(O - E)^2}{E} = \frac{(113 - 123.75)^2}{123.75} + \frac{(113 - 124)^2}{124} + \frac{(110 - 124.26)^2}{124.26} + etc... = 91.73$$

$$df = (r-1)(c-1) = (4-1)(4-1) = 9$$

From the table of  $\chi^2$  critical values, the p-value is less than 0.001, so we reject  $H_0$  and conclude that there is a relationship between fitness level and smoking habits.

### Type I error

The probability of rejecting a true null hypothesis

Equals  $\alpha$

### Type II error

The probability of failing to reject a false null hypothesis

## Probability

### Joint Probability

The probability of events A and B occurring

$$P(A \text{ and } B) = P(A) \times P(B) \text{ when events A and B are independent}$$

### Union of Events

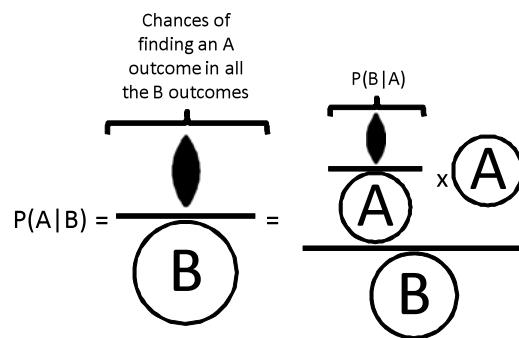
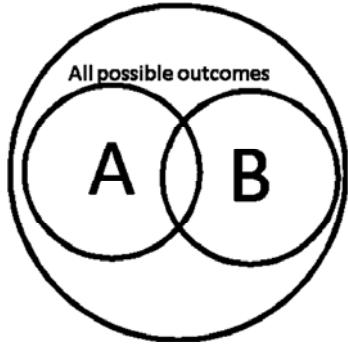
The probability of either event A or event B occurring

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

### Conditional Probability

The probability of event A occurring given that event B has occurred

$$P(A | B) = \frac{P(A \text{ and } B)}{P(B)} \quad \text{or} \quad P(A | B) = \frac{P(B | A) \times P(A)}{P(B)}$$



### Example 1:

Assume that eye color is an autosomally inherited trait controlled by one gene with two alleles. Brown is dominant to blue. A brown-eyed man with genotype Bb and a blue-eyed woman have three children. The first has blue eyes. What is the probability that all three children have blue eyes?

Without considering the first child, the probability that the couple has three children with blue eyes is  $0.5 \times 0.5 \times 0.5 = 0.125 = P(A \text{ and } B) = P(2 \text{ children} = bb \text{ and } 1\text{st child } bb)$

With his parents, the probability that the 1<sup>st</sup> child is bb is:  $P(B) = P(1\text{st child} = bb) = 0.5$

$$\text{Therefore, } P(2 \text{ children} = bb \mid 1\text{st child } bb) = P(A \mid B) = \frac{P(A \text{ and } B)}{P(B)} = \frac{0.125}{0.5} = 0.25$$

### **Example 2:**

Based on an analysis of her pedigree, it is determined that a woman has a 70% chance of being Zz and a 30% chance of being ZZ for a sex-linked trait, where Z is dominant to z. If she now has a son with the Z phenotype, what is the probability of her being Zz?

We're looking for:  $P(W=Zz \mid S=Z)$

But it's hard to find  $P(W=Zz \text{ and } S=Z)$  because the two events are not independent. Instead, let us use:

$$P(A \mid B) = \frac{P(B \mid A) \times P(A)}{P(B)}$$

$P(S = Z \mid W = Zz) = 0.5$  (50% chance of passing on the Z allele)

$P(W = Zz) = 0.7$  (given)

$P(S = Z) = (0.7 \times 0.5) + (0.3 \times 1) = 0.65$  (son can be Z from the woman being either Zz or ZZ)

$$P(W = Zz \mid S = Z) = \frac{0.5 \times 0.7}{0.65} = 0.538$$

## **Multiple Experiments**

### **Binomial distribution**

For when you are not concerned about the order of the events, only that they occur

$$P(X = m) = \frac{n! \times p^m \times (1-p)^{(n-m)}}{m! \times (n-m)!}$$

for  $m$  outcomes of event X in  $n$  total trials with  $p$ =probability of X occurring once

### **Example:**

What is the probability that a couple has one boy out of five children?

$$P(1 \text{ boy of } 5 \text{ children}) = \frac{5! \times 0.5^1 \times 0.5^4}{1! \times (4)!} = 0.15625$$

### **Poisson distribution**

The binomial distribution works for a small number of trials but as  $n$  gets too large, the factorials become unwieldy.

The Poisson distribution is an estimate of the binomial distribution for large  $n$ .

$$P(X = m) = \frac{e^{-np} \times (n \times p)^m}{m!}$$

Note:  $np$  is also known as the number of expected outcomes for event X

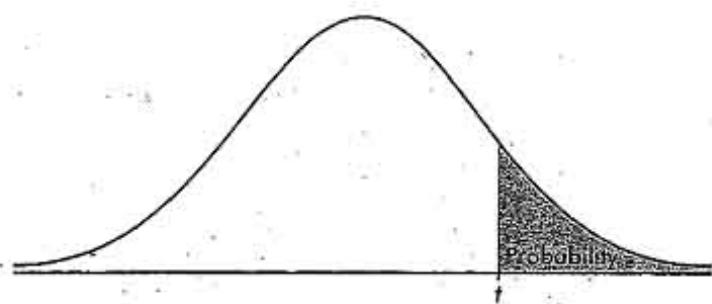
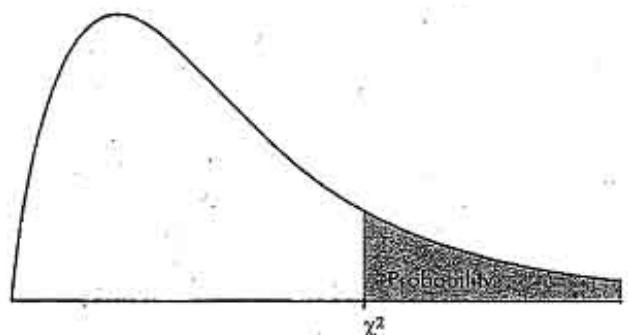
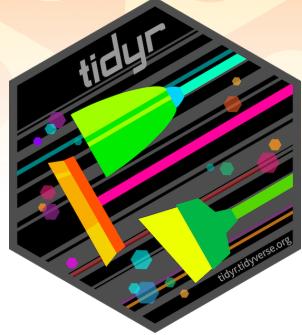


TABLE B:  $t$ -DISTRIBUTION CRITICAL VALUES

$\chi^2$  CRITICAL VALUESTABLE C:  $\chi^2$  CRITICAL VALUES

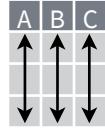
df	Tail probability $p$										
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001
1	1.32	1.64	2.07	2.71	3.84	5.02	5.41	6.63	7.88	9.14	10.83
2	2.77	3.22	3.79	4.61	5.99	7.38	7.82	9.21	10.60	11.98	13.82
3	4.11	4.64	5.32	6.25	7.81	9.35	9.84	11.34	12.84	14.32	16.27
4	5.39	5.99	6.74	7.78	9.49	11.14	11.67	13.28	14.86	16.42	18.47
5	6.63	7.29	8.12	9.24	11.07	12.83	13.39	15.09	16.75	18.39	20.51
6	7.84	8.56	9.45	10.64	12.59	14.45	15.03	16.81	18.55	20.25	22.46
7	9.04	9.80	10.75	12.02	14.07	16.01	16.62	18.48	20.28	22.04	24.32
8	10.22	11.03	12.03	13.36	15.51	17.53	18.17	20.09	21.95	23.77	26.12
9	11.39	12.24	13.29	14.68	16.92	19.02	19.68	21.67	23.59	25.46	27.88
10	12.55	13.44	14.53	15.99	18.31	20.48	21.16	23.21	25.19	27.11	29.59
11	13.70	14.63	15.77	17.28	19.68	21.92	22.62	24.72	26.76	28.73	31.26
12	14.85	15.81	16.99	18.55	21.03	23.34	24.05	26.22	28.30	30.32	32.91
13	15.98	16.98	18.20	19.81	22.36	24.74	25.47	27.69	29.82	31.88	34.53
14	17.12	18.15	19.41	21.06	23.68	26.12	26.87	29.14	31.32	33.43	36.12
15	18.25	19.31	20.60	22.31	25.00	27.49	28.26	30.58	32.80	34.95	37.70
16	19.37	20.47	21.79	23.54	26.30	28.85	29.63	32.00	34.27	36.46	39.25
17	20.49	21.61	22.98	24.77	27.59	30.19	31.00	33.41	35.72	37.95	40.79
18	21.60	22.76	24.16	25.99	28.87	31.53	32.35	34.81	37.16	39.42	42.31
19	22.72	23.90	25.33	27.20	30.14	32.85	33.69	36.19	38.58	40.88	43.82
20	23.83	25.04	26.50	28.41	31.41	34.17	35.02	37.57	40.00	42.34	45.31
21	24.93	26.17	27.66	29.62	32.67	35.48	36.34	38.93	41.40	43.78	46.80
22	26.04	27.30	28.82	30.81	33.92	36.78	37.66	40.29	42.80	45.20	48.27
23	27.14	28.43	29.98	32.01	35.17	38.08	38.97	41.64	44.18	46.62	49.73
24	28.24	29.55	31.13	33.20	36.42	39.36	40.27	42.98	45.56	48.03	51.18
25	29.34	30.68	32.28	34.38	37.65	40.65	41.57	44.31	46.93	49.44	52.62
26	30.43	31.79	33.43	35.56	38.89	41.92	42.86	45.64	48.29	50.83	54.05
27	31.53	32.91	34.57	36.74	40.11	43.19	44.14	46.96	49.64	52.22	55.48
28	32.62	34.03	35.71	37.92	41.34	44.46	45.42	48.28	50.99	53.59	56.89
29	33.71	35.14	36.85	39.09	42.56	45.72	46.69	49.59	52.34	54.97	58.30
30	34.80	36.25	37.99	40.26	43.77	46.98	47.96	50.89	53.67	56.33	59.70
40	45.62	47.27	49.24	51.81	55.76	59.34	60.44	63.69	66.77	69.70	73.40
50	56.33	58.16	60.35	63.17	67.50	71.42	72.61	76.15	79.49	82.66	86.66
60	66.98	68.97	71.34	74.40	79.08	83.30	84.58	88.38	91.95	95.34	99.61
80	88.13	90.41	93.11	96.58	101.9	106.6	108.1	112.3	116.3	120.1	124.8
100	109.1	111.7	114.7	118.5	124.3	129.6	131.1	135.8	140.2	144.3	149.4

# Data tidying with `tidyr` :: CHEAT SHEET



**Tidy data** is a way to organize tabular data in a consistent data structure across packages.

A table is tidy if:



Each **variable** is in its own **column**

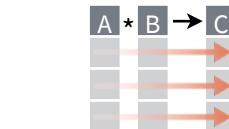
&



Each **observation**, or **case**, is in its own row



Access **variables** as **vectors**



Preserve **cases** in vectorized operations

## Tibbles

### AN ENHANCED DATA FRAME

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with `[`, a vector with `[[` and `$`.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

`options(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf)` Control default display settings.

`View()` or `glimpse()` View the entire data set.

### CONSTRUCT A TIBBLE

**tibble(...)** Construct by columns.

`tibble(x = 1:3, y = c("a", "b", "c"))`

Both make this tibble

A tibble: 3 × 2  
`x` <int> <chr>  
 1 1 a  
 2 2 b  
 3 3 c

**as\_tibble(x, ...)** Convert a data frame to a tibble.

**enframe(x, name = "name", value = "value")**

Convert a named vector to a tibble. Also `deframe()`.

**is\_tibble(x)** Test whether x is a tibble.

## Reshape Data

- Pivot data to reorganize values into a new layout.

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M
C	1999	212K	1T
C	2000	213K	1T

## Split Cells

- Use these functions to split or combine cells into individual, isolated values.

table5

country	century	year
A	19	99
A	20	00
B	19	99
B	20	00



country	year
A	1999
A	2000
B	1999
B	2000

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M



country	year	rate
A	1999	0.7K
A	1999	19M
A	2000	2K
A	2000	20M
B	1999	37K
B	1999	172M
B	2000	80K
B	2000	174M

**pivot\_longer**(data, cols, names\_to = "name", values\_to = "value", values\_drop\_na = FALSE)

"Lengthen" data by collapsing several columns into two. Column names move to a new names\_to column and values to a new values\_to column.

`pivot_longer(table4a, cols = 2:3, names_to = "year", values_to = "cases")`

**pivot\_wider**(data, names\_from = "name", values\_from = "value")

The inverse of `pivot_longer()`. "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

`pivot_wider(table2, names_from = type, values_from = count)`

## Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	

→ A 1 3

A 2 1

B 1 2

B 2 2

**expand**(data, ...) Create a new tibble with all possible combinations of the values listed in ... Drop other variables.

`expand(mtcars, cyl, gear, carb)`

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	
			NA

→ A 1 3

A 2 1 NA

B 1 4

B 2 3

**complete**(data, ..., fill = list()) Add missing possible combinations of values of variables listed in ... Fill remaining variables with NA.

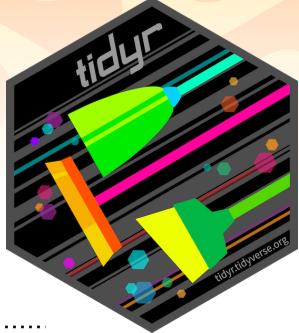
`complete(mtcars, cyl, gear, carb)`

x	x1	x2
A	1	
B	NA	
C	NA	
D	3	
E	NA	

→ A 1

D 3

B NA



# Nested Data

A **nested data frame** stores individual tables as a list-column of data frames within a larger organizing data frame. List-columns can also be lists of vectors or lists of varying data types.

Use a nested data frame to:

- Preserve relationships between observations and subsets of data. Preserve the type of the variables being nested (factors and datetimes aren't coerced to character).
- Manipulate many sub-tables at once with **purrr** functions like `map()`, `map2()`, or `pmap()` or with **dplyr** `rowwise()` grouping.

## CREATE NESTED DATA

**nest(data, ...)** Moves groups of cells into a list-column of a data frame. Use alone or with `dplyr::group_by()`:

1. Group the data frame with `group_by()` and use `nest()` to move the groups into a list-column.

```
n_storms <- storms %>%
  group_by(name) %>%
  nest()
```

2. Use `nest(new_col = c(x, y))` to specify the columns to group using `dplyr::select()` syntax.

```
n_storms <- storms %>%
  nest(data = c(year:long))
```

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

Index list-columns with `[[[]]]`. `n_storms$data[[1]]`

## CREATE TIBBLES WITH LIST-COLUMNS

**tibble::tribble(...)** Makes list-columns when needed.

```
tibble(~max, ~seq,
      3, 1:3,
      4, 1:4,
      5, 1:5)
```

max	seq
3	<int [3]>
4	<int [4]>
5	<int [5]>

**tibble::tibble(...)** Saves list input as list-columns.

```
tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))
```

**tibble::enframe(x, name="name", value="value")**

Converts multi-level list to a tibble with list-cols.  
`enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')`

## OUTPUT LIST-COLUMNS FROM OTHER FUNCTIONS

**dplyr::mutate(), transmute(), and summarise()** will output list-columns if they return a list.

```
mtcars %>%
  group_by(cyl) %>%
  summarise(q = list(quantile(mpg)))
```

## RESHAPE NESTED DATA

**unnest(data, cols, ..., keep\_empty = FALSE)** Flatten nested columns back to regular columns. The inverse of `nest()`.  
`n_storms %>% unnest(data)`

**unnest\_longer(data, col, values\_to = NULL, indices\_to = NULL)**  
Turn each element of a list-column into a row.

```
starwars %>%
  select(name, films) %>%
  unnest_longer(films)
```

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

"cell" contents

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	data
Amy	<tibble [50x3]>
Bob	<tibble [50x3]>
Zeta	<tibble [50x3]>

name	films
Luke	<chr [5]>
C-3PO	<chr [6]>
R2-D2	<chr[7]>

→

name	films
Luke	The Empire Strik...
Luke	Revenge of the S...
Luke	Return of the Jed...
C-3PO	The Empire Strik...
C-3PO	Attack of the Cl...
C-3PO	The Phantom M...
R2-D2	The Empire Strik...
R2-D2	Attack of the Cl...
R2-D2	The Phantom M...

**unnest\_wider(data, col)** Turn each element of a list-column into a regular column.

```
starwars %>%
  select(name, films) %>%
  unnest_wider(films)
```

name	films
Luke	<chr [5]>
C-3PO	<chr [6]>
R2-D2	<chr[7]>

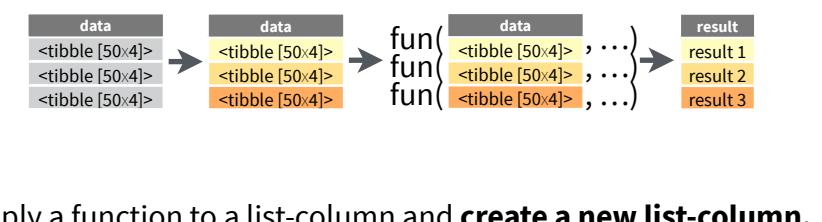
→

name	..1	..2	..3
Luke	The Empire...	Revenge of...	Return of...
C-3PO	The Empire...	Attack of...	The Phantom...
R2-D2	The Empire...	Attack of...	The Phantom...

## TRANSFORM NESTED DATA

A vectorized function takes a vector, transforms each element in parallel, and returns a vector of the same length. By themselves vectorized functions cannot work with lists, such as list-columns.

**dplyr::rowwise(.data, ...)** Group data so that each row is one group, and within the groups, elements of list-columns appear directly (accessed with `[[ ]]`, not as lists of length one. **When you use `rowwise()`, dplyr functions will seem to apply functions to list-columns in a vectorized fashion.**



Apply a function to a list-column and **create a new list-column**.

`n_storms %>% rowwise() %>% mutate(n = list(dim(data)))`

**dim()** returns two values per row  
wrap with `list` to tell `mutate` to create a list-column

Apply a function to a list-column and **create a regular column**.

`n_storms %>% rowwise() %>% mutate(n = nrow(data))`

**nrow()** returns one integer per row

**Collapsing multiple list-columns into a single list-column.**

`starwars %>% rowwise() %>% mutate(transport = list(append(vehicles, starships)))`

**append()** returns a list for each row, so col type must be list

Apply a function to **multiple list-columns**.

`starwars %>% rowwise() %>% mutate(n_transports = length(c(vehicles, starships)))`

**length()** returns one integer per row

See **pur**

# String manipulation with stringr :: CHEAT SHEET



The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

## Detect Matches

	<b>str_detect(string, pattern, negate = FALSE)</b> Detect the presence of a pattern match in a string. Also <b>str_like()</b> . str_detect(fruit, "a")
	<b>str_starts(string, pattern, negate = FALSE)</b> Detect the presence of a pattern match at the beginning of a string. Also <b>str_ends()</b> . str_starts(fruit, "a")
	<b>str_which(string, pattern, negate = FALSE)</b> Find the indexes of strings that contain a pattern match. str_which(fruit, "a")
	<b>str_locate(string, pattern)</b> Locate the positions of pattern matches in a string. Also <b>str_locate_all()</b> . str_locate(fruit, "a")
	<b>str_count(string, pattern)</b> Count the number of matches in a string. str_count(fruit, "a")

## Subset Strings

	<b>str_sub(string, start = 1L, end = -1L)</b> Extract substrings from a character vector. str_sub(fruit, 1, 3); str_sub(fruit, -2)
	<b>str_subset(string, pattern, negate = FALSE)</b> Return only the strings that contain a pattern match. str_subset(fruit, "p")
	<b>str_extract(string, pattern)</b> Return the first pattern match found in each string, as a vector. Also <b>str_extract_all()</b> to return every pattern match. str_extract(fruit, "[aeiou]")
	<b>str_match(string, pattern)</b> Return the first pattern match found in each string, as a matrix with a column for each () group in pattern. Also <b>str_match_all()</b> . str_match(sentences, "(a the) ([^ +])")

## Manage Lengths

	<b>str_length(string)</b> The width of strings (i.e. number of code points, which generally equals the number of characters). str_length(fruit)
	<b>str_pad(string, width, side = c("left", "right", "both"), pad = " ")</b> Pad strings to constant width. str_pad(fruit, 17)
	<b>str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...")</b> Truncate the width of strings, replacing content with ellipsis. str_trunc(sentences, 6)
	<b>str_trim(string, side = c("both", "left", "right"))</b> Trim whitespace from the start and/or end of a string. str_trim(str_pad(fruit, 17))
	<b>str_squish(string)</b> Trim whitespace from each end and collapse multiple spaces into single spaces. str_squish(str_pad(fruit, 17, "both"))

## Mutate Strings

	<b>str_sub()</b> <- value. Replace substrings by identifying the substrings with str_sub() and assigning into the results. str_sub(fruit, 1, 3) <- "str"
	<b>str_replace(string, pattern, replacement)</b> Replace the first matched pattern in each string. Also <b>str_remove()</b> . str_replace(fruit, "p", "-")
	<b>str_replace_all(string, pattern, replacement)</b> Replace all matched patterns in each string. Also <b>str_remove_all()</b> . str_replace_all(fruit, "p", "-")
	<b>str_to_lower(string, locale = "en")<sup>1</sup></b> Convert strings to lower case. str_to_lower(sentences)
	<b>str_to_upper(string, locale = "en")<sup>1</sup></b> Convert strings to upper case. str_to_upper(sentences)
	<b>str_to_title(string, locale = "en")<sup>1</sup></b> Convert strings to title case. Also <b>str_to_sentence()</b> . str_to_title(sentences)

## Join and Split

	<b>str_c(..., sep = "", collapse = NULL)</b> Join multiple strings into a single string. str_c(letters, LETTERS)
	<b>str_flatten(string, collapse = "")</b> Combines into a single string, separated by collapse. str_flatten(fruit, ",")
	<b>str_dup(string, times)</b> Repeat strings times times. Also <b>str_unique()</b> to remove duplicates. str_dup(fruit, times = 2)
	<b>str_split_fixed(string, pattern, n)</b> Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also <b>str_split()</b> to return a list of substrings and <b>str_split_n()</b> to return the nth substring. str_split_fixed(sentences, " ", n=3)
	<b>str_glue(..., .sep = "", .envir = parent.frame())</b> Create a string from strings and {expressions} to evaluate. str_glue("Pi is {pi}")
	<b>str_glue_data(.x, ..., .sep = "", .envir = parent.frame(), .na = "NA")</b> Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. str_glue_data(mtcars, "[rownames(mtcars)] has {hp} hp")

## Order Strings

	<b>str_order(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)<sup>1</sup></b> Return the vector of indexes that sorts a character vector. fruit[str_order(fruit)]
	<b>str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)<sup>1</sup></b> Sort a character vector. str_sort(fruit)

## Helpers

	<b>str_conv(string, encoding)</b> Override the encoding of a string. str_conv(fruit, "ISO-8859-1")
	<b>str_view_all(string, pattern, match = NA)</b> View HTML rendering of all regex matches. Also <b>str_view()</b> to see only the first match. str_view_all(sentences, "[aeiou]")
	<b>str_equal(x, y, locale = "en", ignore_case = FALSE, ...)<sup>1</sup></b> Determine if two strings are equivalent. str_equal(c("a", "b"), c("a", "c"))
	<b>str_wrap(string, width = 80, indent = 0, exdent = 0)</b> Wrap strings into nicely formatted paragraphs. str_wrap(sentences, 20)

<sup>1</sup> See [bit.ly/ISO639-1](https://bit.ly/ISO639-1) for a complete list of locales.



# Shiny :: CHEAT SHEET



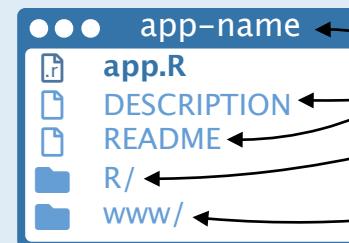
## Building an App

A **Shiny** app is a web page (**ui**) connected to a computer running a live R session (**server**).



Users can manipulate the UI, which will cause the server to update the UI's displays (by running R code).

Save your template as **app.R**. Keep your app in a directory along with optional extra files.



Launch apps stored in a directory with **runApp(<path to directory>)**.

To generate the template, type **shinyapp** and press **Tab** in the RStudio IDE or go to **File > New Project > New Directory > Shiny Web Application**

```
# app.R
library(shiny)

ui <- fluidPage(
  numericInput(inputId = "n",
    "Sample size", value = 25),
  plotOutput(outputId = "hist")
)

server <- function(input, output, session) {
  output$hist <- renderPlot({
    hist(rnorm(input$n))
  })
}

shinyApp(ui = ui, server = server)
```

In **ui** nest R functions to build an HTML interface

Customize the UI with **Layout Functions**

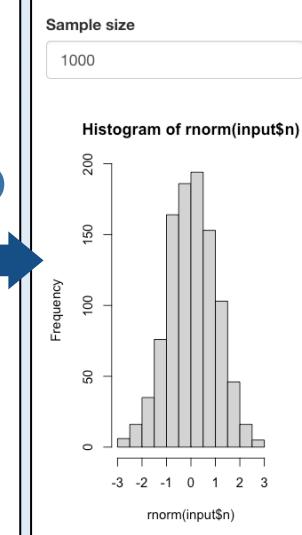
Add Inputs with **\*Input()** functions

Add Outputs with **\*Output()** functions

Tell the **server** how to render outputs and respond to inputs with R

Wrap code in **render\*()** functions before saving to output

Refer to UI inputs with **input\$<id>** and outputs with **output\$<id>**



Call **shinyApp()** to combine **ui** and **server** into an interactive app!

See annotated examples of Shiny apps by running **runExample(<example name>)**. Run **runExample()** with no arguments for a list of example names.

## Share

Share your app in three ways:

1. **Host it on shinyapps.io**, a cloud based service from RStudio. To deploy Shiny apps:

Create a free or professional account at [shinyapps.io](https://shinyapps.io)

Click the Publish icon in RStudio IDE, or run: **rsconnect::deployApp("<path to directory>")**

2. **Purchase RStudio Connect**, a publishing platform for R and Python. [rstudio.com/products/connect/](https://rstudio.com/products/connect/)

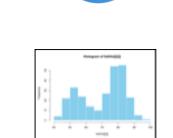
3. **Build your own Shiny Server** [rstudio.com/products/shiny/shiny-server/](https://rstudio.com/products/shiny/shiny-server/)

## Outputs

**render\*()** and **\*Output()** functions work together to add R output to the UI.



**DT::renderDataTable(expr, options, searchDelay, callback, escape, env, quoted, outputArgs)**



**renderImage(expr, env, quoted, deleteFile, outputArgs)**



**renderPrint(expr, env, quoted, width, outputArgs)**



**renderTable(expr, striped, hover, bordered, spacing, width, align, rownames, colnames, digits, na, ..., env, quoted, outputArgs)**

**renderText(expr, env, quoted, outputArgs, sep)**

**renderUI(expr, env, quoted, outputArgs)**

**dataTableOutput(outputId)**

**imageOutput(outputId, width, height, click, dblclick, hover, brush, inline)**

**plotOutput(outputId, width, height, click, dblclick, hover, brush, inline)**

**verbatimTextOutput(outputId, placeholder)**

**tableOutput(outputId)**

**textOutput(outputId, container, inline)**

**uiOutput(outputId, inline, container, ...)**

**htmlOutput(outputId, inline, container, ...)**

## Inputs

Collect values from the user.

Access the current value of an input object with **input\$<inputId>**. Input values are **reactive**.

Action

**actionButton(inputId, label, icon, width, ...)**

Link

**actionLink(inputId, label, icon, ...)**

checkbox

**checkboxGroupInput(inputId, label, choices, selected, inline, width, choiceNames, choiceValues)**

checkbox

**checkboxInput(inputId, label, value, width)**

date

**dateInput(inputId, label, value, min, max, format, startview, weekstart, language, width, autoclose, datesdisabled, daysofweekdisabled)**

dateRange

**dateRangeInput(inputId, label, start, end, min, max, format, startview, weekstart, language, separator, width, autoclose)**

file

**fileInput(inputId, label, multiple, accept, width, buttonLabel, placeholder)**

numeric

**numericInput(inputId, label, value, min, max, step, width)**

password

**passwordInput(inputId, label, value, width, placeholder)**

radio

**radioButtons(inputId, label, choices, selected, inline, width, choiceNames, choiceValues)**

select

**selectInput(inputId, label, choices, selected, multiple, selectize, width, size)**  
Also **selectizeInput()**

slider

**sliderInput(inputId, label, min, max, value, step, round, format, locale, ticks, animate, width, sep, pre, post, timeFormat, timezone, dragRange)**

submit

**submitButton(text, icon, width)**  
(Prevent reactions for entire app)

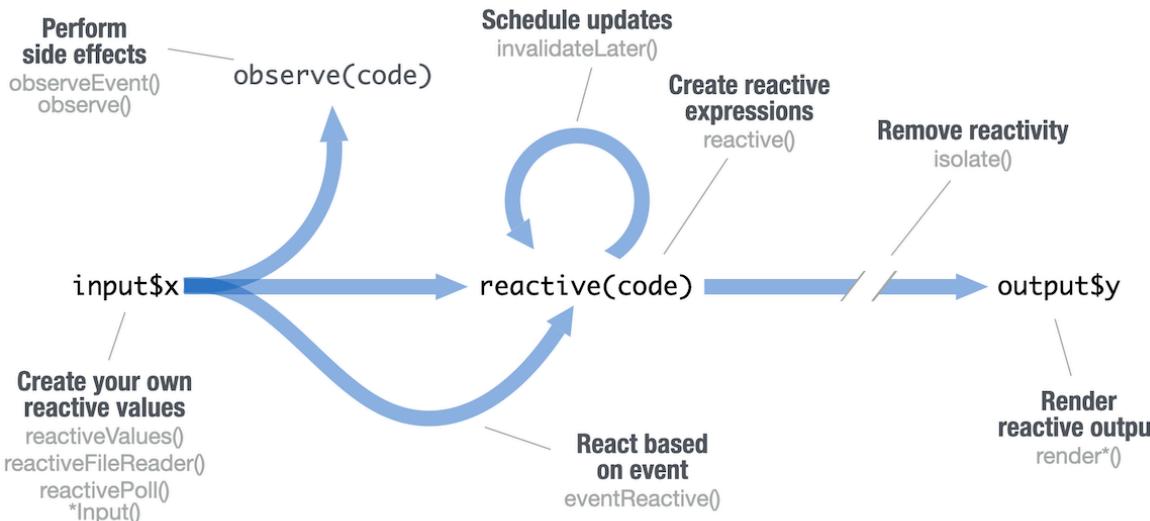
text

**textInput(inputId, label, value, width, placeholder)** Also **textAreaInput()**

These are the core output types. See [htmlwidgets.org](http://htmlwidgets.org) for many more options.

# Reactivity

Reactive values work together with reactive functions. Call a reactive value from within the arguments of one of these functions to avoid the error **Operation not allowed without an active reactive context**.



## CREATE YOUR OWN REACTIVE VALUES

```
# *Input() example
ui <- fluidPage(
  textInput("a", "", "A")
)
```

```
#reactiveValues example
server <-
  function(input, output){
    rv <- reactiveValues()
    rv$number <- 5
  }
```

### \*Input() functions (see front page)

Each input function creates a reactive value stored as **input\$<inputId>**.

### reactiveValues(...)

Creates a list of reactive values whose values you can set.

## CREATE REACTIVE EXPRESSIONS

```
library(shiny)
ui <- fluidPage(
 textInput("a", "", "A"),
 textInput("z", "", "Z"),
  textOutput("b"))
server <-
  function(input, output){
    re <- reactive({
      paste(input$a, input$z)
    })
    output$b <- renderText({
      re()
    })
  }
shinyApp(ui, server)
```

### reactive(x, env, quoted, label, domain)

#### Reactive expressions:

- cache their value to reduce computation
- can be called elsewhere
- notify dependencies when invalidated
- Call the expression with function syntax, e.g. **re()**.

## PERFORM SIDE EFFECTS

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  actionButton("go", "Go"))
server <-
  function(input, output){
    observeEvent(input$go, {
      print(input$a)
    })
  }
shinyApp(ui, server)
```

### observeEvent(eventExpr,

**handlerExpr**, **event.env**, **event.quoted**, **handler.env**, **handler.quoted**, ..., **label**, **suspended**, **priority**, **domain**, **autoDestroy**, **ignoreNULL**, **ignoreInit**, **once**)

Runs code in 2nd argument when reactive values in 1st argument change. See **observe()** for alternative.

## REACT BASED ON EVENT

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  actionButton("go", "Go"),
  textOutput("b"))
server <-
  function(input, output){
    re <- eventReactive(
      input$go, {input$a})
    output$b <- renderText({
      re()
    })
  }
shinyApp(ui, server)
```

### eventReactive(eventExpr, valueExpr, event.env, event.quoted, value.env, value.quoted, ..., label, domain, ignoreNULL, ignoreInit)

Creates reactive expression with code in 2nd argument that only invalidates when reactive values in 1st argument change.

## REMOVE REACTIVITY

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  textOutput("b"))
server <-
  function(input, output){
    output$b <- renderText({
      isolate({input$a})
    })
  }
shinyApp(ui, server)
```

### isolate(expr)

Runs a code block. Returns a **non-reactive** copy of the results.

# UI

An app's UI is an HTML document.

Use Shiny's functions to assemble this HTML with R.

```
fluidPage(
  textInput("a", ""))
## <div class="container-fluid">
##   <div class="form-group shiny-input-container">
##     <label for="a"></label>
##     <input id="a" type="text"
##           class="form-control" value="" />
##   </div>
## </div>
```

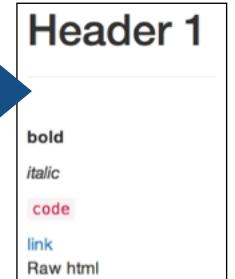
Returns HTML

 Add static HTML elements with **tags**, a list of functions that parallel common HTML tags, e.g. **tags\$a()**. Unnamed arguments will be passed into the tag; named arguments will become tag attributes.

Run **names(tags)** for a complete list.  
**tags\$h1("Header")** → `<h1>Header</h1>`

The most common tags have wrapper functions. You do not need to prefix their names with **tags\$**

```
ui <- fluidPage(
  h1("Header 1"),
  hr(),
  br(),
  p(strong("bold")),
  p(em("italic")),
  p(code("code")),
  a(href="", "link"),
  HTML("<p>Raw html</p>"))
)
```



 To include a CSS file, use **includeCSS()**, or 1. Place the file in the **www** subdirectory 2. Link to it with:

```
tags$head(tags$link(rel = "stylesheet",
  type = "text/css", href = "<file name>"))
```

 To include JavaScript, use **includeScript()** or 1. Place the file in the **www** subdirectory 2. Link to it with:

```
tags$head(tags$script(src = "<file name>"))
```

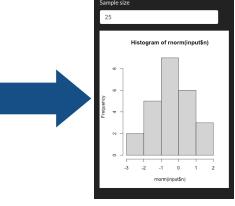
 To include an image:

1. Place the file in the **www** subdirectory 2. Link to it with **img(src = "<file name>")**

# Themes

Use the **bslib** package to add existing themes to your Shiny app ui, or make your own.

```
library(bslib)
ui <- fluidPage(
  theme = bs_theme(
    bootswatch = "darkly",
    ...
  )
)
```



**bootswatch\_themes()** Get a list of themes.

# Layouts

Combine multiple elements into a "single element" that has its own properties with a panel function, e.g.



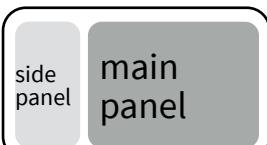
```
wellPanel(
  dateInput("a", ""),
  submitButton()
)
```

```
absolutePanel()
conditionalPanel()
fixedPanel()
headerPanel()
inputPanel()
mainPanel()
```

```
navlistPanel()
sidebarPanel()
tabPanel()
tabsetPanel()
titlePanel()
wellPanel()
```

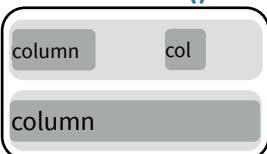
Organize panels and elements into a layout with a layout function. Add elements as arguments of the layout functions.

### sidebarLayout()



```
ui <- fluidPage(
  sidebarLayout(
    sidebarPanel(),
    mainPanel()
  )
)
```

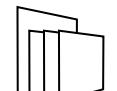
### fluidRow()



```
ui <- fluidPage(
  fluidRow(column(width = 4),
    column(width = 2, offset = 3)),
  fluidRow(column(width = 12))
)
```

Also **flowLayout()**, **splitLayout()**, **verticalLayout()**, **fixedPage()**, and **fixedRow()**.

Layer tabPanels on top of each other, and navigate between them, with:



```
ui <- fluidPage( tabsetPanel(
```

```
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
)
```

```
ui <- fluidPage( navlistPanel(
```

```
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
)
```

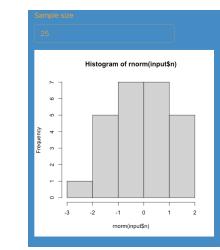
```
ui <- navbarPage(title = "Page",
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
```



Build your own theme by customizing individual arguments.

**bs\_theme(bg = "#558AC5", fg = "#F9B02D", ...)**

?**bs\_theme** for a full list of arguments.



**bs\_themer()** Place within the server function to use the interactive theming widget.

# rmarkdown :: CHEAT SHEET

## What is rmarkdown?



**.Rmd files** • Develop your code and ideas side-by-side in a single document. Run code as individual chunks or as an entire document.

**Dynamic Documents** • Knit together plots, tables, and results with narrative text. Render to a variety of formats like HTML, PDF, MS Word, or MS Powerpoint.

**Reproducible Research** • Upload, link to, or attach your report to share. Anyone can read or run your code to reproduce your work.

## Workflow

- 1 Open a **new .Rmd file** in the RStudio IDE by going to *File > New File > R Markdown*.
- 2 **Embed code** in chunks. Run code by line, by chunk, or all at once.
- 3 **Write text** and add tables, figures, images, and citations. Format with Markdown syntax or the RStudio Visual Markdown Editor.
- 4 **Set output format(s) and options** in the YAML header. Customize themes or add parameters to execute or add interactivity with Shiny.
- 5 **Save and render** the whole document. Knit periodically to preview your work as you write.
- 6 **Share your work!**

## Embed Code with knitr

### CODE CHUNKS

Surround code chunks with `{{r}}` and `{{` or use the Insert Code Chunk button. Add a chunk label and/or chunk options inside the curly braces after **r**.

```
```{r chunk-label, include=FALSE}
summary(mtcars)
```
```

### SET GLOBAL OPTIONS

Set options for the entire document in the first chunk.

```
```{r include=FALSE}
knitr::opts_chunk$message = FALSE
```
```

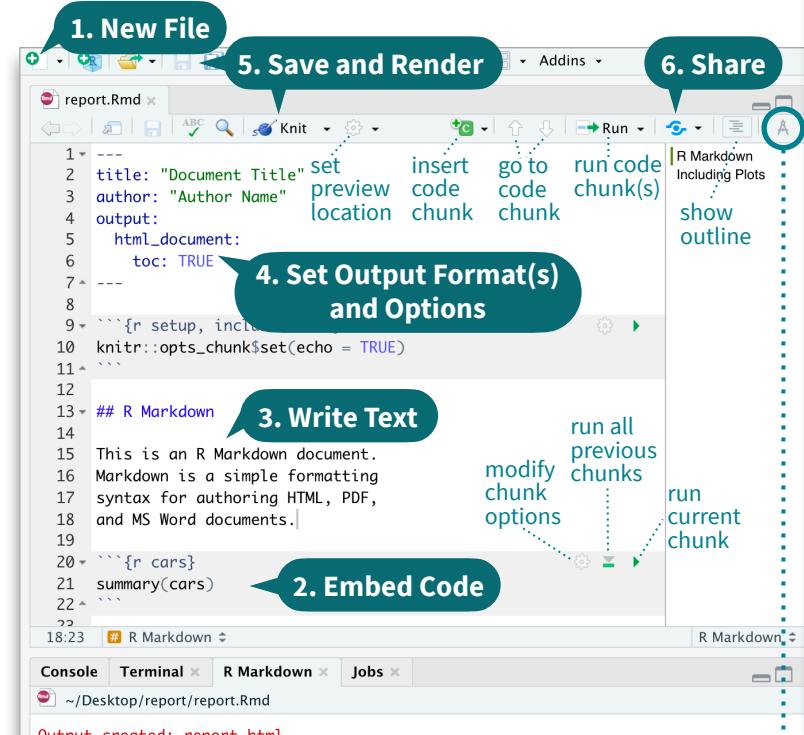
### INLINE CODE

Insert `r <code>` into text sections. Code is evaluated at render and results appear as text.

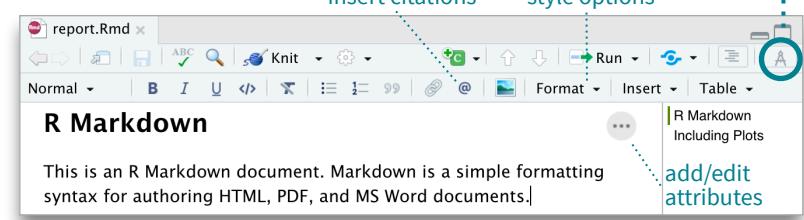
"Built with `r getRversion()`" --> "Built with 4.1.0"



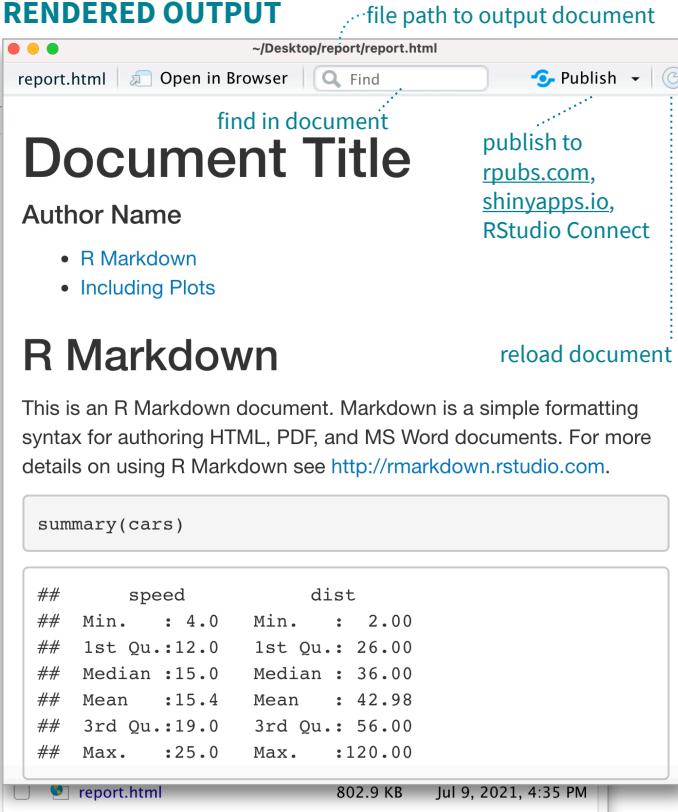
### SOURCE EDITOR



### VISUAL EDITOR



### RENDERED OUTPUT



## Write with Markdown

The syntax on the left renders as the output on the right.

Plain text.

Plain text.

End a line with two spaces to start a new paragraph.

End a line with two spaces to start a new paragraph.

Also end with a backslash\ to make a new line.

Also end with a backslash\ to make a new line.

**italics\*** and **\*\*bold\*\***

**italics** and **bold**

superscript<sup>2</sup>/subscript<sub>2</sub>

superscript<sup>2</sup>/subscript<sub>2</sub>

~~strikethrough~~

strikethrough

escaped: `\*` \`

escaped: \* \`

endash: --, emdash: ---

endash: -, emdash: --

### Header 1 Header 2

...

Header 6

- unordered list

• item 2

- item 2a (indent 1 tab)

• item 2b

1. ordered list

2. item 2

- item 2a (indent 1 tab)

• item 2b

<link url>

[This is a link.](link url)

[This is another link][id].

This is another link.

<http://www.rstudio.com/>

This is a link.

This is another link.



Caption.

verbatim code

multiple lines of verbatim code

> block quotes

block quotes

equation:  $e^{i\pi} + 1 = 0$

equation block:

$$E = mc^2$$

horizontal rule:

| Right | Left | Default | Center |
|-------|------|---------|--------|
| 12    | 12   | 12      | 12     |
| 123   | 123  | 123     | 123    |
| 1     | 1    | 1       | 1      |

### HTML Tabsets

```
# Results {.tabset}
## Plots text
text
```

## Tables
more text

### Results

|       |        |
|-------|--------|
| Plots | Tables |
|-------|--------|

text





# Set Output Formats and their Options in YAML

Use the document's YAML header to set an **output format** and customize it with **output options**.

```
---
```

```
title: "My Document"
author: "Author Name"
output:
  html_document:
    toc: TRUE
---
```

**Indent format 2 characters,  
indent options 4 characters**

| OUTPUT FORMAT           | CREATES                      |
|-------------------------|------------------------------|
| html_document           | .html                        |
| pdf_document*           | .pdf                         |
| word_document           | Microsoft Word (.docx)       |
| powerpoint_presentation | Microsoft Powerpoint (.pptx) |
| odt_document            | OpenDocument Text            |
| rtf_document            | Rich Text Format             |
| md_document             | Markdown                     |
| github_document         | Markdown for Github          |
| ioslides_presentation   | ioslides HTML slides         |
| slidy_presentation      | Slidy HTML slides            |
| beamer_presentation*    | Beamer slides                |

\* Requires LaTeX, use `tinytex::install_tinytex()`  
Also see `flexdashboard`, `bookdown`, `distill`, and `blogdown`.

| IMPORTANT OPTIONS   | DESCRIPTION  | HTML    | PDF | MS Word | MS PPT |
|---------------------|--|---------|-----|---------|--------|
| anchor_sections     | Show section anchors on mouse hover (TRUE or FALSE)                                    | X       |     |         |        |
| citation_package    | The LaTeX package to process citations ("default", "natbib", "biblatex")               | X       |     |         |        |
| code_download       | Give readers an option to download the .Rmd source code (TRUE or FALSE)                | X       |     |         |        |
| code_folding        | Let readers to toggle the display of R code ("none", "hide", or "show")                | X       |     |         |        |
| css                 | CSS or SCSS file to use to style document (e.g. "style.css")                           | X       |     |         |        |
| dev                 | Graphics device to use for figure output (e.g. "png", "pdf")                           | X X     |     |         |        |
| df_print            | Method for printing data frames ("default", "kable", "tibble", "paged")                | X X X X |     |         |        |
| fig_caption         | Should figures be rendered with captions (TRUE or FALSE)                               | X X X X |     |         |        |
| highlight           | Syntax highlighting ("tango", "pygments", "kate", "zenburn", "textmate")               | X X X   |     |         |        |
| includes            | File of content to place in doc ("in_header", "before_body", "after_body")             | X X     |     |         |        |
| keep_md             | Keep the Markdown .md file generated by knitting (TRUE or FALSE)                       | X X X X |     |         |        |
| keep_tex            | Keep the intermediate TEX file used to convert to PDF (TRUE or FALSE)                  | X       |     |         |        |
| latex_engine        | LaTeX engine for producing PDF output ("pdflatex", "xelatex", or "lualatex")           | X       |     |         |        |
| reference_docx/_doc | docx/pptx file containing styles to copy in the output (e.g. "file.docx", "file.pptx") | X X     |     |         |        |
| theme               | Theme options (see Bootswatch and Custom Themes below)                                 | X       |     |         |        |
| toc                 | Add a table of contents at start of document (TRUE or FALSE)                           | X X X X |     |         |        |
| toc_depth           | The lowest level of headings to add to table of contents (e.g. 2, 3)                   | X X X X |     |         |        |
| toc_float           | Float the table of contents to the left of the main document content (TRUE or FALSE)   | X       |     |         |        |

Use `?<output format>` to see all of a format's options, e.g. `?html_document`

## More Header Options

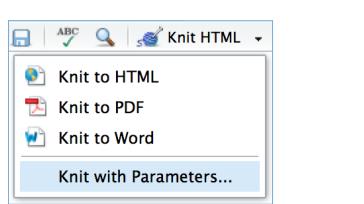
### PARAMETERS

Parameterize your documents to reuse with new inputs (e.g., data, values, etc.).

1. **Add parameters** in the header as sub-values of `params`.
2. **Call parameters** in code using `params$<name>`.
3. **Set parameters** with Knit with Parameters or the `params` argument of `render()`.

### REUSABLE TEMPLATES

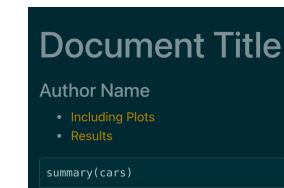
1. **Create a new package** with a `inst/rmarkdown/templates` directory.
2. **Add a folder** containing `template.yaml` (below) and `skeleton.Rmd` (template contents).
3. **Install** the package to access template by going to **File > New R Markdown > From Template**.



### BOOTSWATCH THEMES

Customize HTML documents with Bootswatch themes from the `bslib` package using the theme output option.

Use `bslib::bootswatch_themes()` to list available themes.



```
---
```

```
title: "Document Title"
author: "Author Name"
output:
  html_document:
    theme:
      bootswatch: solar
---
```

### CUSTOM THEMES

Customize individual HTML elements using `bslib` variables. Use `?bs_theme` to see more variables.

```
---
```

```
output:
  html_document:
    theme:
      bg: "#121212"
      fg: "#E4E4E4"
      base_font:
        google: "Prompt"
---
```

More on `bslib` at [pkgs.rstudio.com/bslib/](https://pkgs.rstudio.com/bslib/).

### STYLING WITH CSS AND SCSS

Add CSS and SCSS to your document by adding a path to a file with the `css` option in the YAML header.

```
---
```

```
title: "My Document"
author: "Author Name"
output:
  html_document:
    css: "style.css"
---
```

Apply CSS styling by writing HTML tags directly or:

- Use markdown to apply style attributes inline.

Bracketed Span  
A [green]{.my-color} word.

A green word.

Fenced Div  
:::{.my-color}  
All of these words  
are green.  
:::

All of these words  
are green.

- Use the Visual Editor. Go to **Format > Div/Span** and add CSS styling directly with Edit Attributes.

.my-css-tag ...  
This is a div with some text in it.

## Render

When you render a document, rmarkdown:

1. Runs the code and embeds results and text into an .md file with knitr.
2. Converts the .md file into the output format with Pandoc.



**Save**, then **Knit** to preview the document output. The resulting HTML/PDF/MS Word/etc. document will be created and saved in the same directory as the .Rmd file.

Use `rmarkdown::render()` to render/knit in the R console. See `?render` for available options.

## Share

### Publish on RStudio Connect

to share R Markdown documents securely, schedule automatic updates, and interact with parameters in real time.

[rstudio.com/products/connect/](https://rstudio.com/products/connect/)



### INTERACTIVITY

Turn your report into an interactive Shiny document in 4 steps:

1. Add `runtime: shiny` to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render with `rmarkdown::run()` or click **Run Document** in RStudio IDE.

```
---
```

```
output: html_document
runtime: shiny
---
```

```
```{r, echo = FALSE}
numericInput("n",
  "How many cars?", 5)
renderTable({
  head(cars, input$n)
})
```

	speed	dist
1	4.00	2.00
2	4.00	10.00
3	7.00	4.00
4	7.00	22.00
5	8.00	16.00

Also see Shiny Prerendered for better performance.

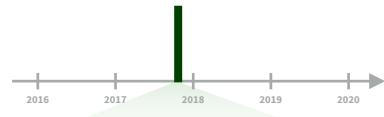
[rmarkdown.rstudio.com/authoring\\_shiny\\_prerendered](https://rmarkdown.rstudio.com/authoring_shiny_prerendered)

Embed a complete app into your document with `shiny::shinyAppDir()`. More at [bookdown.org/yihui/rmarkdown/shiny-embedded.html](https://bookdown.org/yihui/rmarkdown/shiny-embedded.html).

# Dates and times with lubridate :: CHEAT SHEET



## Date-times



2017-11-28 12:00:00

2017-11-28 12:00:00

A **date-time** is a point on the timeline, stored as the number of seconds since 1970-01-01 00:00:00 UTC

```
dt <- as_datetime(1511870400)  
## "2017-11-28 12:00:00 UTC"
```

### PARSE DATE-TIMES (Convert strings or numbers to date-times)

1. Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
2. Use the function below whose name replicates the order. Each accepts a tz argument to set the time zone, e.g. ymd(x, tz = "UTC").

2017-11-28T14:02:00

ymd\_hms(), ymd\_hm(), ymd\_h().  
ymd\_hms("2017-11-28T14:02:00")

2017-22-12 10:00:00

ydm\_hms(), ydm\_hm(), ydm\_h().  
ydm\_hms("2017-22-12 10:00:00")

11/28/2017 1:02:03

mdy\_hms(), mdy\_hm(), mdy\_h().  
mdy\_hms("11/28/2017 1:02:03")

1 Jan 2017 23:59:59

dmy\_hms(), dmy\_hm(), dmy\_h().  
dmy\_hms("1 Jan 2017 23:59:59")

20170131

ymd(), ydm(). ymd(20170131)

July 4th, 2000

mdy(), myd(). mdy("July 4th, 2000")

4th of July '99

dmy(), dym(). dmy("4th of July '99")

2001: Q3

yq() Q for quarter. yq("2001: Q3")

07-2020

my(), ym(). my("07-2020")

2:01

hms::hms() Also lubridate::hms(), hm() and ms(), which return periods.\* hms::hms(sec = 0, min = 1, hours = 2, roll = FALSE)

2017.5

date\_decimal(decimal, tz = "UTC")  
date\_decimal(2017.5)

now(zone = "") Current time in tz (defaults to system tz). now()

today(zone = "") Current date in a tz (defaults to system tz). today()

fast.strptime() Faster strftime.

fast.strptime('9/1/01', '%y/%m/%d')

parse\_date\_time() Easier strftime.

parse\_date\_time("9/1/01", "ymd")

2017-11-28

A **date** is a day stored as the number of days since 1970-01-01

```
d <- as_date(17498)  
## "2017-11-28"
```

12:00:00

An hms is a **time** stored as the number of seconds since 00:00:00

```
t <- hms::as.hms(85)  
## 00:01:25
```

### GET AND SET COMPONENTS

Use an accessor function to get a component. Assign into an accessor function to change a component in place.

```
d ## "2017-11-28"  
day(d) ## 28  
day(d) <- 1  
d ## "2017-11-01"
```

2018-01-31 11:59:59

**date(x)** Date component. date(dt)

2018-01-31 11:59:59

**year(x)** Year. year(dt)  
**isoyear(x)** The ISO 8601 year.  
**epiyear(x)** Epidemiological year.

2018-01-31 11:59:59

**month(x, label, abbr)** Month. month(dt)

2018-01-31 11:59:59

**day(x)** Day of month. day(dt)  
**wday(x, label, abbr)** Day of week.  
**qday(x)** Day of quarter.

2018-01-31 11:59:59

**hour(x)** Hour. hour(dt)

2018-01-31 11:59:59

**minute(x)** Minutes. minute(dt)

2018-01-31 11:59:59

**second(x)** Seconds. second(dt)

2018-01-31 11:59:59 UTC

**tz(x)** Time zone. tz(dt)

2018-01-31 11:59:59

**week(x)** Week of the year. week(dt)  
**isoweek()** ISO 8601 week.  
**epiweek()** Epidemiological week.

2018-01-31 11:59:59

**quarter(x)** Quarter. quarter(dt)

2018-01-31 11:59:59

**semester(x, with\_year = FALSE)** Semester. semester(dt)

2018-01-31 11:59:59

**am(x)** Is it in the am? am(dt)  
**pm(x)** Is it in the pm? pm(dt)

2018-01-31 11:59:59

**dst(x)** Is it daylight savings? dst(dt)

2018-01-31 11:59:59

**leap\_year(x)** Is it a leap year?  
leap\_year(dt)

2018-01-31 11:59:59

**update(object, ..., simple = FALSE)**  
update(dt, mday = 2, hour = 1)



## Round Date-times



**floor\_date(x, unit = "second")**  
Round down to nearest unit.  
floor\_date(dt, unit = "month")

**round\_date(x, unit = "second")**  
Round to nearest unit.  
round\_date(dt, unit = "month")

**ceiling\_date(x, unit = "second", change\_on\_boundary = NULL)**  
Round up to nearest unit.  
ceiling\_date(dt, unit = "month")

Valid units are second, minute, hour, day, week, month, bimonth, quarter, season, halfyear and year.

**rollback(dates, roll\_to\_first = FALSE, preserve\_hms = TRUE)** Roll back to last day of previous month. Also **rollforward()**. rollback(dt)

## Stamp Date-times

**stamp()** Derive a template from an example string and return a new function that will apply the template to date-times. Also **stamp\_date()** and **stamp\_time()**.

1. Derive a template, create a function  
sf <- stamp("Created Sunday, Jan 17, 1999 3:34")

2. Apply the template to dates  
sf(ymd("2010-04-05"))  
## [1] "Created Monday, Apr 05, 2010 00:00"

**Tip:** use a date with day > 12

## Time Zones

R recognizes ~600 time zones. Each encodes the time zone, Daylight Savings Time, and historical calendar variations for an area. R assigns one time zone per vector.

Use the **UTC** time zone to avoid Daylight Savings.

**OlsonNames()** Returns a list of valid time zone names. OlsonNames()

**Sys.timezone()** Gets current time zone.

5:00 Mountain 6:00 Central 7:00 Eastern  
4:00 Pacific

PT MT CT ET  
7:00 Pacific 7:00 Eastern

7:00 Mountain 7:00 Central  
7:00 Eastern

**with\_tz(time, tzzone = "")** Get the same date-time in a new time zone (a new clock time). Also **local\_time(dt, tz, units)**. with\_tz(dt, "US/Pacific")

**force\_tz(time, tzzone = "")** Get the same clock time in a new time zone (a new date-time). Also **force\_tzs()**. force\_tz(dt, "US/Pacific")





# Math with Date-times

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

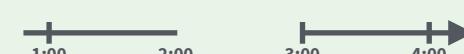
A normal day

```
nor <- ymd_hms("2018-01-01 01:30:00",tz="US/Eastern")
```



The start of daylight savings (spring forward)

```
gap <- ymd_hms("2018-03-11 01:30:00",tz="US/Eastern")
```



The end of daylight savings (fall back)

```
lap <- ymd_hms("2018-11-04 00:30:00",tz="US/Eastern")
```



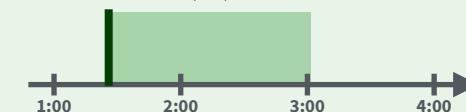
Leap years and leap seconds

```
leap <- ymd("2019-03-01")
```

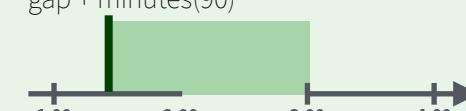


**Periods** track changes in clock times, which ignore time line irregularities.

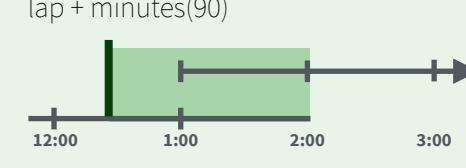
nor + minutes(90)



gap + minutes(90)



lap + minutes(90)

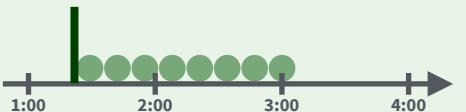


leap + years(1)



**Durations** track the passage of physical time, which deviates from clock time when irregularities occur.

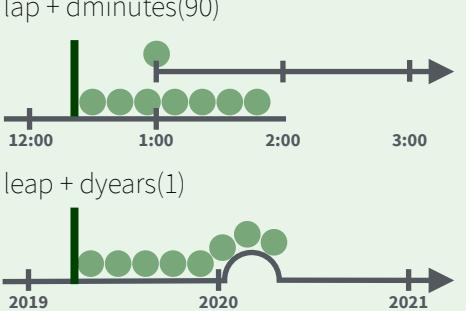
nor + dminutes(90)



gap + dminutes(90)



lap + dminutes(90)



leap + dyears(1)



**Intervals** represent specific intervals of the timeline, bounded by start and end date-times.

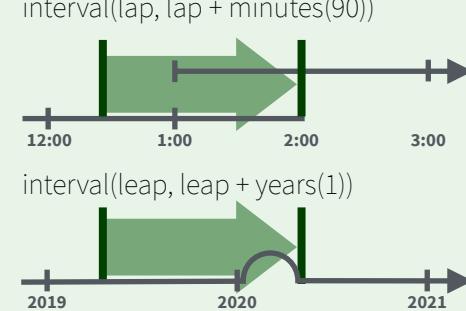
interval(nor, nor + minutes(90))



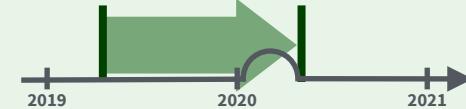
interval(gap, gap + minutes(90))



interval(lap, lap + minutes(90))



interval(leap, leap + years(1))



Not all years are 365 days due to **leap days**.

Not all minutes are 60 seconds due to **leap seconds**.

It is possible to create an imaginary date by adding **months**, e.g. February 31st

```
jan31 <- ymd(20180131)
```

```
jan31 + months(1)
```

```
## "NA"
```

%m+% and %m-% will roll imaginary dates to the last day of the previous month.

```
jan31 %m+% months(1)
```

```
## "2018-02-28"
```

**add\_with\_rollback(e1, e2, roll\_to\_first = TRUE)** will roll imaginary dates to the first day of the new month.

```
add_with_rollback(jan31, months(1), roll_to_first = TRUE)
```

```
## "2018-03-01"
```

## PERIODS

Add or subtract periods to model events that happen at specific clock times, like the NYSE opening bell.

Make a period with the name of a time unit **pluralized**, e.g.

```
p <- months(3) + days(12)
```

```
p
```

```
"3m 12d 0H 0M 0S"
```

Number of months Number of days etc.

**years(x = 1)** x years.

**months(x)** x months.

**weeks(x = 1)** x weeks.

**days(x = 1)** x days.

**hours(x = 1)** x hours.

**minutes(x = 1)** x minutes.

**seconds(x = 1)** x seconds.

**milliseconds(x = 1)** x milliseconds.

**microseconds(x = 1)** x microseconds.

**nanoseconds(x = 1)** x nanoseconds.

**picoseconds(x = 1)** x picoseconds.

**period(num = NULL, units = "second", ...)**

An automation friendly period constructor.  
period(5, unit = "years")

**as.period(x, unit)** Coerce a timespan to a period, optionally in the specified units. Also **is.period()**. as.period(i)

**period\_to\_seconds(x)** Convert a period to the "standard" number of seconds implied by the period. Also **seconds\_to\_period()**.  
period\_to\_seconds(p)

## DURATIONS

Add or subtract durations to model physical processes, like battery life. Durations are stored as seconds, the only time unit with a consistent length.

**Diftimes** are a class of durations found in base R.

Make a duration with the name of a period prefixed with a **d**, e.g.

```
dd <- ddays(14)
```

```
dd
```

```
"1209600s (~2 weeks)"
```

Exact length in seconds Equivalent in common units

**dyears(x = 1)** 31536000x seconds.

**dmonths(x = 1)** 2629800x seconds.

**dweeks(x = 1)** 604800x seconds.

**ddays(x = 1)** 86400x seconds.

**dhours(x = 1)** 3600x seconds.

**dminutes(x = 1)** 60x seconds.

**dseconds(x = 1)** x seconds.

**dmilliseconds(x = 1)** x × 10<sup>-3</sup> seconds.

**dmicroseconds(x = 1)** x × 10<sup>-6</sup> seconds.

**dnanoseconds(x = 1)** x × 10<sup>-9</sup> seconds.

**dpicoseconds(x = 1)** x × 10<sup>-12</sup> seconds.

**duration(num = NULL, units = "second", ...)**

An automation friendly duration constructor. duration(5, unit = "years")

**as.duration(x, ...)** Coerce a timespan to a duration. Also **is.duration()**, **is.difftime()**. as.duration(i)

**make\_difftime(x)** Make difftime with the specified number of units. make\_difftime(99999)

## INTERVALS

Divide an interval by a duration to determine its physical length, divide an interval by a period to determine its implied length in clock time.

Make an interval with **interval()** or %--%, e.g.

```
i <- interval(ymd("2017-01-01"), d)
```

```
## 2017-01-01 UTC--2017-11-28 UTC
```

```
j <- d %--% ymd("2017-12-31")
```

```
## 2017-11-28 UTC--2017-12-31 UTC
```



Start Date End Date

# Leaflet Cheat Sheet



an open-source JavaScript library for mobile-friendly interactive maps

## Quick Start

### Installation

Use `install.packages("leaflet")` to install the package or directly from Github `devtools::install_github("rstudio/leaflet")`.

### First Map

```
m <- leaflet() %>%
  addTiles() %>%
  addMarkers(lng = 174.768, lat = -36.852, popup = "The birthplace of R")
# add a single point layer
```



## Map Widget

### Initialization

<code>m &lt;- leaflet(options = leafletOptions(...))</code>	Initial geographic center of the map
<code>center</code>	Initial map zoom level
<code>zoom</code>	Minimum zoom level of the map
<code>minZoom</code>	Maximum zoom level of the map
<code>maxZoom</code>	

### Map Methods

```
m %>% setView(lng, lat, zoom, options = list())
# Set the view of the map (center and zoom level)
m %>% fitBounds(lng1, lat1, lng2, lat2)
# Fit the view into the rectangle [lng1, lat1] - [lng2, lat2]
m %>% clearBounds()
# Clear the bound, automatically determine from the map elements
```

### Data Object

Both `leaflet()` and the `map` layers have an optional data parameter that is designed to receive spatial data with the following formats:

#### Base R

The arguments of all layers take normal R objects:

```
df <- data.frame(lat = ..., lng = ...)
```

```
leaflet(df) %>% addTiles() %>% addCircles()
```

library(sp) Useful functions:

SpatialPoints, SpatialLines, SpatialPolygons, ...

library(maps) Build a map of states with colors:

```
mapStates <- map("state", fill = TRUE, plot = FALSE)
```

```
leaflet(mapStates) %>% addTiles() %>%
```

```
addPolygons(fillColor = topo.colors(10, alpha = NULL), stroke = FALSE)
```

## Markers

Use markers to call out points, express locations with latitude/longitude coordinates, appear as icons or as circles.

Data come from vectors or assigned data frame, or `sp` package objects.

### Icon Markers

Regular Icons: default and simple

```
addMarkers(lng, lat, popup, label) add basic icon markers
```

```
makeIcon(Icons(iconUrl, iconWidth, iconHeight, iconAnchorX, iconAnchorY,
  shadowUrl, shadowWidth, shadowHeight, ...)) customize marker icons
```

```
iconList() create a list of icons
```

Awesome Icons: customizable with colors and icons

```
addAwesomeMarkers, makeAwesomeIcon, awesomeIcons, awesomeIconList
```

Marker Clusters: option of `addMarkers()`

```
clusterOptions = markerClusterOptions()
```

```
freezeAtZoom Freeze the cluster at assigned zoom level
```

### Circle Markers

```
addCircleMarkers(color, radius, stroke, opacity, ...)
```

Customize their color, radius, stroke, opacity

## Popups and Labels

`addPopups(lng, lat, ...content..., options)` Add standalone popups

```
options = popupOptions(closeButton=FALSE)
```

`addMarkers(..., popup, ...)` Show popups with markers or shapes

`addMarkers(..., label, labelOptions...)` Show labels with markers or shapes

```
labelOptions = labelOptions(noHide, textOnly, textSize, direction, style)
```

`addLabelOnlyMarkers()` Add labels without markers

## Lines and Shapes

### Polygons and Polylines

`addPolygons(color, weight=1, smoothFactor=0.5, opacity=1.0, fillOpacity=0.5,`  
`fillColor= ~colorQuantile("YlOrRd", ALAND)(ALAND), highlightOptions, ... )`

`highlightOptions(color, weight=2, bringToFront=TRUE)` highlight shapes

Use `rmapshaper::ms_simplify` to simplify complex shapes

`Circles addCircles(lng, lat, weight=1, radius, ...)`

`Rectangles addRectangles(lng1, lat1, lng2, lat2, fillColor="transparent", ... )`

## Basemaps

`addTiles()`

Default Tiles

Use `addTiles()` to add a custom map tile URL template, use `addWMSTiles()` to add WMS (Web Map Service) tiles

`providers$Stamen.Toner, CartoDB.Positron, Esri.NatGeoWorldMap`

Third-Party Tiles

`addProviderTiles()`

## GeoJSON and TopoJSON

There are two options to use the GeoJSON/TopoJSON data:

\* To read into `sp` objects with the `geojsonio` or `rgdal` package:  
`geojsonio::geojson_read(..., what="sp") rgdal::readOGR(..., "OGRGeoJSON")`

\* Or to use the `addGeoJSON()` and `addTopoJSON()` functions:  
`addTopoJSON/addGeoJSON(... weight, color, fill, opacity, fillOpacity...)`

Styles can also be tuned separately with a `style: {}` object.  
Other packages including `RJSONIO` and `jsonlite` can help fast parse or generate the data needed.

## Shiny Integration

To integrate a Leaflet map into an app:

\* In the UI, call `leafletOutput("name")`

\* On the server side, assign a `renderLeaflet(...)` call to the output

\* Inside the `renderLeaflet` expression, return a Leaflet map object

### Modification

To modify an existing map or add incremental changes to the map, you can use `leafletProxy()`. This should be performed in an observer on the server side.

Other useful functions to edit your map:

`fitBounds(o, 0, 11, 11)` similar to `setView`

fit the view to within these bounds

`addCircles(1:10, 1:10, layerId = LETTERS[1:10])`

create circles with layerIds of "A", "B", "C"...

`removeShape(c("B", "F"))` remove some of the circles

`clearShapes()` clear all circles (and other shapes)

### Inputs/Events

#### Object Events

Object event names generally use this pattern:

`inputs$MAPID_OBJCATEGORY_EVENTNAME`.

Triger an event changes the value of the Shiny input at this variable.

Valid values for `OBJCATEGORY` are `marker`, `shape`, `geojson` and `topojson`.

Valid values for `EVENTNAME` are `click`, `mouseover` and `mouseout`.

All of these events are set to either `NULL` if the event has never happened, or a `list()` that includes:

\* `lat` The latitude of the object, if available; otherwise, the mouse cursor

\* `lng` The longitude of the object, if available; otherwise, the mouse cursor

\* `id` The layerId, if any

GeoJSON events also include additional properties:

\* `featureId` The feature ID, if any

\* `properties` The feature properties

#### Map Events

`inputs$MAPID_click`

when the map background or basemap is clicked  
`value -- a list with lat and lng`

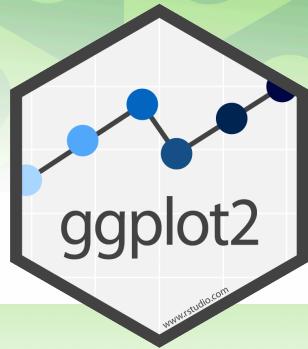
`inputs$MAPID_bounds`

provide the lat/long bounds of the visible map area  
`value -- a list with north, east, south and west`

`inputs$MAPID_zoom`

an integer indicates the zoom level

# Data visualization with ggplot2 :: CHEAT SHEET



## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data** set, a **coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
  <GEOM_FUNCTION>(mapping = aes(<Mappings>),
  stat = <STAT>, position = <POSITION>) +
  <COORDINATE_FUNCTION> +
  <FACET_FUNCTION> +
  <SCALE_FUNCTION> +
  <THEME_FUNCTION>
```

required  
Not required, sensible defaults supplied

`ggplot(data = mpg, aes(x = cty, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

`last_plot()` Returns the last plot.

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

## Aes Common aesthetic values.

**color** and **fill** - string ("red", "#RRGGBB")

**linetype** - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

**lineend** - string ("round", "butt", or "square")

**linejoin** - string ("round", "mitre", or "bevel")

**size** - integer (line width in mm)

**shape** - integer/shape name or a single character ("a")

## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables.  
Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
```

- a + geom\_blank()** and **a + expand\_limits()**  
Ensure limits include values across all plots.
- b + geom\_curve(aes(yend = lat + 1, xend = long + 1, curvature = 1))** - x, yend, alpha, angle, curvature, linetype, size
- a + geom\_path(lineend = "butt", linejoin = "round", linemitre = 1)** - x, y, alpha, color, group, linetype, size
- a + geom\_polygon(aes(alpha = 50))** - x, y, alpha, color, fill, group, subgroup, linetype, size
- b + geom\_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom\_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

- b + geom\_abline(aes(intercept = 0, slope = 1))**
- b + geom\_hline(aes(yintercept = lat))**
- b + geom\_vline(aes(xintercept = long))**
- b + geom\_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom\_spoke(aes(angle = 1:1155, radius = 1))**

### ONE VARIABLE continuous

- ```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```
- c + geom\_area(stat = "bin")** - x, y, alpha, color, fill, linetype, size
  - c + geom\_density(kernel = "gaussian")** - x, y, alpha, color, fill, group, linetype, size, weight
  - c + geom\_dotplot()** - x, y, alpha, color, fill
  - c + geom\_freqpoly()** - x, y, alpha, color, group, linetype, size
  - c + geom\_histogram(binwidth = 5)** - x, y, alpha, color, fill, linetype, size, weight
  - c2 + geom\_qq(aes(sample = hwy))** - x, y, alpha, color, fill, linetype, size, weight

### discrete

```
d <- ggplot(mpg, aes(fl))
```

- d + geom\_bar()** - x, alpha, color, fill, linetype, size, weight

### TWO VARIABLES both continuous

```
e <- ggplot(mpg, aes(cty, hwy))
```

- e + geom\_label(aes(label = cty), nudge\_x = 1, nudge\_y = 1)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust
- e + geom\_point()** - x, y, alpha, color, fill, shape, size, stroke
- e + geom\_quantile()** - x, y, alpha, color, group, linetype, size, weight
- e + geom\_rug(sides = "bl")** - x, y, alpha, color, linetype, size
- e + geom\_smooth(method = lm)** - x, y, alpha, color, fill, group, linetype, size, weight
- e + geom\_text(aes(label = cty), nudge\_x = 1, nudge\_y = 1)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

### one discrete, one continuous

```
f <- ggplot(mpg, aes(class, hwy))
```

- f + geom\_col()** - x, y, alpha, color, fill, group, linetype, size
- f + geom\_boxplot()** - x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight
- f + geom\_dotplot(binaxis = "y", stackdir = "center")** - x, y, alpha, color, fill, group
- f + geom\_violin(scale = "area")** - x, y, alpha, color, fill, group, linetype, size, weight

### both discrete

```
g <- ggplot(diamonds, aes(cut, color))
```

- g + geom\_count()** - x, y, alpha, color, fill, shape, size, stroke
- e + geom\_jitter(height = 2, width = 2)** - x, y, alpha, color, fill, shape, size

### THREE VARIABLES

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))
```

- l + geom\_contour(aes(z = z))** - x, y, z, alpha, color, group, linetype, size, weight
- l + geom\_contour\_filled(aes(fill = z))** - x, y, alpha, color, fill, group, linetype, size, subgroup
- l + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)** - x, y, alpha, fill
- l + geom\_tile(aes(fill = z))** - x, y, alpha, color, fill, linetype, size, width

### continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
```

- h + geom\_bin2d(binwidth = c(0.25, 500))** - x, y, alpha, color, fill, linetype, size, weight
- h + geom\_density\_2d()** - x, y, alpha, color, group, linetype, size
- h + geom\_hex()** - x, y, alpha, color, fill, size

### continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

- i + geom\_area()** - x, y, alpha, color, fill, linetype, size
- i + geom\_line()** - x, y, alpha, color, group, linetype, size
- i + geom\_step(direction = "hv")** - x, y, alpha, color, group, linetype, size

### visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

- j + geom\_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size
- j + geom\_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width  
Also **geom\_errorbarh()**.
- j + geom\_linerange()** - x, ymin, ymax, alpha, color, group, linetype, size
- j + geom\_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

### maps

```
data <- data.frame(murder = USArrests$Murder, state = tolower(rownames(USArrests)))
```

```
map <- map_data("state")
```

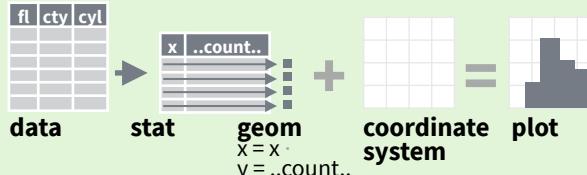
```
k <- ggplot(data, aes(fill = murder))
```

- k + geom\_map(aes(map\_id = state), map = map) + expand\_limits(x = map\$long, y = map\$lat)**  
map\_id, alpha, color, fill, linetype, size

# Stats

An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop).



Visualize a stat by changing the default stat of a geom function, `geom_bar(stat="count")` or by using a stat function, `stat_count(geom="bar")`, which calls a default geom to make a layer (equivalent to a geom function). Use `..name..` syntax to map stat variables to aesthetics.



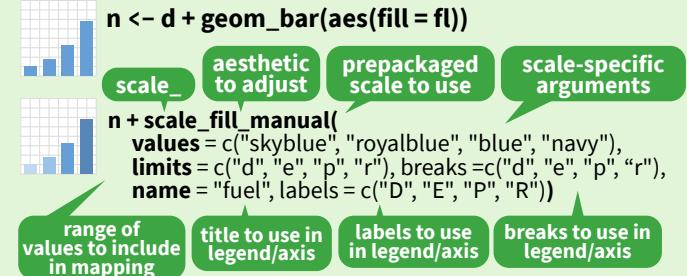
```

c + stat_bin(binwidth = 1, boundary = 10)
x, y | ..count.., ..ncount.., ..density.., ..ndensity..
c + stat_count(width = 1) x, y | ..count.., ..prop..
c + stat_density(adjust = 1, kernel = "gaussian")
x, y | ..count.., ..density.., ..scaled..
e + stat_bin_2d(bins = 30, drop = T)
x, y, fill | ..count.., ..density..
e + stat_bin_hex(bins = 30) x, y, fill | ..count.., ..density..
e + stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | ..level..
e + stat_ellipse(level = 0.95, segments = 51, type = "t")
l + stat_contour(aes(z = z)) x, y, z, order | ..level..
l + stat_summary_hex(aes(z = z), bins = 30, fun = max)
x, y, z, fill | ..value..
l + stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..
f + stat_boxplot(coef = 1.5)
x, y | ..lower.., ..middle.., ..upper.., ..width.., ..ymin.., ..ymax..
f + stat_ydensity(kernel = "gaussian", scale = "area") x, y
| ..density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..
e + stat_ecdf(n = 40) x, y | ..x.., ..y..
e + stat_quantile(quantiles = c(0.1, 0.9),
formula = y ~ log(x), method = "rq") x, y | ..quantile..
e + stat_smooth(method = "lm", formula = y ~ x, se = T,
level = 0.95) x, y | ..se.., ..x.., ..y.., ..ymin.., ..ymax..
ggplot() + xlim(-5, 5) + stat_function(fun = dnorm,
n = 20, geom = "point") x | ..x.., ..y..
ggplot() + stat_qq(aes(sample = 1:100))
x, y, sample | ..sample.., ..theoretical..
e + stat_sum() x, y, size | ..n.., ..prop..
e + stat_summary(fun.data = "mean_cl_boot")
h + stat_summary_bin(fun = "mean", geom = "bar")
e + stat_identity()
e + stat_unique()
  
```

# Scales

Override defaults with `scales` package.

**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



## GENERAL PURPOSE SCALES

Use with most aesthetics

`scale_*_continuous()` - Map cont' values to visual ones.  
`scale_*_discrete()` - Map discrete values to visual ones.  
`scale_*_binned()` - Map continuous values to discrete bins.  
`scale_*_identity()` - Use data values as visual ones.  
`scale_*_manual(values = c())` - Map discrete values to manually chosen visual ones.  
`scale_*_date(date_labels = "%m/%d")`,  
`date_breaks = "2 weeks"` - Treat data values as dates.  
`scale_*_datetime()` - Treat data values as date times.  
 Same as `scale_*_date()`. See `?strptime` for label formats.

## X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

`scale_x_log10()` - Plot x on log10 scale.  
`scale_x_reverse()` - Reverse the direction of the x axis.  
`scale_x_sqrt()` - Plot x on square root scale.

## COLOR AND FILL SCALES (DISCRETE)

`n + scale_fill_brewer(palette = "Blues")`  
 For palette choices:  
`RColorBrewer::display.brewer.all()`  
`n + scale_fill_grey(start = 0.2,`  
`end = 0.8, na.value = "red")`

## COLOR AND FILL SCALES (CONTINUOUS)

`o <- c + geom_dotplot(aes(fill = ..x..))`  
`o + scale_fill_distiller(palette = "Blues")`  
`o + scale_fill_gradient(low = "red", high = "yellow")`  
`o + scale_fill_gradient2(low = "red", high = "blue",`  
`mid = "white", midpoint = 25)`  
`o + scale_fill_gradientn(colors = topo.colors(6))`  
 Also: `rainbow()`, `heat.colors()`, `terrain.colors()`,  
`cm.colors()`, `RColorBrewer::brewer.pal()`

## SHAPE AND SIZE SCALES

`p <- e + geom_point(aes(shape = fl, size = cyl))`  
`p + scale_shape() + scale_size()`  
`p + scale_shape_manual(values = c(3:7))`  
`0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25`  
`□○△×+×◊▽★◆⊕⊖⊕⊖□○△○○□○△○○`  
`p + scale_radius(range = c(1,6))`  
`p + scale_size_area(max_size = 6)`

# Coordinate Systems

`r <- d + geom_bar()`

`r + coord_cartesian(xlim = c(0, 5))` - xlim, ylim  
 The default cartesian coordinate system.

`r + coord_fixed(ratio = 1/2)`  
 ratio, xlim, ylim - Cartesian coordinates with fixed aspect ratio between x and y units.

`ggplot(mpg, aes(y = fl)) + geom_bar()`  
 Flip cartesian coordinates by switching x and y aesthetic mappings.

`r + coord_polar(theta = "x", direction=1)`  
 theta, start, direction - Polar coordinates.

`r + coord_trans(y = "sqrt")` - x, y, xlim, ylim  
 Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

`π + coord_quickmap()`  
`π + coord_map(projection = "ortho", orientation = c(41, -74, 0))` - projection, xlim, ylim  
 Map projections from the mapproj package (mercator (default), azequalarea, lagrange, etc.).

## Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s <- ggplot(mpg, aes(fl, fill = drv))`

`s + geom_bar(position = "dodge")`  
 Arrange elements side by side.

`s + geom_bar(position = "fill")`  
 Stack elements on top of one another, normalize height.

`e + geom_point(position = "jitter")`  
 Add random noise to X and Y position of each element to avoid overplotting.

`e + geom_label(position = "nudge")`  
 Nudge labels away from points.

`s + geom_bar(position = "stack")`  
 Stack elements on top of one another.

Each position adjustment can be recast as a function with manual `width` and `height` arguments:  
`s + geom_bar(position = position_dodge(width = 1))`

## Themes

`r + theme_bw()`  
 White background with grid lines.

`r + theme_gray()`  
 Grey background (default theme).

`r + theme_dark()`  
 Dark for contrast.

`r + theme_classic()`  
`r + theme_light()`

`r + theme_linedraw()`  
`r + theme_minimal()`

`r + theme_void()`  
 Empty theme.

`r + theme()` Customize aspects of the theme such as axis, legend, panel, and facet properties.  
`r + ggtitle("Title") + theme(plot.title.position = "plot")`  
`r + theme(panel.background = element_rect(fill = "blue"))`

# Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t <- ggplot(mpg, aes(cty, hwy)) + geom_point()`

`t + facet_grid(cols = vars(fl))`  
 Facet into columns based on fl.

`t + facet_grid(rows = vars(year))`  
 Facet into rows based on year.

`t + facet_grid(rows = vars(year), cols = vars(fl))`  
 Facet into both rows and columns.

`t + facet_wrap(vars(fl))`  
 Wrap facets into a rectangular layout.

Set `scales` to let axis limits vary across facets.

`t + facet_grid(rows = vars(drv), cols = vars(fl), scales = "free")`

x and y axis limits adjust to individual facets:  
`"free_x"` - x axis limits adjust  
`"free_y"` - y axis limits adjust

Set `labeler` to adjust facet label:

`t + facet_grid(cols = vars(fl), labeler = label_both)`

`fl: c fl: d fl: e fl: p fl: r`

`t + facet_grid(rows = vars(fl), labeler = label_bquote(alpha ^ .(fl)))`

`αc αd αe αp αr`

## Labels and Legends

Use `labs()` to label the elements of your plot.

`t + labs(x = "New x axis label", y = "New y axis label", title = "Add a title above the plot", subtitle = "Add a subtitle below title", caption = "Add a caption below plot", alt = "Add alt text to the plot", <AES> = "New <AES> legend title")`

`t + annotate(geom = "text", x = 8, y = 9, label = "A")`  
 Places a geom with manually selected aesthetics.

`p + guides(x = guide_axis(n.dodge = 2))` Avoid crowded or overlapping labels with `guide_axis(n.dodge` or `angle`).

`n + guides(fill = "none")` Set legend type for each aesthetic: colorbar, legend, or none (no legend).

`n + theme(legend.position = "bottom")`  
 Place legend at "bottom", "top", "left", or "right".

`n + scale_fill_discrete(name = "Title", labels = c("A", "B", "C", "D", "E"))`  
 Set legend title and labels with a scale function.

## Zooming

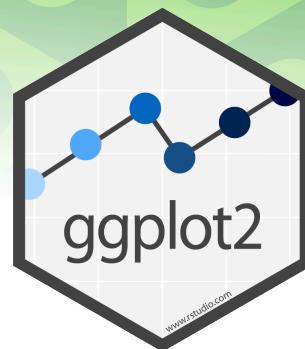
Without clipping (preferred):

`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`

With clipping (removes unseen data points):

`t + xlim(0, 100) + ylim(10, 20)`

`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`



# Data transformation with dplyr :: CHEAT SHEET



dplyr functions work with pipes and expect **tidy data**. In tidy data:



Each **variable** is in its own **column**



Each **observation**, or **case**, is in its own **row**



`x %>% f(y)` becomes `f(x, y)`

## Summarise Cases

Apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



`summarise(.data, ...)`  
Compute table of summaries.  
`summarise(mtcars, avg = mean(mpg))`

`count(.data, ..., wt = NULL, sort = FALSE, name = NULL)` Count number of rows in each group defined by the variables in ... Also **tally()**.  
`count(mtcars, cyl)`

## Group Cases

Use **group\_by(.data, ..., .add = FALSE, .drop = TRUE)** to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.

`mtcars %>%  
group_by(cyl) %>%  
summarise(avg = mean(mpg))`

Use **rowwise(.data, ...)** to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyverse cheat sheet for list-column workflow.

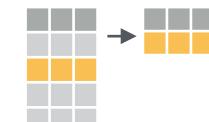
`starwars %>%  
rowwise() %>%  
mutate(film_count = length(films))`

**ungroup(x, ...)** Returns ungrouped copy of table.  
`ungroup(g_mtcars)`

## Manipulate Cases

### EXTRACT CASES

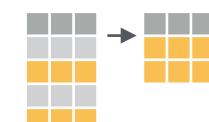
Row functions return a subset of rows as a new table.



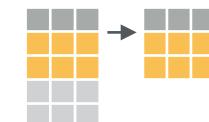
**filter(.data, ..., .preserve = FALSE)** Extract rows that meet logical criteria.  
`filter(mtcars, mpg > 20)`



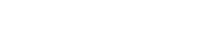
**distinct(.data, ..., .keep\_all = FALSE)** Remove rows with duplicate values.  
`distinct(mtcars, gear)`



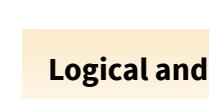
**slice(.data, ..., .preserve = FALSE)** Select rows by position.  
`slice(mtcars, 10:15)`



**slice\_sample(.data, ..., n, prop, weight\_by = NULL, replace = FALSE)** Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows.  
`slice_sample(mtcars, n = 5, replace = TRUE)`



**slice\_min(.data, order\_by, ..., n, prop, with\_ties = TRUE)** and **slice\_max()** Select rows with the lowest and highest values.  
`slice_min(mtcars, mpg, prop = 0.25)`



**slice\_head(.data, ..., n, prop)** and **slice\_tail()** Select the first or last rows.  
`slice_head(mtcars, n = 5)`

### Logical and boolean operators to use with filter()

|                 |                   |                    |                       |                   |                    |                    |
|-----------------|-------------------|--------------------|-----------------------|-------------------|--------------------|--------------------|
| <code>==</code> | <code>&lt;</code> | <code>&lt;=</code> | <code>is.na()</code>  | <code>%in%</code> | <code> </code>     | <code>xor()</code> |
| <code>!=</code> | <code>&gt;</code> | <code>&gt;=</code> | <code>!is.na()</code> | <code>!</code>    | <code>&amp;</code> |                    |

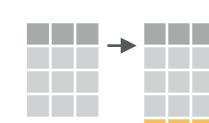
See **?base::Logic** and **?Comparison** for help.

### ARRANGE CASES



**arrange(.data, ..., .by\_group = FALSE)** Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.  
`arrange(mtcars, mpg)`  
`arrange(mtcars, desc(mpg))`

### ADD CASES



**add\_row(.data, ..., .before = NULL, .after = NULL)** Add one or more rows to a table.  
`add_row(cars, speed = 1, dist = 1)`

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.



**pull(.data, var = -1, name = NULL, ...)** Extract column values as a vector, by name or index.  
`pull(mtcars, wt)`



**select(.data, ...)** Extract columns as a table.  
`select(mtcars, mpg, wt)`



**relocate(.data, ..., .before = NULL, .after = NULL)** Move columns to new position.  
`relocate(mtcars, mpg, cyl, .after = last_col())`

### Use these helpers with select() and across()

e.g. `select(mtcars, mpg:cyl)`

`contains(match)`

`ends_with(match)`

`starts_with(match)`

`num_range(prefix, range)`

`all_of(x)/any_of(x, ..., vars)`

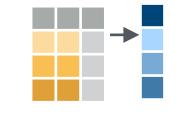
- e.g. `-gear`

`matches(match)` **everything()**

### MANIPULATE MULTIPLE VARIABLES AT ONCE



**across(.cols, .funs, ..., .names = NULL)** Summarise or mutate multiple columns in the same way.  
`summarise(mtcars, across(everything(), mean))`



**c\_across(.cols)** Compute across columns in row-wise data.  
`transmute(rowwise(UKgas), total = sum(c_across(1:2)))`

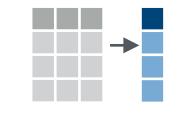
### MAKE NEW VARIABLES

Apply **vectorized functions** to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).



**vectorized function**

**mutate(.data, ..., .keep = "all", .before = NULL, .after = NULL)** Compute new column(s). Also **add\_column()**, **add\_count()**, and **add\_tally()**.  
`mutate(mtcars, gpm = 1 / mpg)`



**transmute(.data, ...)** Compute new column(s), drop others.  
`transmute(mtcars, gpm = 1 / mpg)`



**rename(.data, ...)** Rename columns. Use **rename\_with()** to rename with a function.  
`rename(cars, distance = dist)`



# Vectorized Functions

## TO USE WITH MUTATE ()

**mutate()** and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function →

## OFFSET

dplyr::lag() - offset elements by 1  
dplyr::lead() - offset elements by -1

## CUMULATIVE AGGREGATE

dplyr::cumall() - cumulative all()  
dplyr::cumany() - cumulative any()  
  **cummax()** - cumulative max()  
dplyr::cummean() - cumulative mean()  
  **cummin()** - cumulative min()  
  **cumprod()** - cumulative prod()  
  **cumsum()** - cumulative sum()

## RANKING

dplyr::cume\_dist() - proportion of all values <=  
dplyr::dense\_rank() - rank w ties = min, no gaps  
dplyr::min\_rank() - rank with ties = min  
dplyr::ntile() - bins into n bins  
dplyr::percent\_rank() - min\_rank scaled to [0,1]  
dplyr::row\_number() - rank with ties = "first"

## MATH

+, -, \*, /, ^, %/%, %% - arithmetic ops  
log(), log2(), log10() - logs  
<, <=, >, >=, !=, == - logical comparisons  
dplyr::between() - x >= left & x <= right  
dplyr::near() - safe == for floating point numbers

## MISCELLANEOUS

dplyr::case\_when() - multi-case if\_else()  
starwars %>%  
  mutate(type = case\_when(  
    height > 200 | mass > 200 ~ "large",  
    species == "Droid" ~ "robot",  
    TRUE ~ "other")  
  )  
dplyr::coalesce() - first non-NA values by element across a set of vectors  
dplyr::if\_else() - element-wise if() + else()  
dplyr::na\_if() - replace specific values with NA  
  pmax() - element-wise max()  
  pmin() - element-wise min()

# Summary Functions

## TO USE WITH SUMMARISE ()

**summarise()** applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function →

## COUNT

dplyr::n() - number of values/rows  
dplyr::n\_distinct() - # of uniques  
  **sum(!is.na())** - # of non-NAs

## POSITION

**mean()** - mean, also **mean(!is.na())**  
**median()** - median

## LOGICAL

**mean()** - proportion of TRUE's  
**sum()** - # of TRUE's

## ORDER

dplyr::first() - first value  
dplyr::last() - last value  
dplyr::nth() - value in nth location of vector

## RANK

**quantile()** - nth quantile  
**min()** - minimum value  
**max()** - maximum value

## SPREAD

**IQR()** - Inter-Quartile Range  
**mad()** - median absolute deviation  
**sd()** - standard deviation  
**var()** - variance

# Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

A B C → C A B  
1 a t → 1 a t  
2 b u → 2 b u  
3 c v → 3 c v  
tibble::rownames\_to\_column()  
Move row names into col.  
a <- rownames\_to\_column(mtcars,  
var = "C")

A B C → A B  
1 a t → 1 a  
2 b u → 2 b  
3 c v → 3 c  
tibble::column\_to\_rownames()  
Move col into row names.  
column\_to\_rownames(a, var = "C")

Also **tibble::has\_rownames()** and **tibble::remove\_rownames()**.

# Combine Tables

## COMBINE VARIABLES

|       |       |
|-------|-------|
| X     | y     |
| A B C | E F G |
| a t 1 | a t 3 |
| b u 2 | b u 2 |
| c v 3 | d w 1 |

**bind\_cols(..., .name\_repair)** Returns tables placed side by side as a single table. Column lengths must be equal. Columns will NOT be matched by id (to do that look at Relational Data below), so be sure to check that both tables are ordered the way you want before binding.

## RELATIONAL DATA

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

|                                |   |
|--------------------------------|---|
| A B C D                        | <b>left_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")</b> Join matching values from y to x. |
| a t 1 3<br>b u 2 2<br>c v 3 NA |   |

|                                |  |
|--------------------------------|--|
| A B C D                        | <b>right_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")</b> Join matching values from x to y. |
| a t 1 3<br>b u 2 2<br>d w NA 1 |  |

|                    |  |
|--------------------|--|
| A B C D            | <b>inner_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")</b> Join data. Retain only rows with matches. |
| a t 1 3<br>b u 2 2 |  |

|  |   |
|--|---|
| A B C D                                    | <b>full_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")</b> Join data. Retain all values, all rows. |
| a t 1 3<br>b u 2 2<br>c v 3 NA<br>d w NA 1 |   |

## COLUMN MATCHING FOR JOINS

|               |
|---------------|
| A B x C B y D |
| a t 1 t 3     |
| b u 2 u 2     |
| c v 3 NA NA   |

Use **by = c("col1", "col2", ...)** to specify one or more common columns to match on.  
left\_join(x, y, by = "A")

|                   |
|-------------------|
| A x B x C A y B y |
| a t 1 d w         |
| b u 2 b u         |
| c v 3 a t         |

Use a named vector, **by = c("col1" = "col2")**, to match on columns that have different names in each table.  
left\_join(x, y, by = c("C" = "D"))

|               |
|---------------|
| A1 B1 C A2 B2 |
| a t 1 d w     |
| b u 2 b u     |
| c v 3 a t     |

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.  
left\_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

## COMBINE CASES

|       |       |
|-------|-------|
| X     | y     |
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | d w 4 |

**bind\_rows(..., id = NULL)**  
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured).

|       |         |
|-------|---------|
| X     | y       |
| A B C | A B C   |
| a t 1 | a t 1   |
| b u 2 | b u 2   |
| c v 3 | y d w 4 |

Use a "Filtering Join" to filter one table against the rows of another.

|       |       |
|-------|-------|
| X     | y     |
| A B C | A B D |
| a t 1 | a t 3 |
| b u 2 | b u 2 |
| c v 3 | d w 1 |

**semi\_join(x, y, by = NULL, copy = FALSE, ..., na\_matches = "na")** Return rows of x that have a match in y. Use to see what will be included in a join.

|       |
|-------|
| A B C |
| c v 3 |

**anti\_join(x, y, by = NULL, copy = FALSE, ..., na\_matches = "na")** Return rows of x that do not have a match in y. Use to see what will not be included in a join.

Use a "Nest Join" to inner join one table to another into a nested data frame.

|       |                |
|-------|----------------|
| A B C | y              |
| a t 1 | <tibble [1x2]> |
| b u 2 | <tibble [1x2]> |
| c v 3 | <tibble [1x2]> |

**nest\_join(x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, ...)** Join data, nesting matches from y in a single new data frame column.

## SET OPERATIONS

|       |
|-------|
| A B C |
| c v 3 |

**intersect(x, y, ...)**



Rows that appear in both x and y.

|       |
|-------|
| A B C |
| a t 1 |
| b u 2 |

**setdiff(x, y, ...)**



Rows that appear in x but not y.



## QUERYING DATA FROM A TABLE

**SELECT c1, c2 FROM t;**

Query data in columns c1, c2 from a table

**SELECT \* FROM t;**

Query all rows and columns from a table

**SELECT c1, c2 FROM t**

**WHERE condition;**

Query data and filter rows with a condition

**SELECT DISTINCT c1 FROM t**

**WHERE condition;**

Query distinct rows from a table

**SELECT c1, c2 FROM t**

**ORDER BY c1 ASC [DESC];**

Sort the result set in ascending or descending order

**SELECT c1, c2 FROM t**

**ORDER BY c1**

**LIMIT n OFFSET offset;**

Skip offset of rows and return the next n rows

**SELECT c1, aggregate(c2)**

**FROM t**

**GROUP BY c1;**

Group rows using an aggregate function

**SELECT c1, aggregate(c2)**

**FROM t**

**GROUP BY c1**

**HAVING condition;**

Filter groups using HAVING clause

## QUERYING FROM MULTIPLE TABLES

**SELECT c1, c2**

**FROM t1**

**INNER JOIN t2 ON condition;**

Inner join t1 and t2

**SELECT c1, c2**

**FROM t1**

**LEFT JOIN t2 ON condition;**

Left join t1 and t2

**SELECT c1, c2**

**FROM t1**

**RIGHT JOIN t2 ON condition;**

Right join t1 and t2

**SELECT c1, c2**

**FROM t1**

**FULL OUTER JOIN t2 ON condition;**

Perform full outer join

**SELECT c1, c2**

**FROM t1**

**CROSS JOIN t2;**

Produce a Cartesian product of rows in tables

**SELECT c1, c2**

**FROM t1, t2;**

Another way to perform cross join

**SELECT c1, c2**

**FROM t1 A**

**INNER JOIN t2 B ON condition;**

Join t1 to itself using INNER JOIN clause

## USING SQL OPERATORS

**SELECT c1, c2 FROM t1**

**UNION [ALL]**

**SELECT c1, c2 FROM t2;**

Combine rows from two queries

**SELECT c1, c2 FROM t1**

**INTERSECT**

**SELECT c1, c2 FROM t2;**

Return the intersection of two queries

**SELECT c1, c2 FROM t1**

**MINUS**

**SELECT c1, c2 FROM t2;**

Subtract a result set from another result set

**SELECT c1, c2 FROM t1**

**WHERE c1 [NOT] LIKE pattern;**

Query rows using pattern matching %, \_

**SELECT c1, c2 FROM t**

**WHERE c1 [NOT] IN value\_list;**

Query rows in a list

**SELECT c1, c2 FROM t**

**WHERE c1 BETWEEN low AND high;**

Query rows between two values

**SELECT c1, c2 FROM t**

**WHERE c1 IS [NOT] NULL;**

Check if values in a table is NULL or not



## MANAGING TABLES

```
CREATE TABLE t (
    id INT PRIMARY KEY,
    name VARCHAR NOT NULL,
    price INT DEFAULT 0
);
```

Create a new table with three columns

```
DROP TABLE t;
```

Delete the table from the database

```
ALTER TABLE t ADD column;
```

Add a new column to the table

```
ALTER TABLE t DROP COLUMN c;
```

Drop column c from the table

```
ALTER TABLE t ADD constraint;
```

Add a constraint

```
ALTER TABLE t DROP constraint;
```

Drop a constraint

```
ALTER TABLE t1 RENAME TO t2;
```

Rename a table from t1 to t2

```
ALTER TABLE t1 RENAME c1 TO c2;
```

Rename column c1 to c2

```
TRUNCATE TABLE t;
```

Remove all data in a table

## USING SQL CONSTRAINTS

```
CREATE TABLE t(
    c1 INT, c2 INT, c3 VARCHAR,
    PRIMARY KEY (c1,c2)
);
```

Set c1 and c2 as a primary key

```
CREATE TABLE t1(
    c1 INT PRIMARY KEY,
    c2 INT,
    FOREIGN KEY (c2) REFERENCES t2(c2)
);
```

Set c2 column as a foreign key

```
CREATE TABLE t(
    c1 INT, c2 INT,
    UNIQUE(c2,c3)
);
```

Make the values in c1 and c2 unique

```
CREATE TABLE t(
    c1 INT, c2 INT,
    CHECK(c1 > 0 AND c1 >= c2)
);
```

Ensure c1 > 0 and values in c1 >= c2

```
CREATE TABLE t(
    c1 INT PRIMARY KEY,
    c2 VARCHAR NOT NULL
);
```

Set values in c2 column not NULL

## MODIFYING DATA

```
INSERT INTO t(column_list)
VALUES(value_list);
```

Insert one row into a table

```
INSERT INTO t(column_list)
VALUES (value_list),
       (value_list), ....;
```

Insert multiple rows into a table

```
INSERT INTO t1(column_list)
SELECT column_list
FROM t2;
```

Insert rows from t2 into t1

```
UPDATE t
SET c1 = new_value;
```

Update new value in the column c1 for all rows

```
UPDATE t
SET c1 = new_value,
    c2 = new_value
WHERE condition;
```

Update values in the column c1, c2 that match the condition

```
DELETE FROM t;
```

Delete all data in a table

```
DELETE FROM t
WHERE condition;
```

Delete subset of rows in a table



## MANAGING VIEWS

**CREATE VIEW** v(c1,c2)

**AS**

**SELECT** c1, c2

**FROM** t;

Create a new view that consists of c1 and c2

**CREATE VIEW** v(c1,c2)

**AS**

**SELECT** c1, c2

**FROM** t;

**WITH [CASCADED | LOCAL] CHECK OPTION;**

Create a new view with check option

**CREATE RECURSIVE VIEW** v

**AS**

select-statement -- *anchor part*

**UNION [ALL]**

select-statement; -- *recursive part*

Create a recursive view

**CREATE TEMPORARY VIEW** v

**AS**

**SELECT** c1, c2

**FROM** t;

Create a temporary view

**DROP VIEW** view\_name;

Delete a view

## MANAGING INDEXES

**CREATE INDEX** idx\_name

**ON** t(c1,c2);

Create an index on c1 and c2 of the table t

**CREATE UNIQUE INDEX** idx\_name

**ON** t(c3,c4);

Create a unique index on c3, c4 of the table t

**DROP INDEX** idx\_name;

Drop an index

## SQL AGGREGATE FUNCTIONS

**AVG** returns the average of a list

**COUNT** returns the number of elements of a list

**SUM** returns the total of a list

**MAX** returns the maximum value in a list

**MIN** returns the minimum value in a list

## MANAGING TRIGGERS

**CREATE OR MODIFY TRIGGER** trigger\_name

**WHEN EVENT**

**ON** table\_name **TRIGGER\_TYPE**

**EXECUTE stored\_procedure;**

Create or modify a trigger

### WHEN

- **BEFORE** – invoke before the event occurs
- **AFTER** – invoke after the event occurs

### EVENT

- **INSERT** – invoke for INSERT
- **UPDATE** – invoke for UPDATE
- **DELETE** – invoke for DELETE

### TRIGGER\_TYPE

- **FOR EACH ROW**
- **FOR EACH STATEMENT**

**CREATE TRIGGER** before\_insert\_person

**BEFORE INSERT**

**ON** person **FOR EACH ROW**

**EXECUTE stored\_procedure;**

Create a trigger invoked before a new row is inserted into the person table

**DROP TRIGGER** trigger\_name;

Delete a specific trigger

# Super VIP Cheatsheet: Machine Learning

Afshine AMIDI and Shervine AMIDI

October 6, 2018

## Contents

### 1 Supervised Learning

**2**

|       |  |   |
|-------|--|---|
| 1.1   | Introduction to Supervised Learning . . . . .    | 2 |
| 1.2   | Notations and general concepts . . . . .         | 2 |
| 1.3   | Linear models . . . . .                          | 2 |
| 1.3.1 | Linear regression . . . . .                      | 2 |
| 1.3.2 | Classification and logistic regression . . . . . | 3 |
| 1.3.3 | Generalized Linear Models . . . . .              | 3 |
| 1.4   | Support Vector Machines . . . . .                | 3 |
| 1.5   | Generative Learning . . . . .                    | 4 |
| 1.5.1 | Gaussian Discriminant Analysis . . . . .         | 4 |
| 1.5.2 | Naive Bayes . . . . .                            | 4 |
| 1.6   | Tree-based and ensemble methods . . . . .        | 4 |
| 1.7   | Other non-parametric approaches . . . . .        | 4 |
| 1.8   | Learning Theory . . . . .                        | 5 |

### 2 Unsupervised Learning

**6**

|       |   |   |
|-------|---|---|
| 2.1   | Introduction to Unsupervised Learning . . . . . | 6 |
| 2.2   | Clustering . . . . .                            | 6 |
| 2.2.1 | Expectation-Maximization . . . . .              | 6 |
| 2.2.2 | $k$ -means clustering . . . . .                 | 6 |
| 2.2.3 | Hierarchical clustering . . . . .               | 6 |
| 2.2.4 | Clustering assessment metrics . . . . .         | 6 |
| 2.3   | Dimension reduction . . . . .                   | 7 |
| 2.3.1 | Principal component analysis . . . . .          | 7 |
| 2.3.2 | Independent component analysis . . . . .        | 7 |

### 3 Deep Learning

**8**

|     |  |   |
|-----|--|---|
| 3.1 | Neural Networks . . . . .                    | 8 |
| 3.2 | Convolutional Neural Networks . . . . .      | 8 |
| 3.3 | Recurrent Neural Networks . . . . .          | 8 |
| 3.4 | Reinforcement Learning and Control . . . . . | 9 |

|          |   |           |
|----------|---|-----------|
| <b>4</b> | <b>Machine Learning Tips and Tricks</b>                 | <b>10</b> |
| 4.1      | Metrics . . . . .                                       | 10        |
| 4.1.1    | Classification . . . . .                                | 10        |
| 4.1.2    | Regression . . . . .                                    | 10        |
| 4.2      | Model selection . . . . .                               | 11        |
| 4.3      | Diagnostics . . . . .                                   | 11        |
| <b>5</b> | <b>Refreshers</b>                                       | <b>12</b> |
| 5.1      | Probabilities and Statistics . . . . .                  | 12        |
| 5.1.1    | Introduction to Probability and Combinatorics . . . . . | 12        |
| 5.1.2    | Conditional Probability . . . . .                       | 12        |
| 5.1.3    | Random Variables . . . . .                              | 13        |
| 5.1.4    | Jointly Distributed Random Variables . . . . .          | 13        |
| 5.1.5    | Parameter estimation . . . . .                          | 14        |
| 5.2      | Linear Algebra and Calculus . . . . .                   | 14        |
| 5.2.1    | General notations . . . . .                             | 14        |
| 5.2.2    | Matrix operations . . . . .                             | 15        |
| 5.2.3    | Matrix properties . . . . .                             | 15        |
| 5.2.4    | Matrix calculus . . . . .                               | 16        |

## 1 Supervised Learning

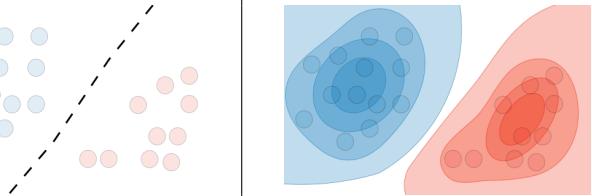
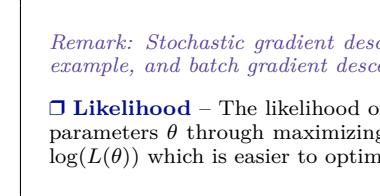
### 1.1 Introduction to Supervised Learning

Given a set of data points  $\{x^{(1)}, \dots, x^{(m)}\}$  associated to a set of outcomes  $\{y^{(1)}, \dots, y^{(m)}\}$ , we want to build a classifier that learns how to predict  $y$  from  $x$ .

**Type of prediction** – The different types of predictive models are summed up in the table below:

|          | Regression        | Classifier                            |
|----------|-------------------|---------------------------------------|
| Outcome  | Continuous        | Class                                 |
| Examples | Linear regression | Logistic regression, SVM, Naive Bayes |

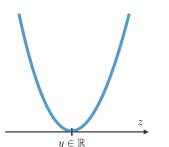
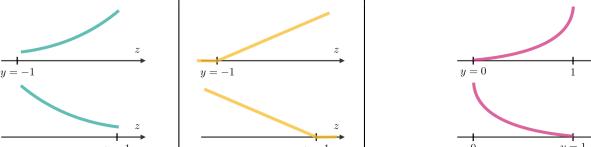
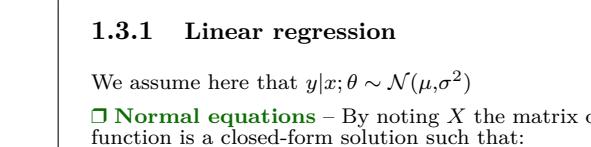
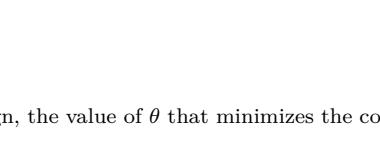
**Type of model** – The different models are summed up in the table below:

|                | Discriminative model   | Generative model  |
|----------------|--|---|
| Goal           | Directly estimate $P(y x)$   | Estimate $P(x y)$ to deduce $P(y x)$  |
| What's learned | Decision boundary  | Probability distributions of the data   |
| Illustration   |  |  |
| Examples       | Regressions, SVMs  | GDA, Naive Bayes  |

### 1.2 Notations and general concepts

**Hypothesis** – The hypothesis is noted  $h_\theta$  and is the model that we choose. For a given input data  $x^{(i)}$ , the model prediction output is  $h_\theta(x^{(i)})$ .

**Loss function** – A loss function is a function  $L : (z,y) \in \mathbb{R} \times Y \mapsto L(z,y) \in \mathbb{R}$  that takes as inputs the predicted value  $z$  corresponding to the real data value  $y$  and outputs how different they are. The common loss functions are summed up in the table below:

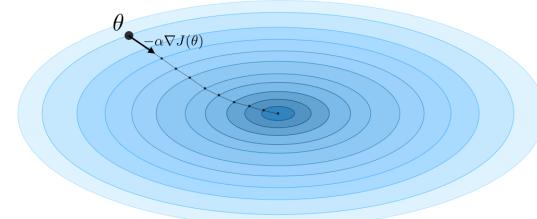
| Least squared   | Logistic  | Hinge  | Cross-entropy   |
|---|---|--|---|
| $\frac{1}{2}(y - z)^2$  | $\log(1 + \exp(-yz))$   | $\max(0, 1 - yz)$  | $-\left[y \log(z) + (1 - y) \log(1 - z)\right]$                                       |
|  |  |  |  |
| Linear regression   | Logistic regression   | SVM  | Neural Network  |

**Cost function** – The cost function  $J$  is commonly used to assess the performance of a model, and is defined with the loss function  $L$  as follows:

$$J(\theta) = \sum_{i=1}^m L(h_\theta(x^{(i)}), y^{(i)})$$

**Gradient descent** – By noting  $\alpha \in \mathbb{R}$  the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function  $J$  as follows:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$



Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.

**Likelihood** – The likelihood of a model  $L(\theta)$  given parameters  $\theta$  is used to find the optimal parameters  $\theta$  through maximizing the likelihood. In practice, we use the log-likelihood  $\ell(\theta) = \log(L(\theta))$  which is easier to optimize. We have:

$$\theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$

**Newton's algorithm** – The Newton's algorithm is a numerical method that finds  $\theta$  such that  $\ell'(\theta) = 0$ . Its update rule is as follows:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule:

$$\theta \leftarrow \theta - (\nabla_{\theta}^2 \ell(\theta))^{-1} \nabla_{\theta} \ell(\theta)$$

### 1.3 Linear models

#### 1.3.1 Linear regression

We assume here that  $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$

**Normal equations** – By noting  $X$  the matrix design, the value of  $\theta$  that minimizes the cost function is a closed-form solution such that:

$$\theta = (X^T X)^{-1} X^T y$$

□ **LMS algorithm** – By noting  $\alpha$  the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a training set of  $m$  data points, which is also known as the Widrow-Hoff learning rule, is as follows:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)}$$

*Remark: the update rule is a particular case of the gradient ascent.*

□ **LWR** – Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by  $w^{(i)}(x)$ , which is defined with parameter  $\tau \in \mathbb{R}$  as:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

### 1.3.2 Classification and logistic regression

□ **Sigmoid function** – The sigmoid function  $g$ , also known as the logistic function, is defined as follows:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in ]0,1[$$

□ **Logistic regression** – We assume here that  $y|x; \theta \sim \text{Bernoulli}(\phi)$ . We have the following form:

$$\phi = p(y=1|x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

*Remark: there is no closed form solution for the case of logistic regressions.*

□ **Softmax regression** – A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes. By convention, we set  $\theta_K = 0$ , which makes the Bernoulli parameter  $\phi_i$  of each class  $i$  equal to:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

### 1.3.3 Generalized Linear Models

□ **Exponential family** – A class of distributions is said to be in the exponential family if it can be written in terms of a natural parameter, also called the canonical parameter or link function,  $\eta$ , a sufficient statistic  $T(y)$  and a log-partition function  $a(\eta)$  as follows:

$$p(y; \eta) = b(y) \exp(\eta T(y) - a(\eta))$$

*Remark: we will often have  $T(y) = y$ . Also,  $\exp(-a(\eta))$  can be seen as a normalization parameter that will make sure that the probabilities sum to one.*

Here are the most common exponential distributions summed up in the following table:

| Distribution | $\eta$                                 | $T(y)$ | $a(\eta)$                                  | $b(y)$  |
|--------------|--|--------|--|---|
| Bernoulli    | $\log\left(\frac{\phi}{1-\phi}\right)$ | $y$    | $\log(1 + \exp(\eta))$                     | 1   |
| Gaussian     | $\mu$                                  | $y$    | $\frac{\eta^2}{2}$                         | $\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$ |
| Poisson      | $\log(\lambda)$                        | $y$    | $e^\eta$                                   | $\frac{1}{y!}$  |
| Geometric    | $\log(1 - \phi)$                       | $y$    | $\log\left(\frac{e^\eta}{1-e^\eta}\right)$ | 1   |

□ **Assumptions of GLMs** – Generalized Linear Models (GLM) aim at predicting a random variable  $y$  as a function of  $x \in \mathbb{R}^{n+1}$  and rely on the following 3 assumptions:

$$(1) \quad y|x; \theta \sim \text{ExpFamily}(\eta) \quad (2) \quad h_\theta(x) = E[y|x; \theta] \quad (3) \quad \eta = \theta^T x$$

*Remark: ordinary least squares and logistic regression are special cases of generalized linear models.*

### 1.4 Support Vector Machines

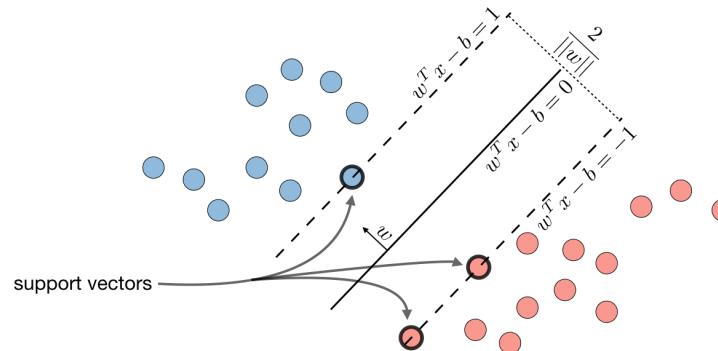
The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

□ **Optimal margin classifier** – The optimal margin classifier  $h$  is such that:

$$h(x) = \text{sign}(w^T x - b)$$

where  $(w, b) \in \mathbb{R}^n \times \mathbb{R}$  is the solution of the following optimization problem:

$$\min \frac{1}{2} \|w\|^2 \quad \text{such that} \quad y^{(i)}(w^T x^{(i)} - b) \geq 1$$



*Remark: the line is defined as  $w^T x - b = 0$ .*

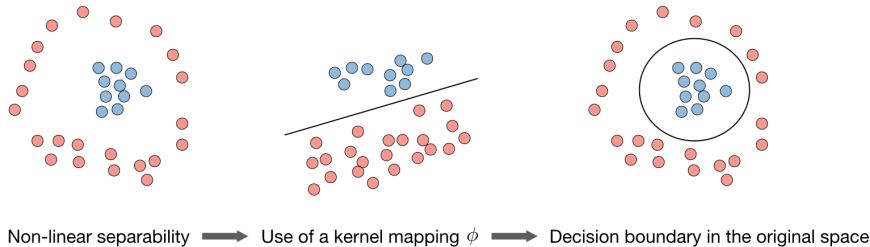
□ **Hinge loss** – The hinge loss is used in the setting of SVMs and is defined as follows:

$$L(z, y) = [1 - yz]_+ = \max(0, 1 - yz)$$

□ **Kernel** – Given a feature mapping  $\phi$ , we define the kernel  $K$  to be defined as:

$$K(x,z) = \phi(x)^T \phi(z)$$

In practice, the kernel  $K$  defined by  $K(x,z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$  is called the Gaussian kernel and is commonly used.



*Remark: we say that we use the "kernel trick" to compute the cost function using the kernel because we actually don't need to know the explicit mapping  $\phi$ , which is often very complicated. Instead, only the values  $K(x,z)$  are needed.*

□ **Lagrangian** – We define the Lagrangian  $\mathcal{L}(w,b)$  as follows:

$$\mathcal{L}(w,b) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

*Remark: the coefficients  $\beta_i$  are called the Lagrange multipliers.*

## 1.5 Generative Learning

A generative model first tries to learn how the data is generated by estimating  $P(x|y)$ , which we can then use to estimate  $P(y|x)$  by using Bayes' rule.

### 1.5.1 Gaussian Discriminant Analysis

□ **Setting** – The Gaussian Discriminant Analysis assumes that  $y$  and  $x|y = 0$  and  $x|y = 1$  are such that:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{and} \quad x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$

□ **Estimation** – The following table sums up the estimates that we find when maximizing the likelihood:

| $\hat{\phi}$  | $\hat{\mu}_j \quad (j = 0,1)$   | $\hat{\Sigma}$  |
|---|---|---|
| $\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{y^{(i)}=1\}}$ | $\frac{\sum_{i=1}^m \mathbf{1}_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m \mathbf{1}_{\{y^{(i)}=j\}}}$ | $\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$ |

### 1.5.2 Naive Bayes

□ **Assumption** – The Naive Bayes model supposes that the features of each data point are all independent:

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^n P(x_i|y)$$

□ **Solutions** – Maximizing the log-likelihood gives the following solutions, with  $k \in \{0,1\}$ ,  $l \in [1, L]$

$$P(y = k) = \frac{1}{m} \times \#\{j | y^{(j)} = k\}$$

and

$$P(x_i = l | y = k) = \frac{\#\{j | y^{(j)} = k \text{ and } x_i^{(j)} = l\}}{\#\{j | y^{(j)} = k\}}$$

*Remark: Naive Bayes is widely used for text classification and spam detection.*

## 1.6 Tree-based and ensemble methods

These methods can be used for both regression and classification problems.

□ **CART** – Classification and Regression Trees (CART), commonly known as decision trees, can be represented as binary trees. They have the advantage to be very interpretable.

□ **Random forest** – It is a tree-based technique that uses a high number of decision trees built out of randomly selected sets of features. Contrary to the simple decision tree, it is highly uninterpretable but its generally good performance makes it a popular algorithm.

*Remark: random forests are a type of ensemble methods.*

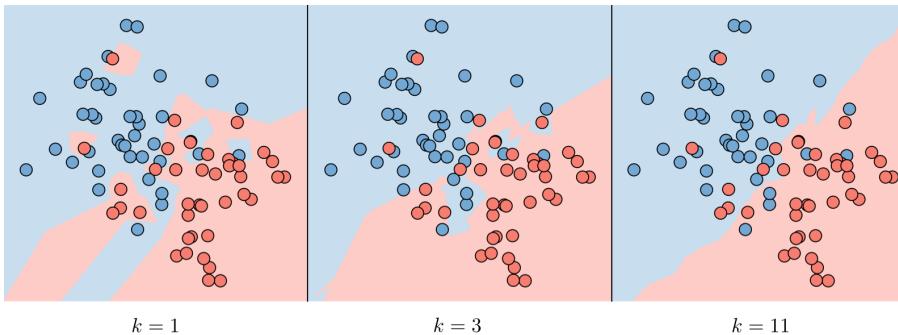
□ **Boosting** – The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

| Adaptive boosting  | Gradient boosting                           |
|--|---|
| - High weights are put on errors to improve at the next boosting step<br>- Known as Adaboost | - Weak learners trained on remaining errors |

## 1.7 Other non-parametric approaches

□  **$k$ -nearest neighbors** – The  $k$ -nearest neighbors algorithm, commonly known as  $k$ -NN, is a non-parametric approach where the response of a data point is determined by the nature of its  $k$  neighbors from the training set. It can be used in both classification and regression settings.

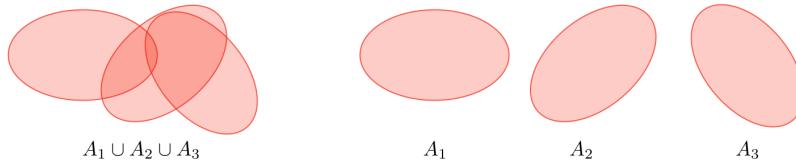
*Remark: The higher the parameter  $k$ , the higher the bias, and the lower the parameter  $k$ , the higher the variance.*



## 1.8 Learning Theory

**□ Union bound** – Let  $A_1, \dots, A_k$  be  $k$  events. We have:

$$P(A_1 \cup \dots \cup A_k) \leq P(A_1) + \dots + P(A_k)$$



**□ Hoeffding inequality** – Let  $Z_1, \dots, Z_m$  be  $m$  iid variables drawn from a Bernoulli distribution of parameter  $\phi$ . Let  $\hat{\phi}$  be their sample mean and  $\gamma > 0$  fixed. We have:

$$P(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

*Remark: this inequality is also known as the Chernoff bound.*

**□ Training error** – For a given classifier  $h$ , we define the training error  $\hat{\epsilon}(h)$ , also known as the empirical risk or empirical error, to be as follows:

$$\hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

**□ Probably Approximately Correct (PAC)** – PAC is a framework under which numerous results on learning theory were proved, and has the following set of assumptions:

- the training and testing sets follow the same distribution
- the training examples are drawn independently

**□ Shattering** – Given a set  $S = \{x^{(1)}, \dots, x^{(d)}\}$ , and a set of classifiers  $\mathcal{H}$ , we say that  $\mathcal{H}$  shatters  $S$  if for any set of labels  $\{y^{(1)}, \dots, y^{(d)}\}$ , we have:

$$\exists h \in \mathcal{H}, \quad \forall i \in \llbracket 1, d \rrbracket, \quad h(x^{(i)}) = y^{(i)}$$

**□ Upper bound theorem** – Let  $\mathcal{H}$  be a finite hypothesis class such that  $|\mathcal{H}| = k$  and let  $\delta$  and the sample size  $m$  be fixed. Then, with probability of at least  $1 - \delta$ , we have:

$$\epsilon(\hat{h}) \leq \left( \min_{h \in \mathcal{H}} \epsilon(h) \right) + 2 \sqrt{\frac{1}{2m} \log \left( \frac{2k}{\delta} \right)}$$

**□ VC dimension** – The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class  $\mathcal{H}$ , noted  $\text{VC}(\mathcal{H})$  is the size of the largest set that is shattered by  $\mathcal{H}$ .

*Remark: the VC dimension of  $\mathcal{H} = \{\text{set of linear classifiers in 2 dimensions}\}$  is 3.*



**□ Theorem (Vapnik)** – Let  $\mathcal{H}$  be given, with  $\text{VC}(\mathcal{H}) = d$  and  $m$  the number of training examples. With probability at least  $1 - \delta$ , we have:

$$\epsilon(\hat{h}) \leq \left( \min_{h \in \mathcal{H}} \epsilon(h) \right) + O \left( \sqrt{\frac{d}{m} \log \left( \frac{m}{d} \right)} + \frac{1}{m} \log \left( \frac{1}{\delta} \right) \right)$$

## 2 Unsupervised Learning

### 2.1 Introduction to Unsupervised Learning

**Motivation** – The goal of unsupervised learning is to find hidden patterns in unlabeled data  $\{x^{(1)}, \dots, x^{(m)}\}$ .

**Jensen's inequality** – Let  $f$  be a convex function and  $X$  a random variable. We have the following inequality:

$$E[f(X)] \geq f(E[X])$$

### 2.2 Clustering

#### 2.2.1 Expectation-Maximization

**Latent variables** – Latent variables are hidden/unobserved variables that make estimation problems difficult, and are often denoted  $z$ . Here are the most common settings where there are latent variables:

| Setting                  | Latent variable $z$   | $x z$                                | Comments  |
|--------------------------|-----------------------|--------------------------------------|---|
| Mixture of $k$ Gaussians | Multinomial( $\phi$ ) | $\mathcal{N}(\mu_j, \Sigma_j)$       | $\mu_j \in \mathbb{R}^n, \phi \in \mathbb{R}^k$ |
| Factor analysis          | $\mathcal{N}(0, I)$   | $\mathcal{N}(\mu + \Lambda z, \psi)$ | $\mu_j \in \mathbb{R}^n$                        |

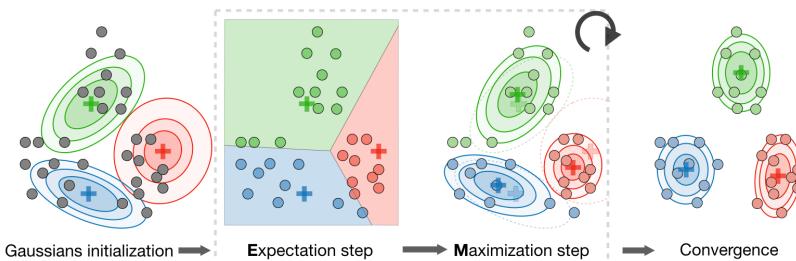
**Algorithm** – The Expectation-Maximization (EM) algorithm gives an efficient method at estimating the parameter  $\theta$  through maximum likelihood estimation by repeatedly constructing a lower-bound on the likelihood (E-step) and optimizing that lower bound (M-step) as follows:

- E-step: Evaluate the posterior probability  $Q_i(z^{(i)})$  that each data point  $x^{(i)}$  came from a particular cluster  $z^{(i)}$  as follows:

$$Q_i(z^{(i)}) = P(z^{(i)}|x^{(i)}; \theta)$$

- M-step: Use the posterior probabilities  $Q_i(z^{(i)})$  as cluster specific weights on data points  $x^{(i)}$  to separately re-estimate each cluster model as follows:

$$\theta_i = \underset{\theta}{\operatorname{argmax}} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log \left( \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) dz^{(i)}$$



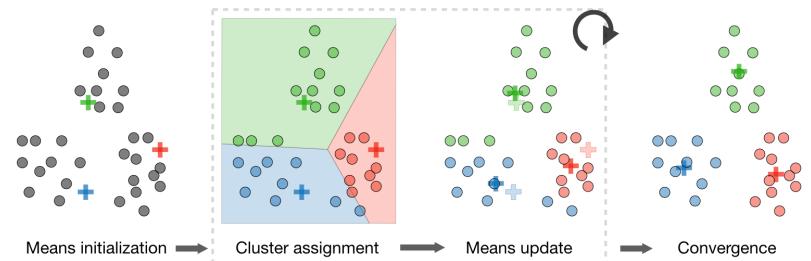
#### 2.2.2 $k$ -means clustering

We note  $c^{(i)}$  the cluster of data point  $i$  and  $\mu_j$  the center of cluster  $j$ .

**Algorithm** – After randomly initializing the cluster centroids  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ , the  $k$ -means algorithm repeats the following step until convergence:

$$c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|^2$$

$$\mu_j = \frac{\sum_{i=1}^m 1_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{c^{(i)}=j\}}}$$



**Distortion function** – In order to see if the algorithm converges, we look at the distortion function defined as follows:

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

#### 2.2.3 Hierarchical clustering

**Algorithm** – It is a clustering algorithm with an agglomerative hierarchical approach that build nested clusters in a successive manner.

**Types** – There are different sorts of hierarchical clustering algorithms that aims at optimizing different objective functions, which is summed up in the table below:

| Ward linkage                     | Average linkage                                 | Complete linkage                                   |
|----------------------------------|---|--|
| Minimize within cluster distance | Minimize average distance between cluster pairs | Minimize maximum distance of between cluster pairs |

#### 2.2.4 Clustering assessment metrics

In an unsupervised learning setting, it is often hard to assess the performance of a model since we don't have the ground truth labels as was the case in the supervised learning setting.

**Silhouette coefficient** – By noting  $a$  and  $b$  the mean distance between a sample and all other points in the same class, and between a sample and all other points in the next nearest cluster, the silhouette coefficient  $s$  for a single sample is defined as follows:

$$s = \frac{b - a}{\max(a, b)}$$

**Calinski-Harabaz index** – By noting  $k$  the number of clusters,  $B_k$  and  $W_k$  the between and within-clustering dispersion matrices respectively defined as

$$B_k = \sum_{j=1}^k n_{c(i)} (\mu_{c(i)} - \mu)(\mu_{c(i)} - \mu)^T, \quad W_k = \sum_{i=1}^m (x^{(i)} - \mu_{c(i)})(x^{(i)} - \mu_{c(i)})^T$$

the Calinski-Harabaz index  $s(k)$  indicates how well a clustering model defines its clusters, such that the higher the score, the more dense and well separated the clusters are. It is defined as follows:

$$s(k) = \frac{\text{Tr}(B_k)}{\text{Tr}(W_k)} \times \frac{N - k}{k - 1}$$

## 2.3 Dimension reduction

### 2.3.1 Principal component analysis

It is a dimension reduction technique that finds the variance maximizing directions onto which to project the data.

**Eigenvalue, eigenvector** – Given a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of  $A$  if there exists a vector  $z \in \mathbb{R}^n \setminus \{0\}$ , called eigenvector, such that we have:

$$Az = \lambda z$$

**Spectral theorem** – Let  $A \in \mathbb{R}^{n \times n}$ . If  $A$  is symmetric, then  $A$  is diagonalizable by a real orthogonal matrix  $U \in \mathbb{R}^{n \times n}$ . By noting  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ , we have:

$$\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$$

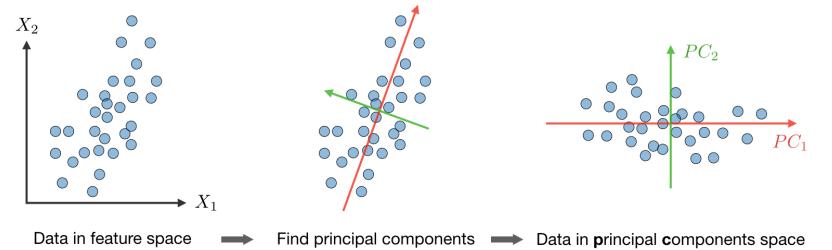
*Remark: the eigenvector associated with the largest eigenvalue is called principal eigenvector of matrix  $A$ .*

**Algorithm** – The Principal Component Analysis (PCA) procedure is a dimension reduction technique that projects the data on  $k$  dimensions by maximizing the variance of the data as follows:

- Step 1: Normalize the data to have a mean of 0 and standard deviation of 1.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{where} \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{and} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

- Step 2: Compute  $\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \in \mathbb{R}^{n \times n}$ , which is symmetric with real eigenvalues.
- Step 3: Compute  $u_1, \dots, u_k \in \mathbb{R}^n$  the  $k$  orthogonal principal eigenvectors of  $\Sigma$ , i.e. the orthogonal eigenvectors of the  $k$  largest eigenvalues.
- Step 4: Project the data on  $\text{span}_{\mathbb{R}}(u_1, \dots, u_k)$ . This procedure maximizes the variance among all  $k$ -dimensional spaces.



### 2.3.2 Independent component analysis

It is a technique meant to find the underlying generating sources.

**Assumptions** – We assume that our data  $x$  has been generated by the  $n$ -dimensional source vector  $s = (s_1, \dots, s_n)$ , where  $s_i$  are independent random variables, via a mixing and non-singular matrix  $A$  as follows:

$$x = As$$

The goal is to find the unmixing matrix  $W = A^{-1}$  by an update rule.

**Bell and Sejnowski ICA algorithm** – This algorithm finds the unmixing matrix  $W$  by following the steps below:

- Write the probability of  $x = As = W^{-1}s$  as:

$$p(x) = \prod_{i=1}^n p_s(w_i^T x) \cdot |W|$$

- Write the log likelihood given our training data  $\{x^{(i)}, i \in [1, m]\}$  and by noting  $g$  the sigmoid function as:

$$l(W) = \sum_{i=1}^m \left( \sum_{j=1}^n \log \left( g'(w_j^T x^{(i)}) \right) + \log |W| \right)$$

Therefore, the stochastic gradient ascent learning rule is such that for each training example  $x^{(i)}$ , we update  $W$  as follows:

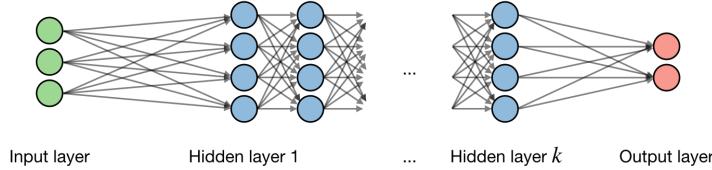
$$W \leftarrow W + \alpha \left( \begin{pmatrix} 1 - 2g(w_1^T x^{(i)}) \\ 1 - 2g(w_2^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{pmatrix} x^{(i)T} + (W^T)^{-1} \right)$$

### 3 Deep Learning

#### 3.1 Neural Networks

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

**Architecture** – The vocabulary around neural networks architectures is described in the figure below:

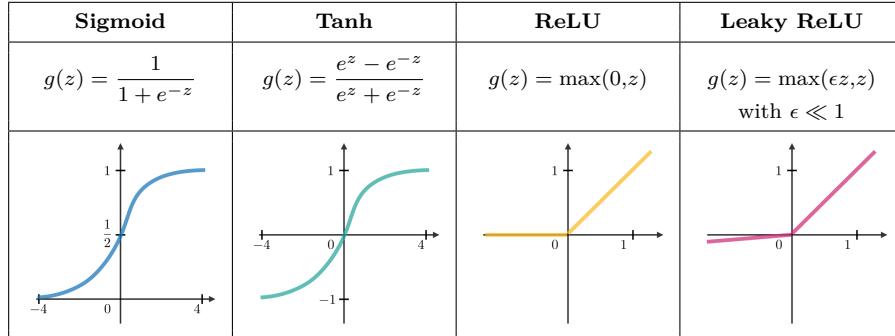


By noting  $i$  the  $i^{th}$  layer of the network and  $j$  the  $j^{th}$  hidden unit of the layer, we have:

$$z_j^{[i]} = w_j^{[i]T} x + b_j^{[i]}$$

where we note  $w$ ,  $b$ ,  $z$  the weight, bias and output respectively.

**Activation function** – Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:



**Cross-entropy loss** – In the context of neural networks, the cross-entropy loss  $L(z, y)$  is commonly used and is defined as follows:

$$L(z, y) = - \left[ y \log(z) + (1 - y) \log(1 - z) \right]$$

**Learning rate** – The learning rate, often noted  $\eta$ , indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

**Backpropagation** – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight  $w$  is computed using chain rule and is of the following form:

$$\frac{\partial L(z, y)}{\partial w} = \frac{\partial L(z, y)}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

As a result, the weight is updated as follows:

$$w \leftarrow w - \eta \frac{\partial L(z, y)}{\partial w}$$

**Updating weights** – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.
- Step 2: Perform forward propagation to obtain the corresponding loss.
- Step 3: Backpropagate the loss to get the gradients.
- Step 4: Use the gradients to update the weights of the network.

**Dropout** – Dropout is a technique meant at preventing overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability  $p$  or kept with probability  $1 - p$ .

#### 3.2 Convolutional Neural Networks

**Convolutional layer requirement** – By noting  $W$  the input volume size,  $F$  the size of the convolutional layer neurons,  $P$  the amount of zero padding, then the number of neurons  $N$  that fit in a given volume is such that:

$$N = \frac{W - F + 2P}{S} + 1$$

**Batch normalization** – It is a step of hyperparameter  $\gamma, \beta$  that normalizes the batch  $\{x_i\}$ . By noting  $\mu_B, \sigma_B^2$  the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \leftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

#### 3.3 Recurrent Neural Networks

**Types of gates** – Here are the different types of gates that we encounter in a typical recurrent neural network:

| Input gate            | Forget gate          | Output gate           | Gate              |
|-----------------------|----------------------|-----------------------|-------------------|
| Write to cell or not? | Erase a cell or not? | Reveal a cell or not? | How much writing? |

**LSTM** – A long short-term memory (LSTM) network is a type of RNN model that avoids the vanishing gradient problem by adding ‘forget’ gates.

### 3.4 Reinforcement Learning and Control

The goal of reinforcement learning is for an agent to learn how to evolve in an environment.

**□ Markov decision processes** – A Markov decision process (MDP) is a 5-tuple  $(S, A, \{P_{sa}\}, \gamma, R)$  where:

- $\mathcal{S}$  is the set of states
- $\mathcal{A}$  is the set of actions
- $\{P_{sa}\}$  are the state transition probabilities for  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$
- $\gamma \in [0, 1[$  is the discount factor
- $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  or  $R : \mathcal{S} \rightarrow \mathbb{R}$  is the reward function that the algorithm wants to maximize

**□ Policy** – A policy  $\pi$  is a function  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maps states to actions.

*Remark: we say that we execute a given policy  $\pi$  if given a state  $s$  we take the action  $a = \pi(s)$ .*

**□ Value function** – For a given policy  $\pi$  and a given state  $s$ , we define the value function  $V^\pi$  as follows:

$$V^\pi(s) = E \left[ R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi \right]$$

**□ Bellman equation** – The optimal Bellman equations characterizes the value function  $V^{\pi^*}$  of the optimal policy  $\pi^*$ :

$$V^{\pi^*}(s) = R(s) + \max_{a \in \mathcal{A}} \gamma \sum_{s' \in \mathcal{S}} P_{sa}(s') V^{\pi^*}(s')$$

*Remark: we note that the optimal policy  $\pi^*$  for a given state  $s$  is such that:*

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{sa}(s') V^*(s')$$

**□ Value iteration algorithm** – The value iteration algorithm is in two steps:

- We initialize the value:

$$V_0(s) = 0$$

- We iterate the value based on the values before:

$$V_{i+1}(s) = R(s) + \max_{a \in \mathcal{A}} \left[ \sum_{s' \in \mathcal{S}} \gamma P_{sa}(s') V_i(s') \right]$$

**□ Maximum likelihood estimate** – The maximum likelihood estimates for the state transition probabilities are as follows:

$$P_{sa}(s') = \frac{\text{\#times took action } a \text{ in state } s \text{ and got to } s'}{\text{\#times took action } a \text{ in state } s}$$

**□ Q-learning** –  $Q$ -learning is a model-free estimation of  $Q$ , which is done as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

## 4 Machine Learning Tips and Tricks

### 4.1 Metrics

Given a set of data points  $\{x^{(1)}, \dots, x^{(m)}\}$ , where each  $x^{(i)}$  has  $n$  features, associated to a set of outcomes  $\{y^{(1)}, \dots, y^{(m)}\}$ , we want to assess a given classifier that learns how to predict  $y$  from  $x$ .

#### 4.1.1 Classification

In a context of a binary classification, here are the main metrics that are important to track to assess the performance of the model.

**Confusion matrix** – The confusion matrix is used to have a more complete picture when assessing the performance of a model. It is defined as follows:

|              |   | Predicted class                       |  |
|--------------|---|---------------------------------------|--|
|              |   | +                                     | -                                      |
| Actual class | + | TP<br>True Positives                  | FN<br>False Negatives<br>Type II error |
|              | - | FP<br>False Positives<br>Type I error | TN<br>True Negatives                   |

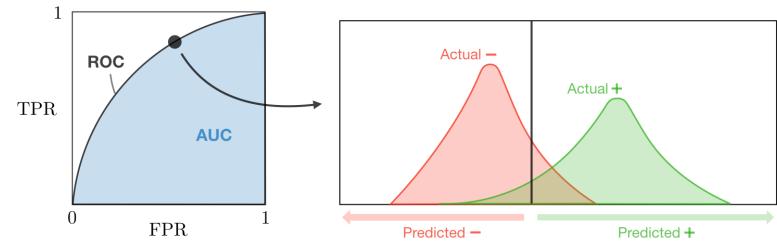
**Main metrics** – The following metrics are commonly used to assess the performance of classification models:

| Metric                | Formula                             | Interpretation                              |
|-----------------------|-------------------------------------|---|
| Accuracy              | $\frac{TP + TN}{TP + TN + FP + FN}$ | Overall performance of model                |
| Precision             | $\frac{TP}{TP + FP}$                | How accurate the positive predictions are   |
| Recall<br>Sensitivity | $\frac{TP}{TP + FN}$                | Coverage of actual positive sample          |
| Specificity           | $\frac{TN}{TN + FP}$                | Coverage of actual negative sample          |
| F1 score              | $\frac{2TP}{2TP + FP + FN}$         | Hybrid metric useful for unbalanced classes |

**ROC** – The receiver operating curve, also noted ROC, is the plot of TPR versus FPR by varying the threshold. These metrics are summed up in the table below:

| Metric                     | Formula              | Equivalent          |
|----------------------------|----------------------|---------------------|
| True Positive Rate<br>TPR  | $\frac{TP}{TP + FN}$ | Recall, sensitivity |
| False Positive Rate<br>FPR | $\frac{FP}{TN + FP}$ | 1-specificity       |

**AUC** – The area under the receiving operating curve, also noted AUC or AUROC, is the area below the ROC as shown in the following figure:



#### 4.1.2 Regression

**Basic metrics** – Given a regression model  $f$ , the following metrics are commonly used to assess the performance of the model:

| Total sum of squares                        | Explained sum of squares                       | Residual sum of squares                    |
|---|--|--|
| $SS_{tot} = \sum_{i=1}^m (y_i - \bar{y})^2$ | $SS_{reg} = \sum_{i=1}^m (f(x_i) - \bar{y})^2$ | $SS_{res} = \sum_{i=1}^m (y_i - f(x_i))^2$ |

**Coefficient of determination** – The coefficient of determination, often noted  $R^2$  or  $r^2$ , provides a measure of how well the observed outcomes are replicated by the model and is defined as follows:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

**Main metrics** – The following metrics are commonly used to assess the performance of regression models, by taking into account the number of variables  $n$  that they take into consideration:

| Mallow's Cp                                 | AIC                  | BIC                       | Adjusted $R^2$                   |
|---|----------------------|---------------------------|----------------------------------|
| $\frac{SS_{res} + 2(n+1)\hat{\sigma}^2}{m}$ | $2[(n+2) - \log(L)]$ | $\log(m)(n+2) - 2\log(L)$ | $1 - \frac{(1-R^2)(m-1)}{m-n-1}$ |

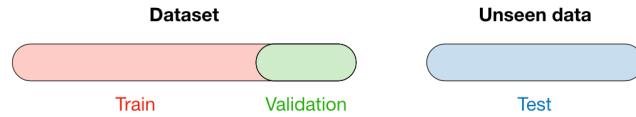
where  $L$  is the likelihood and  $\hat{\sigma}^2$  is an estimate of the variance associated with each response.

## 4.2 Model selection

**Vocabulary** – When selecting a model, we distinguish 3 different parts of the data that we have as follows:

| Training set                                       | Validation set   | Testing set                                |
|--|--|--|
| - Model is trained<br>- Usually 80% of the dataset | - Model is assessed<br>- Usually 20% of the dataset<br>- Also called hold-out or development set | - Model gives predictions<br>- Unseen data |

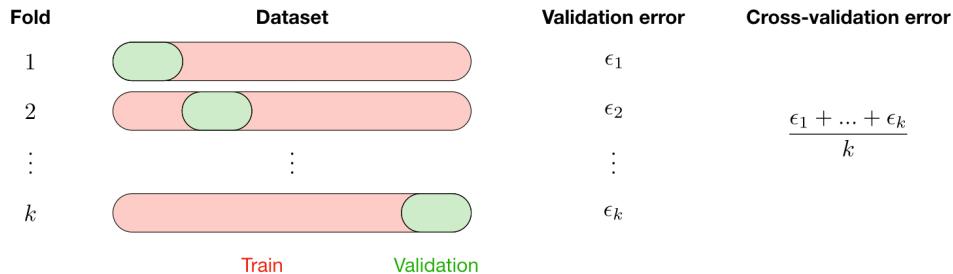
Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:



**Cross-validation** – Cross-validation, also noted CV, is a method that is used to select a model that does not rely too much on the initial training set. The different types are summed up in the table below:

| <i>k</i> -fold   | Leave- <i>p</i> -out  |
|--|---|
| - Training on $k - 1$ folds and assessment on the remaining one<br>- Generally $k = 5$ or 10 | - Training on $n - p$ observations and assessment on the $p$ remaining ones<br>- Case $p = 1$ is called leave-one-out |

The most commonly used method is called *k*-fold cross-validation and splits the training data into *k* folds to validate the model on one fold while training the model on the  $k - 1$  other folds, all of this *k* times. The error is then averaged over the *k* folds and is named cross-validation error.



**Regularization** – The regularization procedure aims at avoiding the model to overfit the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:

| LASSO  | Ridge  | Elastic Net   |
|--|--|---|
| - Shrinks coefficients to 0<br>- Good for variable selection | Makes coefficients smaller                                   | Tradeoff between variable selection and small coefficients  |
|  |  |   |
| $\dots + \lambda \ \theta\ _1$<br>$\lambda \in \mathbb{R}$   | $\dots + \lambda \ \theta\ _2^2$<br>$\lambda \in \mathbb{R}$ | $\dots + \lambda [(1-\alpha)\ \theta\ _1 + \alpha\ \theta\ _2^2]$<br>$\lambda \in \mathbb{R}, \alpha \in [0,1]$ |

**Model selection** – Train model on training set, then evaluate on the development set, then pick best performance model on the development set, and retrain all of that model on the whole training set.

## 4.3 Diagnostics

**Bias** – The bias of a model is the difference between the expected prediction and the correct model that we try to predict for given data points.

**Variance** – The variance of a model is the variability of the model prediction for given data points.

**Bias/variance tradeoff** – The simpler the model, the higher the bias, and the more complex the model, the higher the variance.

|                   | Underfitting   | Just right   | Overfitting  |
|-------------------|--|--|--|
| <b>Symptoms</b>   | - High training error<br>- Training error close to test error<br>- High bias | - Training error slightly lower than test error<br>- High variance | - Low training error<br>- Training error much lower than test error<br>- High variance |
| <b>Regression</b> |  |  |  |

|                |   |   |
|----------------|---|---|
| Classification |   |   |
| Deep learning  |   |   |
| Remedies       | <ul style="list-style-type: none"> <li>- Complexify model</li> <li>- Add more features</li> <li>- Train longer</li> </ul> | <ul style="list-style-type: none"> <li>- Regularize</li> <li>- Get more data</li> </ul> |

□ **Error analysis** – Error analysis is analyzing the root cause of the difference in performance between the current and the perfect models.

□ **Ablative analysis** – Ablative analysis is analyzing the root cause of the difference in performance between the current and the baseline models.

## 5 Refreshers

### 5.1 Probabilities and Statistics

#### 5.1.1 Introduction to Probability and Combinatorics

□ **Sample space** – The set of all possible outcomes of an experiment is known as the sample space of the experiment and is denoted by  $S$ .

□ **Event** – Any subset  $E$  of the sample space is known as an event. That is, an event is a set consisting of possible outcomes of the experiment. If the outcome of the experiment is contained in  $E$ , then we say that  $E$  has occurred.

□ **Axioms of probability** – For each event  $E$ , we denote  $P(E)$  as the probability of event  $E$  occurring. By noting  $E_1, \dots, E_n$  mutually exclusive events, we have the 3 following axioms:

$$(1) \quad 0 \leq P(E) \leq 1 \quad (2) \quad P(S) = 1 \quad (3) \quad P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i)$$

□ **Permutation** – A permutation is an arrangement of  $r$  objects from a pool of  $n$  objects, in a given order. The number of such arrangements is given by  $P(n, r)$ , defined as:

$$P(n, r) = \frac{n!}{(n - r)!}$$

□ **Combination** – A combination is an arrangement of  $r$  objects from a pool of  $n$  objects, where the order does not matter. The number of such arrangements is given by  $C(n, r)$ , defined as:

$$C(n, r) = \frac{P(n, r)}{r!} = \frac{n!}{r!(n - r)!}$$

*Remark: we note that for  $0 \leq r \leq n$ , we have  $P(n, r) \geq C(n, r)$ .*

#### 5.1.2 Conditional Probability

□ **Bayes' rule** – For events  $A$  and  $B$  such that  $P(B) > 0$ , we have:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

*Remark: we have  $P(A \cap B) = P(A)P(B|A) = P(A|B)P(B)$ .*

□ **Partition** – Let  $\{A_i, i \in [1, n]\}$  be such that for all  $i$ ,  $A_i \neq \emptyset$ . We say that  $\{A_i\}$  is a partition if we have:

$$\forall i \neq j, A_i \cap A_j = \emptyset \quad \text{and} \quad \bigcup_{i=1}^n A_i = S$$

*Remark: for any event  $B$  in the sample space, we have  $P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$ .*

**□ Extended form of Bayes' rule** – Let  $\{A_i, i \in [1, n]\}$  be a partition of the sample space. We have:

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{\sum_{i=1}^n P(B|A_i)P(A_i)}$$

**□ Independence** – Two events  $A$  and  $B$  are independent if and only if we have:

$$P(A \cap B) = P(A)P(B)$$

### 5.1.3 Random Variables

**□ Random variable** – A random variable, often noted  $X$ , is a function that maps every element in a sample space to a real line.

**□ Cumulative distribution function (CDF)** – The cumulative distribution function  $F$ , which is monotonically non-decreasing and is such that  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow +\infty} F(x) = 1$ , is defined as:

$$F(x) = P(X \leq x)$$

Remark: we have  $P(a < X \leq b) = F(b) - F(a)$ .

**□ Probability density function (PDF)** – The probability density function  $f$  is the probability that  $X$  takes on values between two adjacent realizations of the random variable.

**□ Relationships involving the PDF and CDF** – Here are the important properties to know in the discrete (D) and the continuous (C) cases.

| Case | CDF $F$                               | PDF $f$                | Properties of PDF                                       |
|------|---------------------------------------|------------------------|---|
| (D)  | $F(x) = \sum_{x_i \leq x} P(X = x_i)$ | $f(x_j) = P(X = x_j)$  | $0 \leq f(x_j) \leq 1$ and $\sum_j f(x_j) = 1$          |
| (C)  | $F(x) = \int_{-\infty}^x f(y)dy$      | $f(x) = \frac{dF}{dx}$ | $f(x) \geq 0$ and $\int_{-\infty}^{+\infty} f(x)dx = 1$ |

**□ Variance** – The variance of a random variable, often noted  $\text{Var}(X)$  or  $\sigma^2$ , is a measure of the spread of its distribution function. It is determined as follows:

$$\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

**□ Standard deviation** – The standard deviation of a random variable, often noted  $\sigma$ , is a measure of the spread of its distribution function which is compatible with the units of the actual random variable. It is determined as follows:

$$\sigma = \sqrt{\text{Var}(X)}$$

**□ Expectation and Moments of the Distribution** – Here are the expressions of the expected value  $E[X]$ , generalized expected value  $E[g(X)]$ ,  $k^{\text{th}}$  moment  $E[X^k]$  and characteristic function  $\psi(\omega)$  for the discrete and continuous cases:

| Case | $E[X]$                             | $E[g(X)]$                             | $E[X^k]$                              | $\psi(\omega)$                                   |
|------|------------------------------------|---------------------------------------|---------------------------------------|--|
| (D)  | $\sum_{i=1}^n x_i f(x_i)$          | $\sum_{i=1}^n g(x_i) f(x_i)$          | $\sum_{i=1}^n x_i^k f(x_i)$           | $\sum_{i=1}^n f(x_i) e^{i\omega x_i}$            |
| (C)  | $\int_{-\infty}^{+\infty} xf(x)dx$ | $\int_{-\infty}^{+\infty} g(x)f(x)dx$ | $\int_{-\infty}^{+\infty} x^k f(x)dx$ | $\int_{-\infty}^{+\infty} f(x) e^{i\omega x} dx$ |

Remark: we have  $e^{i\omega x} = \cos(\omega x) + i \sin(\omega x)$ .

**□ Revisiting the  $k^{\text{th}}$  moment** – The  $k^{\text{th}}$  moment can also be computed with the characteristic function as follows:

$$E[X^k] = \frac{1}{i^k} \left[ \frac{\partial^k \psi}{\partial \omega^k} \right]_{\omega=0}$$

**□ Transformation of random variables** – Let the variables  $X$  and  $Y$  be linked by some function. By noting  $f_X$  and  $f_Y$  the distribution function of  $X$  and  $Y$  respectively, we have:

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right|$$

**□ Leibniz integral rule** – Let  $g$  be a function of  $x$  and potentially  $c$ , and  $a, b$  boundaries that may depend on  $c$ . We have:

$$\frac{\partial}{\partial c} \left( \int_a^b g(x)dx \right) = \frac{\partial b}{\partial c} \cdot g(b) - \frac{\partial a}{\partial c} \cdot g(a) + \int_a^b \frac{\partial g}{\partial c}(x)dx$$

**□ Chebyshev's inequality** – Let  $X$  be a random variable with expected value  $\mu$  and standard deviation  $\sigma$ . For  $k, \sigma > 0$ , we have the following inequality:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

### 5.1.4 Jointly Distributed Random Variables

**□ Conditional density** – The conditional density of  $X$  with respect to  $Y$ , often noted  $f_{X|Y}$ , is defined as follows:

$$f_{X|Y}(x) = \frac{f_{XY}(x,y)}{f_Y(y)}$$

**□ Independence** – Two random variables  $X$  and  $Y$  are said to be independent if we have:

$$f_{XY}(x,y) = f_X(x)f_Y(y)$$

□ **Marginal density and cumulative distribution** – From the joint density probability function  $f_{XY}$ , we have:

| Case | Marginal density                                    | Cumulative function   |
|------|---|---|
| (D)  | $f_X(x_i) = \sum_j f_{XY}(x_i, y_j)$                | $F_{XY}(x, y) = \sum_{x_i \leq x} \sum_{y_j \leq y} f_{XY}(x_i, y_j)$     |
| (C)  | $f_X(x) = \int_{-\infty}^{+\infty} f_{XY}(x, y) dy$ | $F_{XY}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{XY}(x', y') dx' dy'$ |

□ **Distribution of a sum of independent random variables** – Let  $Y = X_1 + \dots + X_n$  with  $X_1, \dots, X_n$  independent. We have:

$$\psi_Y(\omega) = \prod_{k=1}^n \psi_{X_k}(\omega)$$

□ **Covariance** – We define the covariance of two random variables  $X$  and  $Y$ , that we note  $\sigma_{XY}^2$  or more commonly  $\text{Cov}(X, Y)$ , as follows:

$$\text{Cov}(X, Y) \triangleq \sigma_{XY}^2 = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y$$

□ **Correlation** – By noting  $\sigma_X, \sigma_Y$  the standard deviations of  $X$  and  $Y$ , we define the correlation between the random variables  $X$  and  $Y$ , noted  $\rho_{XY}$ , as follows:

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

Remarks: For any  $X, Y$ , we have  $\rho_{XY} \in [-1, 1]$ . If  $X$  and  $Y$  are independent, then  $\rho_{XY} = 0$ .

□ **Main distributions** – Here are the main distributions to have in mind:

| Type | Distribution                                  | PDF   | $\psi(\omega)$                                       | $E[X]$              | $\text{Var}(X)$       |
|------|---|---|--|---------------------|-----------------------|
| (D)  | $X \sim \mathcal{B}(n, p)$<br>Binomial        | $P(X = x) = \binom{n}{x} p^x q^{n-x}$<br>$x \in \llbracket 0, n \rrbracket$                                     | $(pe^{i\omega} + q)^n$                               | $np$                | $npq$                 |
|      | $X \sim \text{Po}(\mu)$<br>Poisson            | $P(X = x) = \frac{\mu^x}{x!} e^{-\mu}$<br>$x \in \mathbb{N}$  | $e^{\mu(e^{i\omega}-1)}$                             | $\mu$               | $\mu$                 |
| (C)  | $X \sim \mathcal{U}(a, b)$<br>Uniform         | $f(x) = \frac{1}{b-a}$<br>$x \in [a, b]$  | $\frac{e^{i\omega b} - e^{i\omega a}}{(b-a)i\omega}$ | $\frac{a+b}{2}$     | $\frac{(b-a)^2}{12}$  |
|      | $X \sim \mathcal{N}(\mu, \sigma)$<br>Gaussian | $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$<br>$x \in \mathbb{R}$ | $e^{i\omega\mu - \frac{1}{2}\omega^2\sigma^2}$       | $\mu$               | $\sigma^2$            |
|      | $X \sim \text{Exp}(\lambda)$<br>Exponential   | $f(x) = \lambda e^{-\lambda x}$<br>$x \in \mathbb{R}_+$   | $\frac{1}{1 - \frac{i\omega}{\lambda}}$              | $\frac{1}{\lambda}$ | $\frac{1}{\lambda^2}$ |

## 5.1.5 Parameter estimation

□ **Random sample** – A random sample is a collection of  $n$  random variables  $X_1, \dots, X_n$  that are independent and identically distributed with  $X$ .

□ **Estimator** – An estimator  $\hat{\theta}$  is a function of the data that is used to infer the value of an unknown parameter  $\theta$  in a statistical model.

□ **Bias** – The bias of an estimator  $\hat{\theta}$  is defined as being the difference between the expected value of the distribution of  $\hat{\theta}$  and the true value, i.e.:

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

Remark: an estimator is said to be unbiased when we have  $E[\hat{\theta}] = \theta$ .

□ **Sample mean and variance** – The sample mean and the sample variance of a random sample are used to estimate the true mean  $\mu$  and the true variance  $\sigma^2$  of a distribution, are noted  $\bar{X}$  and  $s^2$  respectively, and are such that:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad s^2 = \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

□ **Central Limit Theorem** – Let us have a random sample  $X_1, \dots, X_n$  following a given distribution with mean  $\mu$  and variance  $\sigma^2$ , then we have:

$$\bar{X} \underset{n \rightarrow +\infty}{\sim} \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

## 5.2 Linear Algebra and Calculus

### 5.2.1 General notations

□ **Vector** – We note  $x \in \mathbb{R}^n$  a vector with  $n$  entries, where  $x_i \in \mathbb{R}$  is the  $i^{th}$  entry:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$$

□ **Matrix** – We note  $A \in \mathbb{R}^{m \times n}$  a matrix with  $m$  rows and  $n$  columns, where  $A_{i,j} \in \mathbb{R}$  is the entry located in the  $i^{th}$  row and  $j^{th}$  column:

$$A = \begin{pmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Remark: the vector  $x$  defined above can be viewed as a  $n \times 1$  matrix and is more particularly called a column-vector.

□ **Identity matrix** – The identity matrix  $I \in \mathbb{R}^{n \times n}$  is a square matrix with ones in its diagonal and zero everywhere else:

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

Remark: for all matrices  $A \in \mathbb{R}^{n \times n}$ , we have  $A \times I = I \times A = A$ .

**□ Diagonal matrix** – A diagonal matrix  $D \in \mathbb{R}^{n \times n}$  is a square matrix with nonzero values in its diagonal and zero everywhere else:

$$D = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & d_n \end{pmatrix}$$

Remark: we also note  $D$  as  $\text{diag}(d_1, \dots, d_n)$ .

## 5.2.2 Matrix operations

**□ Vector-vector multiplication** – There are two types of vector-vector products:

- inner product: for  $x, y \in \mathbb{R}^n$ , we have:

$$x^T y = \sum_{i=1}^n x_i y_i \in \mathbb{R}$$

- outer product: for  $x \in \mathbb{R}^m, y \in \mathbb{R}^n$ , we have:

$$xy^T = \begin{pmatrix} x_1 y_1 & \cdots & x_1 y_n \\ \vdots & & \vdots \\ x_m y_1 & \cdots & x_m y_n \end{pmatrix} \in \mathbb{R}^{m \times n}$$

**□ Matrix-vector multiplication** – The product of matrix  $A \in \mathbb{R}^{m \times n}$  and vector  $x \in \mathbb{R}^n$  is a vector of size  $\mathbb{R}^m$ , such that:

$$Ax = \begin{pmatrix} a_{r,1}^T x \\ \vdots \\ a_{r,m}^T x \end{pmatrix} = \sum_{i=1}^n a_{c,i} x_i \in \mathbb{R}^m$$

where  $a_{r,i}^T$  are the vector rows and  $a_{c,j}$  are the vector columns of  $A$ , and  $x_i$  are the entries of  $x$ .

**□ Matrix-matrix multiplication** – The product of matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$  is a matrix of size  $\mathbb{R}^{m \times p}$ , such that:

$$AB = \begin{pmatrix} a_{r,1}^T b_{c,1} & \cdots & a_{r,1}^T b_{c,p} \\ \vdots & & \vdots \\ a_{r,m}^T b_{c,1} & \cdots & a_{r,m}^T b_{c,p} \end{pmatrix} = \sum_{i=1}^n a_{c,i} b_{r,i}^T \in \mathbb{R}^{m \times p}$$

where  $a_{r,i}^T, b_{r,i}^T$  are the vector rows and  $a_{c,j}, b_{c,j}$  are the vector columns of  $A$  and  $B$  respectively.

**□ Transpose** – The transpose of a matrix  $A \in \mathbb{R}^{m \times n}$ , noted  $A^T$ , is such that its entries are flipped:

$$\forall i, j, \quad A_{i,j}^T = A_{j,i}$$

Remark: for matrices  $A, B$ , we have  $(AB)^T = B^T A^T$ .

**□ Inverse** – The inverse of an invertible square matrix  $A$  is noted  $A^{-1}$  and is the only matrix such that:

$$AA^{-1} = A^{-1}A = I$$

Remark: not all square matrices are invertible. Also, for matrices  $A, B$ , we have  $(AB)^{-1} = B^{-1}A^{-1}$

**□ Trace** – The trace of a square matrix  $A$ , noted  $\text{tr}(A)$ , is the sum of its diagonal entries:

$$\text{tr}(A) = \sum_{i=1}^n A_{i,i}$$

Remark: for matrices  $A, B$ , we have  $\text{tr}(A^T) = \text{tr}(A)$  and  $\text{tr}(AB) = \text{tr}(BA)$

**□ Determinant** – The determinant of a square matrix  $A \in \mathbb{R}^{n \times n}$ , noted  $|A|$  or  $\det(A)$  is expressed recursively in terms of  $A_{\setminus i, \setminus j}$ , which is the matrix  $A$  without its  $i^{\text{th}}$  row and  $j^{\text{th}}$  column, as follows:

$$\det(A) = |A| = \sum_{j=1}^n (-1)^{i+j} A_{i,j} |A_{\setminus i, \setminus j}|$$

Remark:  $A$  is invertible if and only if  $|A| \neq 0$ . Also,  $|AB| = |A||B|$  and  $|A^T| = |A|$ .

## 5.2.3 Matrix properties

**□ Symmetric decomposition** – A given matrix  $A$  can be expressed in terms of its symmetric and antisymmetric parts as follows:

$$A = \underbrace{\frac{A + A^T}{2}}_{\text{Symmetric}} + \underbrace{\frac{A - A^T}{2}}_{\text{Antisymmetric}}$$

**□ Norm** – A norm is a function  $N : V \rightarrow [0, +\infty]$  where  $V$  is a vector space, and such that for all  $x, y \in V$ , we have:

- $N(x + y) \leq N(x) + N(y)$
- $N(ax) = |a|N(x)$  for  $a$  scalar
- if  $N(x) = 0$ , then  $x = 0$

For  $x \in V$ , the most commonly used norms are summed up in the table below:

| Norm                 | Notation       | Definition                                      | Use case             |
|----------------------|----------------|---|----------------------|
| Manhattan, $L^1$     | $\ x\ _1$      | $\sum_{i=1}^n  x_i $                            | LASSO regularization |
| Euclidean, $L^2$     | $\ x\ _2$      | $\sqrt{\sum_{i=1}^n x_i^2}$                     | Ridge regularization |
| $p$ -norm, $L^p$     | $\ x\ _p$      | $\left(\sum_{i=1}^n x_i^p\right)^{\frac{1}{p}}$ | Hölder inequality    |
| Infinity, $L^\infty$ | $\ x\ _\infty$ | $\max_i  x_i $                                  | Uniform convergence  |

□ **Linearly dependence** – A set of vectors is said to be linearly dependent if one of the vectors in the set can be defined as a linear combination of the others.

*Remark: if no vector can be written this way, then the vectors are said to be linearly independent.*

□ **Matrix rank** – The rank of a given matrix  $A$  is noted  $\text{rank}(A)$  and is the dimension of the vector space generated by its columns. This is equivalent to the maximum number of linearly independent columns of  $A$ .

□ **Positive semi-definite matrix** – A matrix  $A \in \mathbb{R}^{n \times n}$  is positive semi-definite (PSD) and is noted  $A \succeq 0$  if we have:

$$A = A^T \quad \text{and} \quad \forall x \in \mathbb{R}^n, \quad x^T Ax \geq 0$$

*Remark: similarly, a matrix  $A$  is said to be positive definite, and is noted  $A \succ 0$ , if it is a PSD matrix which satisfies for all non-zero vector  $x$ ,  $x^T Ax > 0$ .*

□ **Eigenvalue, eigenvector** – Given a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of  $A$  if there exists a vector  $z \in \mathbb{R}^n \setminus \{0\}$ , called eigenvector, such that we have:

$$Az = \lambda z$$

□ **Spectral theorem** – Let  $A \in \mathbb{R}^{n \times n}$ . If  $A$  is symmetric, then  $A$  is diagonalizable by a real orthogonal matrix  $U \in \mathbb{R}^{n \times n}$ . By noting  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ , we have:

$$\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$$

□ **Singular-value decomposition** – For a given matrix  $A$  of dimensions  $m \times n$ , the singular-value decomposition (SVD) is a factorization technique that guarantees the existence of  $U$   $m \times m$  unitary,  $\Sigma$   $m \times n$  diagonal and  $V$   $n \times n$  unitary matrices, such that:

$$A = U \Sigma V^T$$

## 5.2.4 Matrix calculus

□ **Gradient** – Let  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  be a function and  $A \in \mathbb{R}^{m \times n}$  be a matrix. The gradient of  $f$  with respect to  $A$  is a  $m \times n$  matrix, noted  $\nabla_A f(A)$ , such that:

$$\left( \nabla_A f(A) \right)_{i,j} = \frac{\partial f(A)}{\partial A_{i,j}}$$

*Remark: the gradient of  $f$  is only defined when  $f$  is a function that returns a scalar.*

□ **Hessian** – Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function and  $x \in \mathbb{R}^n$  be a vector. The hessian of  $f$  with respect to  $x$  is a  $n \times n$  symmetric matrix, noted  $\nabla_x^2 f(x)$ , such that:

$$\left( \nabla_x^2 f(x) \right)_{i,j} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

*Remark: the hessian of  $f$  is only defined when  $f$  is a function that returns a scalar.*

□ **Gradient operations** – For matrices  $A, B, C$ , the following gradient properties are worth having in mind:

$$\nabla_A \text{tr}(AB) = B^T$$

$$\nabla_A f(A) = (\nabla_A f(A))^T$$

$$\nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T$$

$$\nabla_A |A| = |A|(A^{-1})^T$$

# VIP Refresher: Linear Algebra and Calculus

Afshin AMIDI and Shervine AMIDI

October 6, 2018

## General notations

**Vector** – We note  $x \in \mathbb{R}^n$  a vector with  $n$  entries, where  $x_i \in \mathbb{R}$  is the  $i^{th}$  entry:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$$

**Matrix** – We note  $A \in \mathbb{R}^{m \times n}$  a matrix with  $m$  rows and  $n$  columns, where  $A_{i,j} \in \mathbb{R}$  is the entry located in the  $i^{th}$  row and  $j^{th}$  column:

$$A = \begin{pmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

*Remark: the vector  $x$  defined above can be viewed as a  $n \times 1$  matrix and is more particularly called a column-vector.*

**Identity matrix** – The identity matrix  $I \in \mathbb{R}^{n \times n}$  is a square matrix with ones in its diagonal and zero everywhere else:

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

*Remark: for all matrices  $A \in \mathbb{R}^{n \times n}$ , we have  $A \times I = I \times A = A$ .*

**Diagonal matrix** – A diagonal matrix  $D \in \mathbb{R}^{n \times n}$  is a square matrix with nonzero values in its diagonal and zero everywhere else:

$$D = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & d_n \end{pmatrix}$$

*Remark: we also note  $D$  as  $\text{diag}(d_1, \dots, d_n)$ .*

## Matrix operations

**Vector-vector multiplication** – There are two types of vector-vector products:

- inner product: for  $x, y \in \mathbb{R}^n$ , we have:

$$x^T y = \sum_{i=1}^n x_i y_i \in \mathbb{R}$$

- outer product: for  $x \in \mathbb{R}^m, y \in \mathbb{R}^n$ , we have:

$$xy^T = \begin{pmatrix} x_1 y_1 & \cdots & x_1 y_n \\ \vdots & & \vdots \\ x_m y_1 & \cdots & x_m y_n \end{pmatrix} \in \mathbb{R}^{m \times n}$$

**Matrix-vector multiplication** – The product of matrix  $A \in \mathbb{R}^{m \times n}$  and vector  $x \in \mathbb{R}^n$  is a vector of size  $\mathbb{R}^m$ , such that:

$$Ax = \begin{pmatrix} a_{r,1}^T x \\ \vdots \\ a_{r,m}^T x \end{pmatrix} = \sum_{i=1}^n a_{c,i} x_i \in \mathbb{R}^m$$

where  $a_{r,i}^T$  are the vector rows and  $a_{c,j}$  are the vector columns of  $A$ , and  $x_i$  are the entries of  $x$ .

**Matrix-matrix multiplication** – The product of matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$  is a matrix of size  $\mathbb{R}^{m \times p}$ , such that:

$$AB = \begin{pmatrix} a_{r,1}^T b_{c,1} & \cdots & a_{r,1}^T b_{c,p} \\ \vdots & & \vdots \\ a_{r,m}^T b_{c,1} & \cdots & a_{r,m}^T b_{c,p} \end{pmatrix} = \sum_{i=1}^n a_{c,i} b_{r,i}^T \in \mathbb{R}^{m \times p}$$

where  $a_{r,i}^T, b_{r,i}^T$  are the vector rows and  $a_{c,j}, b_{c,j}$  are the vector columns of  $A$  and  $B$  respectively.

**Transpose** – The transpose of a matrix  $A \in \mathbb{R}^{m \times n}$ , noted  $A^T$ , is such that its entries are flipped:

$$\forall i, j, \quad A_{i,j}^T = A_{j,i}$$

*Remark: for matrices  $A, B$ , we have  $(AB)^T = B^T A^T$ .*

**Inverse** – The inverse of an invertible square matrix  $A$  is noted  $A^{-1}$  and is the only matrix such that:

$$AA^{-1} = A^{-1}A = I$$

*Remark: not all square matrices are invertible. Also, for matrices  $A, B$ , we have  $(AB)^{-1} = B^{-1}A^{-1}$*

**Trace** – The trace of a square matrix  $A$ , noted  $\text{tr}(A)$ , is the sum of its diagonal entries:

$$\text{tr}(A) = \sum_{i=1}^n A_{i,i}$$

*Remark: for matrices  $A, B$ , we have  $\text{tr}(A^T) = \text{tr}(A)$  and  $\text{tr}(AB) = \text{tr}(BA)$*

**Determinant** – The determinant of a square matrix  $A \in \mathbb{R}^{n \times n}$ , noted  $|A|$  or  $\det(A)$  is expressed recursively in terms of  $A_{\setminus i, \setminus j}$ , which is the matrix  $A$  without its  $i^{th}$  row and  $j^{th}$  column, as follows:

$$\det(A) = |A| = \sum_{j=1}^n (-1)^{i+j} A_{i,j} |A_{\setminus i, \setminus j}|$$

*Remark:*  $A$  is invertible if and only if  $|A| \neq 0$ . Also,  $|AB| = |A||B|$  and  $|A^T| = |A|$ .

## Matrix properties

**Symmetric decomposition** – A given matrix  $A$  can be expressed in terms of its symmetric and antisymmetric parts as follows:

$$A = \underbrace{\frac{A + A^T}{2}}_{\text{Symmetric}} + \underbrace{\frac{A - A^T}{2}}_{\text{Antisymmetric}}$$

**Norm** – A norm is a function  $N : V \rightarrow [0, +\infty[$  where  $V$  is a vector space, and such that for all  $x, y \in V$ , we have:

- $N(x + y) \leq N(x) + N(y)$
- $N(ax) = |a|N(x)$  for  $a$  scalar
- if  $N(x) = 0$ , then  $x = 0$

For  $x \in V$ , the most commonly used norms are summed up in the table below:

| Norm                 | Notation       | Definition                                      | Use case             |
|----------------------|----------------|---|----------------------|
| Manhattan, $L^1$     | $\ x\ _1$      | $\sum_{i=1}^n  x_i $                            | LASSO regularization |
| Euclidean, $L^2$     | $\ x\ _2$      | $\sqrt{\sum_{i=1}^n x_i^2}$                     | Ridge regularization |
| $p$ -norm, $L^p$     | $\ x\ _p$      | $\left(\sum_{i=1}^n x_i^p\right)^{\frac{1}{p}}$ | Hölder inequality    |
| Infinity, $L^\infty$ | $\ x\ _\infty$ | $\max_i  x_i $                                  | Uniform convergence  |

**Linearly dependence** – A set of vectors is said to be linearly dependent if one of the vectors in the set can be defined as a linear combination of the others.

*Remark:* if no vector can be written this way, then the vectors are said to be linearly independent.

**Matrix rank** – The rank of a given matrix  $A$  is noted  $\text{rank}(A)$  and is the dimension of the vector space generated by its columns. This is equivalent to the maximum number of linearly independent columns of  $A$ .

**Positive semi-definite matrix** – A matrix  $A \in \mathbb{R}^{n \times n}$  is positive semi-definite (PSD) and is noted  $A \succeq 0$  if we have:

$$A = A^T \quad \text{and} \quad \forall x \in \mathbb{R}^n, \quad x^T A x \geq 0$$

*Remark:* similarly, a matrix  $A$  is said to be positive definite, and is noted  $A \succ 0$ , if it is a PSD matrix which satisfies for all non-zero vector  $x$ ,  $x^T A x > 0$ .

**Eigenvalue, eigenvector** – Given a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of  $A$  if there exists a vector  $z \in \mathbb{R}^n \setminus \{0\}$ , called eigenvector, such that we have:

$$Az = \lambda z$$

**Spectral theorem** – Let  $A \in \mathbb{R}^{n \times n}$ . If  $A$  is symmetric, then  $A$  is diagonalizable by a real orthogonal matrix  $U \in \mathbb{R}^{n \times n}$ . By noting  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ , we have:

$$\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$$

**Singular-value decomposition** – For a given matrix  $A$  of dimensions  $m \times n$ , the singular-value decomposition (SVD) is a factorization technique that guarantees the existence of  $U$   $m \times m$  unitary,  $\Sigma$   $m \times n$  diagonal and  $V$   $n \times n$  unitary matrices, such that:

$$A = U \Sigma V^T$$

## Matrix calculus

**Gradient** – Let  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  be a function and  $A \in \mathbb{R}^{m \times n}$  be a matrix. The gradient of  $f$  with respect to  $A$  is a  $m \times n$  matrix, noted  $\nabla_A f(A)$ , such that:

$$\left( \nabla_A f(A) \right)_{i,j} = \frac{\partial f(A)}{\partial A_{i,j}}$$

*Remark:* the gradient of  $f$  is only defined when  $f$  is a function that returns a scalar.

**Hessian** – Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function and  $x \in \mathbb{R}^n$  be a vector. The hessian of  $f$  with respect to  $x$  is a  $n \times n$  symmetric matrix, noted  $\nabla_x^2 f(x)$ , such that:

$$\left( \nabla_x^2 f(x) \right)_{i,j} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

*Remark:* the hessian of  $f$  is only defined when  $f$  is a function that returns a scalar.

**Gradient operations** – For matrices  $A, B, C$ , the following gradient properties are worth having in mind:

$$\nabla_A \text{tr}(AB) = B^T$$

$$\nabla_A f(A) = (\nabla_A f(A))^T$$

$$\nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T$$

$$\nabla_A |A| = |A|(A^{-1})^T$$