

IntroSlides with focus on statistics

1

Basic steps in workflow

1. *Define environment*
2. *Import*
3. *Transform*
4. *Explore (general/outlier/distribution) (go back to 3?)*
5. *Classify scale level / distribution (based on 3/4)*
6. *Describe*
7. *Test / Model (may include step 6)*
8. *Report*

2

Define environment

- *Activate packages to use: library() / pacman::p_load()*
- *ggplot theme: theme_set() / theme_update()*
- *flextable settings: set_flextable_defaults()*
- *knitr::opts_chunk\$set()*

```
1 if(!requireNamespace("pacman")){install.packages("pacman")}
2 pacman::p_load(conflicted, tidyverse, wrappedtools, readxl, car, flextable,
3               ggbeeswarm, ggsignif, ggribes, patchwork, ggrepel, easystats)
4
5 conflicts_prefer(dplyr::filter, dplyr::select)
6 theme_set(theme_light(base_size = 20))
7 gdtools::register_gfont('Roboto') # Mono'
```

[1] TRUE

```
1 set_flextable_defaults(
2   theme_fun = theme_zebra, font.size = 18, font.family = 'Roboto',
3   table.layout = 'autofit',
4   padding.bottom = .2, padding.top = .2, padding.left = 2, padding.right = 2)
5
6 knitr::opts_chunk$set(message = FALSE, warning = FALSE, comment = NA, echo = T
```

3

Import

- *read_xlsx() / read_csv() / read_csv2()*
- options related to separators, number formats, ranges etc.
- *rename() / rename_with()*

```
1 rawdata <- read_excel('Data/DOC-20230130-WA0000_.xlsx',
2                       sheet = 1, col_names = TRUE)
```

4

Glimpse at data: Find the problems?

```
1 head(rawdata,n = 15) |> flextable()|>
2   theme_zebra(even_body = 'aquamarine',odd_body = 'antiquewhite')
```

CODE OF CUP	CODE OF SAMPLE	WEIGHT OF EMPTY ALUMINUM(wt)	WEIGHT OF ALMINIUM CUP + SAMPLE (Wt + s)	WEIGHT OF ALUMINIUM CAP + SAMPLE AFTER DRYING (Wt-AL +s+d)	weight of sample before drying (Wts)	weight of sample after drying (Wts+d)	MOISTURE CONTENT (%)
69	D	4.1974	9.3865	4.7000	5.1891	4.6865	90.31431
	D	4.1964	9.2734	4.4670	5.0770	4.8064	94.67008
A	D	4.2108	9.2653	4.6670	5.0545	4.5983	90.97438
114	D	4.2134	9.3146	4.6345	5.1012	4.6801	91.74508
M1	D	4.1856	9.3147	4.6171	5.1291	4.6976	91.58722
a/17	D	4.2090	9.3204	4.5661	5.1114	4.7543	93.01366
8	D	4.1894	9.2661	4.5778	5.0767	4.6883	92.34936
33	D	4.1968	9.2880	4.6057	5.0912	4.6823	91.96849
M	D	4.1535	9.2872	4.6350	5.1337	4.6522	90.62080
E/18/1	D	4.2534	9.2476	4.7403	4.9942	4.5073	90.25069
24/A2	D	4.2066	8.3463	4.5849	4.1397	3.7614	90.86166
13	A	4.1554	9.2384	4.7402	5.0830	4.4982	88.49498
Xp	A	4.1893	9.2495	4.7381	5.0602	4.5114	89.15458
2p/029	A	4.0654	9.2173	4.6940	5.1519	4.5233	87.79868
15	A	4.0641	9.2032	4.8124	5.1391	4.3908	85.43908

5

Rename

```
1 colnames(rawdata)
```

```
[1] "CODE OF CUP"
[2] "CODE OF SAMPLE"
[3] "WEIGHT OF EMPTY ALUMINUM(wt)"
[4] "WEIGHT OF ALMINIUM CUP + SAMPLE (Wt + s)"
[5] "WEIGHT OF ALUMINIUM CAP + SAMPLE AFTER DRYING (Wt-AL +s+d)"
[6] "weight of sample before drying (Wts)"
[7] "weight of sample after drying (Wts+d)"
[8] "MOISTURE CONTENT (%)"
```

```
1 rawdata <- rawdata |>
2   rename(Region=`CODE OF SAMPLE`) |>
3   rename_with(.fn = ~str_replace_all(
4     `
5     c("GTH"="GHT", 'AL.+UM'= 'Cup',
6       'C[UA]P' = 'Cup', '\\(\\w+.*\\)'=''),
7     'Cup Cup'='Cup', ' '= ' ') |>
8     str_to_title() |> str_trim())
9   cn()
```

```
[1] "Code Of Cup"           "Region"
[3] "Weight Of Empty Cup"   "Weight Of Cup + Sample"
[5] "Weight Of Cup + Sample After Drying" "Weight Of Sample Before Drying"
[7] "Weight Of Sample After Drying"      "Moisture Content (%)"
```

6

Transform

- *Change or create columns with `mutate()` / `mutate(across())`*
- *e.g. for log-transformation, creation of factors, text recoding*

```
1 rawdata <- rawdata |>
2   mutate(
3     `Code Of Cup` = case_when(
4       is.na(`Code Of Cup`) ~ paste("sample", row_number()),
5       .default = `Code Of Cup`),
6     `Weight Of Sample After Drying` = `Weight Of Cup + Sample After Drying` -
7     `Weight Of Empty Cup`,
8     `Dry Content (%)` = `Weight Of Sample After Drying` * 100 /
9     `Weight Of Sample Before Drying`,
10    `Moisture Content (%)` = 100 - `Dry Content (%)`)
```

Code Of Cup	Region	Weight Of Empty Cup	Weight Of Cup + Sample	Weight Of Cup + Sample After Drying	Weight Of Sample Before Drying	Weight Of Sample After Drying	Moisture Content (%)	Dry Content (%)
69	D	4.1974	9.3865	4.7000	5.1891	0.5026	90.31431	9.685687
sample 2	D	4.1964	9.2734	4.4670	5.0770	0.2706	94.67008	5.329919
A	D	4.2108	9.2653	4.6670	5.0545	0.4562	90.97438	9.025621
114	D	4.2134	9.3146	4.6345	5.1012	0.4211	91.74508	8.254920
M1	D	4.1856	9.3147	4.6171	5.1291	0.4315	91.58722	8.412782

7

Explore / group variables

Explore (general/outlier/distribution)

- *`ggplot()` + `geom_boxplot()` / `geom_beeswarm()` / `geom_density()`*
- *`ks.test()` / `ksnormal()` / `shapiro.test()`*

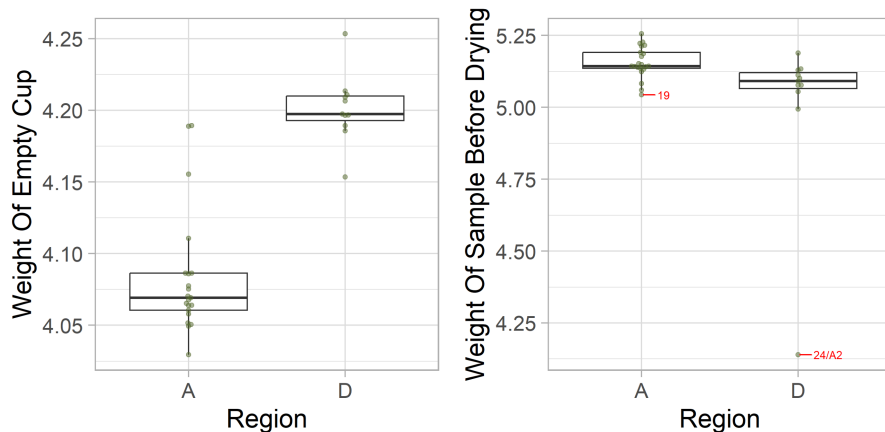
Classify scale level / distribution

- *`gaussvars` / `ordvars` / `factvars`, possibly more...*
- *Store variables accordingly, e.g. `ColSeeker()`*

8

Explore: Outlier

```
1 p1 <- ggplot(data = rawdata, aes(x = `Region`, y = `Weight Of Empty Cup`))+
2   geom_boxplot(outlier.alpha = 0) + #hide outliers, beeswarm will plot them
3   geom_beeswarm(alpha=.5, color="darkolivegreen")
4 p2 <- ggplot(data = rawdata, aes(x = `Region`, y = `Weight Of Sample Before Dr
5   geom_boxplot(outlier.alpha = 0) +
6   geom_beeswarm(alpha=.5, color="darkolivegreen")+
7   geom_text_repel(data=. %>% group_by(Region) %>%
8     filter(`Weight Of Sample Before Drying` %in%
9       boxplot.stats(`Weight Of Sample Before Drying`, coef=1.5
10         aes(label=`Code Of Cup`, y=`Weight Of Sample Before Drying`),
11         nudge_x=0.1, colour="red", size=3, hjust=0)
12 p1|p2
```



9

Handle outliers?

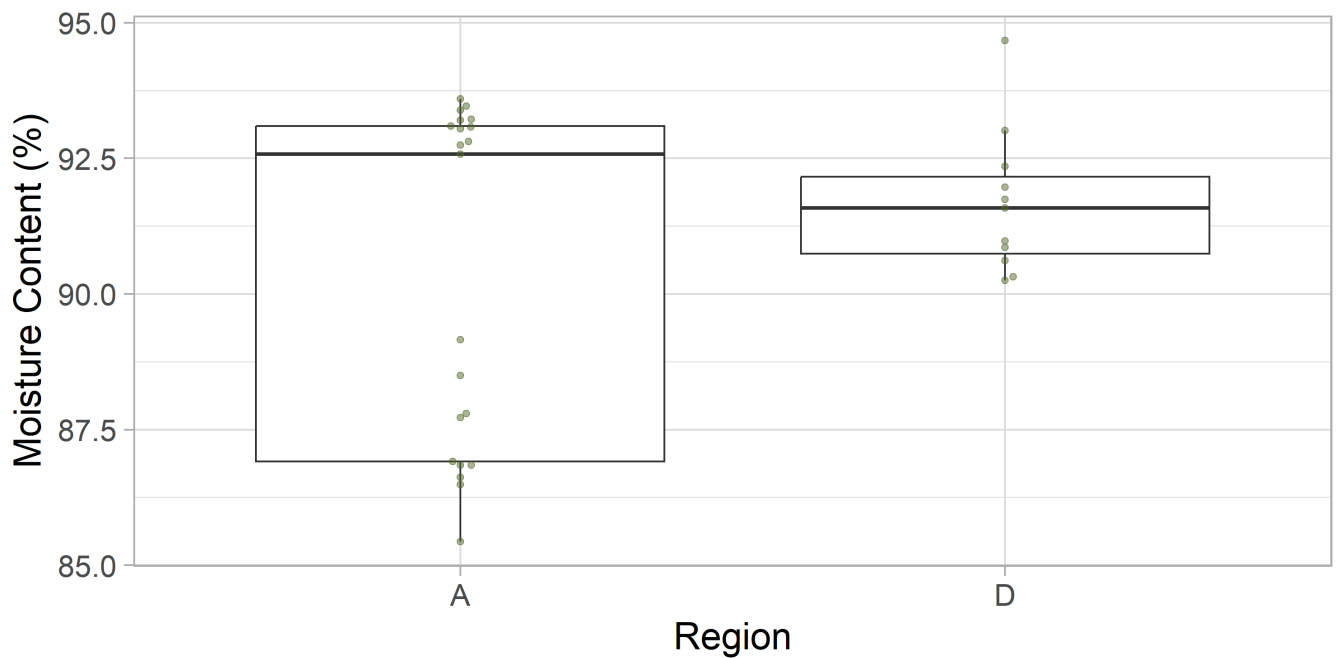
Removal is the worst strategy possible, correct errors, think about distributions, winsorize...



10

Explore: Unexpected

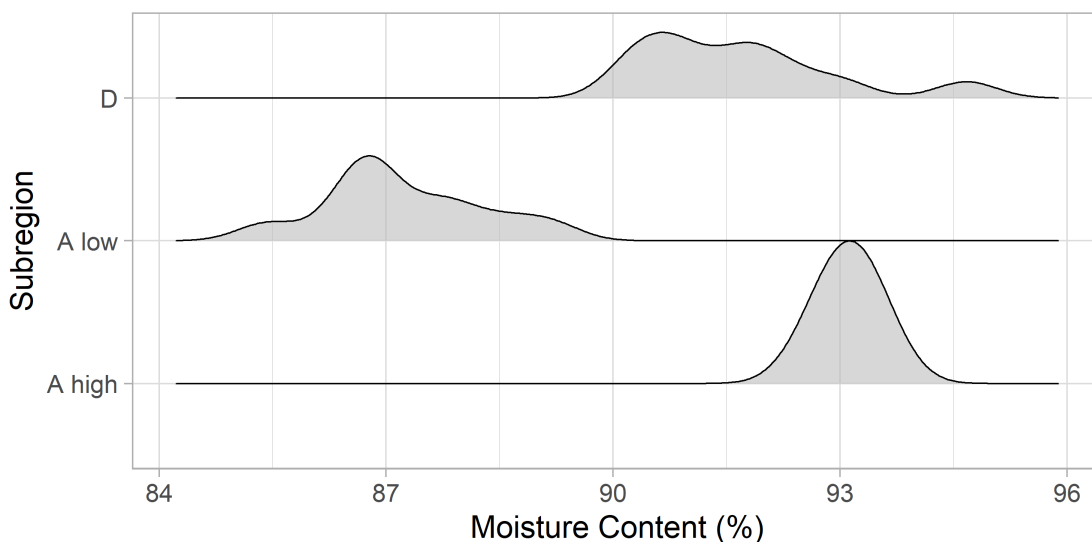
```
1 ggplot(data = rawdata,  
2       aes(x = `Region`,  
3           y = `Moisture Content (%)`))+  
4   geom_boxplot(outlier.alpha = 0) +  
5   geom_beeswarm(alpha=.5, color="darkolivegreen")
```



11

Transform Subregions?

```
1 rawdata <- mutate(rawdata, Subregion = case_when(  
2   `Region`=='D' ~ 'D',  
3   `Region`=='A' & `Moisture Content (%)` > 90 ~ 'A high',  
4   `Region`=='A' & `Moisture Content (%)` <= 90 ~ 'A low') |>  
5   factor())  
6 # Test for Region A is redundant here, but more verbose.  
7  
8 ggplot(data = rawdata, aes(x = `Moisture Content (%)`, y=Subregion))+  
9   geom_density_ridges(alpha=.5, scale=1)
```



12

Explore: Normal distribution 1

- Gaussian Normal distribution is required for many statistical procedures
- Common tests are graphical exploration, Shapiro-Wilk-test and Kolmogorov-Smirnov-test

```
1 p_normal <-  
2   shapiro.test(x = rawdata$`Moisture Content (%)`)  
3 p_normal
```

Shapiro-Wilk normality test

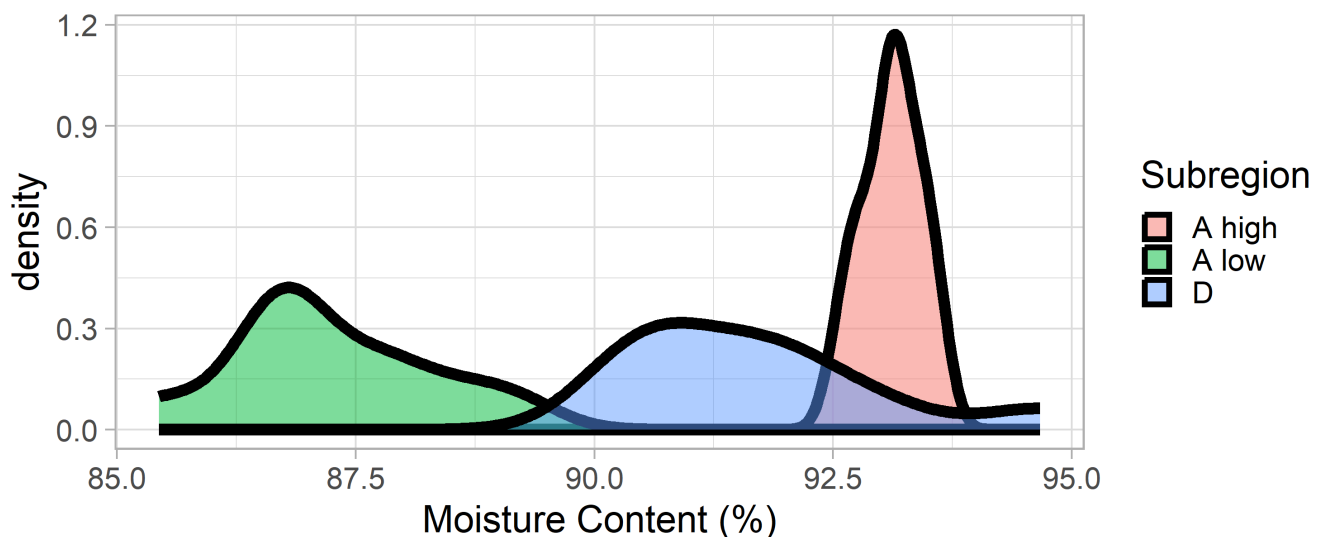
```
data: rawdata$`Moisture Content (%)`  
W = 0.89133, p-value = 0.003752
```

13

```
1 ggplot(data = rawdata, aes(x = `Moisture Content (%)`, fill = `Subregion`)) +  
2   geom_density(linewidth = 3, alpha = .5) +  
3   labs(title = paste('p (Shapiro) global',  
4     formatP(pIn = p_normal$p.value, pretext = T)),  
5     subtitle = rawdata |> group_by(Subregion) |>  
6     summarize(pNormal = shapiro.test(`Moisture Content (%)`)$p.value |> formatP())  
7     unite('p(Normal)', sep = ': p=') |> pull(1) |> paste(collapse = '; '))
```

p (Shapiro) global = 0.004

A high: p=0.923; A low: p=0.756; D: p=0.207



14

Explore: Normal distribution 2

```
1 rawdata |>
2   group_by(`Subregion`) |>
3   summarize(across(.cols = where(is.numeric),
4                     .fns = ~ksnormal(.x) |> #computing p-value
5                               formatP(mark = T))) |> #formatting p-value
6   pivot_longer(cols = -1, names_to = 'Variable',
7                values_to = 'pKS') |> #intermediate, all p-values in 1 column
8   pivot_wider(names_from = `Subregion`,
9                values_from = pKS) #spreading across subregions
```

A tibble: 7 × 4

Variable <chr>	`A high` <chr>	`A low` <chr>	D <chr>
1 Weight Of Empty Cup	0.989 n.s.	0.444 n.s.	0.697 n.s.
2 Weight Of Cup + Sample	0.785 n.s.	0.980 n.s.	0.019 *
3 Weight Of Cup + Sample After Drying	0.900 n.s.	0.710 n.s.	0.969 n.s.
4 Weight Of Sample Before Drying	0.196 n.s.	0.999 n.s.	0.072 +
5 Weight Of Sample After Drying	0.976 n.s.	0.555 n.s.	1.000 n.s.
6 Moisture Content (%)	0.975 n.s.	0.733 n.s.	0.954 n.s.
7 Dry Content (%)	0.975 n.s.	0.733 n.s.	0.954 n.s.

15

Explore: Group variables by type/distribution

Scale level determines what statistics are appropriate

Typical scale levels are

- nominal/categorical/factorial/qualitative: just different groups (species, eye color, genotype, treatment)
- ordered categories: few groups with inherent order (quality bad<medium<good, pain between 0 and 10)
- ordinal measures: many different values, natural order, no distribution assumption (satisfaction on a scale from 0 to 100)
- measures following a Normal distribution
- possibly measures from other known distributions (beta, log-normal, poisson...), often treated as ordinal

16

Make type decision obvious/reproducible

```
1 gaussvars <- ColSeeker(data=rawdata,namepattern = c('Weight','Content'))
2 gaussvars
```

\$index

```
[1] 3 4 5 6 7 8 9
```

\$names

```
[1] "Weight Of Empty Cup" "Weight Of Cup + Sample"
[3] "Weight Of Cup + Sample After Drying" "Weight Of Sample Before Drying"
[5] "Weight Of Sample After Drying" "Moisture Content (%)"
[7] "Dry Content (%)"
```

\$bticked

```
[1] "`Weight Of Empty Cup`"
[2] "`Weight Of Cup + Sample`"
[3] "`Weight Of Cup + Sample After Drying`"
[4] "`Weight Of Sample Before Drying`"
[5] "`Weight Of Sample After Drying`"
[6] "`Moisture Content (%)`"
[7] "`Dry Content (%)`"
```

```
1 ordvars <- ColSeeker(namepattern='Weight.+Sample', exclude = 'After')
2 ordvars$names
```

```
[1] "Weight Of Cup + Sample" "Weight Of Sample Before Drying"
```

```
1 factvars <- ColSeeker(namepattern='region',casesensitive = FALSE)
2 factvars$bticked
```

```
[1] "`Region`" "`Subregion`"
```

17

Model

Describe

- *mean()* / *sd()* / *meansd()*
- *median()* / *quantile()* / *median_quart()*
- *table()* / *prop.table()* / *cat_desc_stats()*

Test

- *t.test()* / *lm()* / *[Aa]nova()* / *compare2numvars()*
- *wilcox.test()*
- *fisher.test()* / *glm(family=binomial)*

18

Model: Describe

Sample size n: per variable, if there are NAs

Mean: central tendency, the expected *typical* value

$$\frac{\sum x}{n}$$

Variance: measure for variability/heterogeneity of data

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

Standard deviation SD: the *typical* weighted deviation from the mean

$$\sqrt{Var}$$

Standard error of the mean SEM: how reliable is the mean *estimate*, what would be the expected SD of means from repeated experiments?

$$\frac{SD}{\sqrt{n}}$$

Median: Split between lower/upper 50% of data

Quartiles: Split at 25%/50%/75% of data (more general: **Quantiles**, e.g. **Percentiles**), used in boxplot

various computational approaches

```

1 desc_gauss <- rawdata |>
2   summarize(across(.cols = gaussvars$names,
3                     .fns = meansd))
4 desc_gauss

```

```

# A tibble: 1 × 7
  `Weight Of Empty Cup` `Weight Of Cup + Sample` Weight Of Cup + Sample After ...1
  <chr>                <chr>                <chr>
1 4.1 ± 0.1            9.2 ± 0.2            4.6 ± 0.1
# i abbreviated name: 1 `Weight Of Cup + Sample After Drying`
# i 4 more variables: `Weight Of Sample Before Drying` <chr>,
#   `Weight Of Sample After Drying` <chr>, `Moisture Content (%)` <chr>,
#   `Dry Content (%)` <chr>

```

```

1 desc_ord <- rawdata |>
2   summarize(across(ordvars$names, .fns=~median_quart(.x, roundDig = 3))) |>
3   pivot_longer(everything(),
4                 names_to = 'Measure', values_to = 'Median[1Q/3Q]')
5 desc_ord

```

```

# A tibble: 2 × 2
  Measure                `Median[1Q/3Q]`
  <chr>                <chr>
1 Weight Of Cup + Sample 9.25 (9.22/9.29)
2 Weight Of Sample Before Drying 5.14 (5.09/5.18)

```

21

Descriptive Stats should match distribution and data



22

Model: Test

Tests require hypotheses



23

Null hypothesis ?

- Working hypothesis: This is what you expect!
E.g. treatment is lowering blood pressure more than placebo, transgenic animals become obese, bio reactor A is more efficient than B, concentration of substance is correlated with speed of reaction ...
- Null hypothesis: This is what you test!
No difference / relation, BP under therapy = BP under placebo

24

4 possibilities:

- Null hypothesis correct, test false positive (case A): alpha-error
- Null hypothesis correct, test correct negative (case B)
- Null hypothesis false, test false negative (case C): beta-error
- Null hypothesis false, test correct positive (case D)

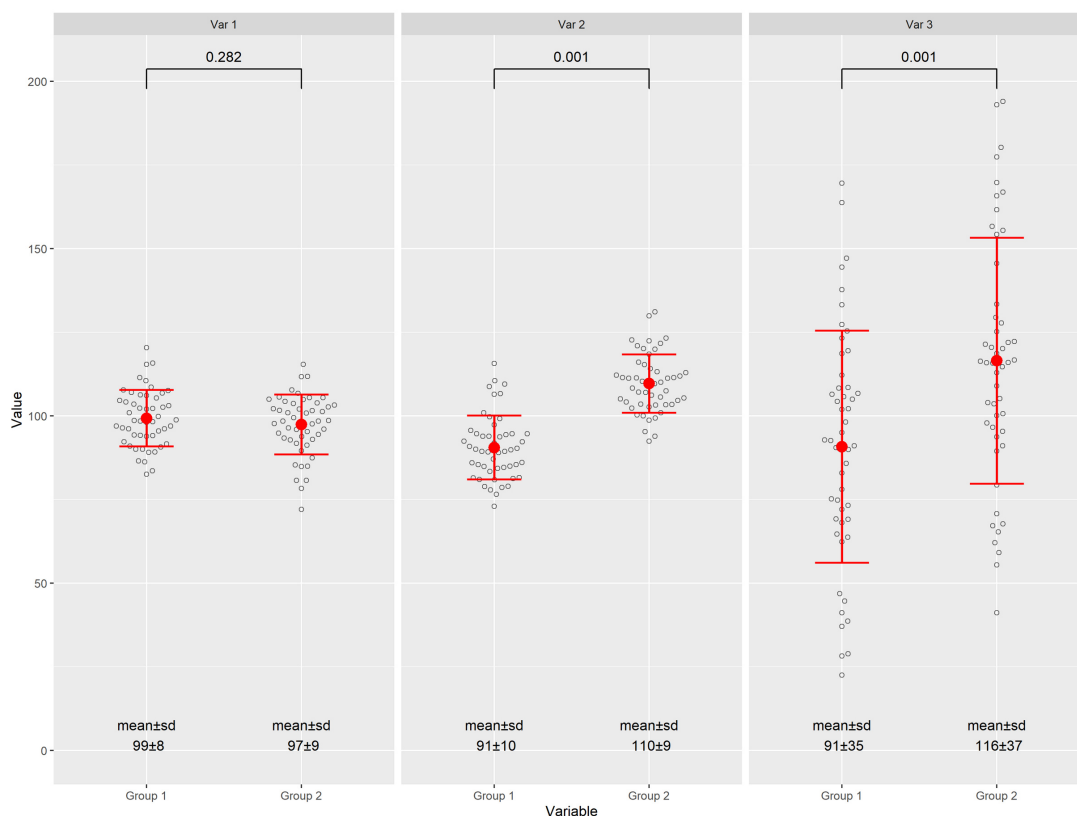
Significance: NOT probability of case A, but probability of your data given the NULL hypothesis, calculated from your data, conventionally <0.05

Power: Probability of case D, *estimated* based on assumptions about effects and sample size, *calculation* would require knowledge of true difference, conventionally set at 0.80

25

Test functions

t-test / Wilcoxon-test (aka Mann-Whitney U-test)



26

t-test

- Assumptions: Continuous data with Normal distribution
- 1 or 2 (independent or dependent) samples with/without equal variances
- how big is the mean difference relative to uncertainty?
$$t = (\text{mean}_1 - \text{mean}_2) / \text{SEM}$$
- t follows a t-distribution, allows estimation of probability of t under the NULL hypothesis

Wilcoxon-test

- nonparametric, no distribution is assumed
- based on rank-transformed data
- insensitive to extreme values

27

Test examples: *single variables*

```
1 #t-Test with test for equal variances
2 t_out <- t.test(formula='Moisture Content (%)'~'Region', data=rawdata,
3                 var.equal=var.test(
4                     formula='Moisture Content (%)'~'Region',
5                     data=rawdata)$p.value>.05)
6 t_out
```

Welch Two Sample t-test

```
data: Moisture Content (%) by Region
t = -1.7274, df = 29.239, p-value = 0.09465
alternative hypothesis: true difference in means between group A and group D is not equal to 0
95 percent confidence interval:
 -2.9624835  0.2490486
sample estimates:
mean in group A mean in group D
 90.31199      91.66870
```

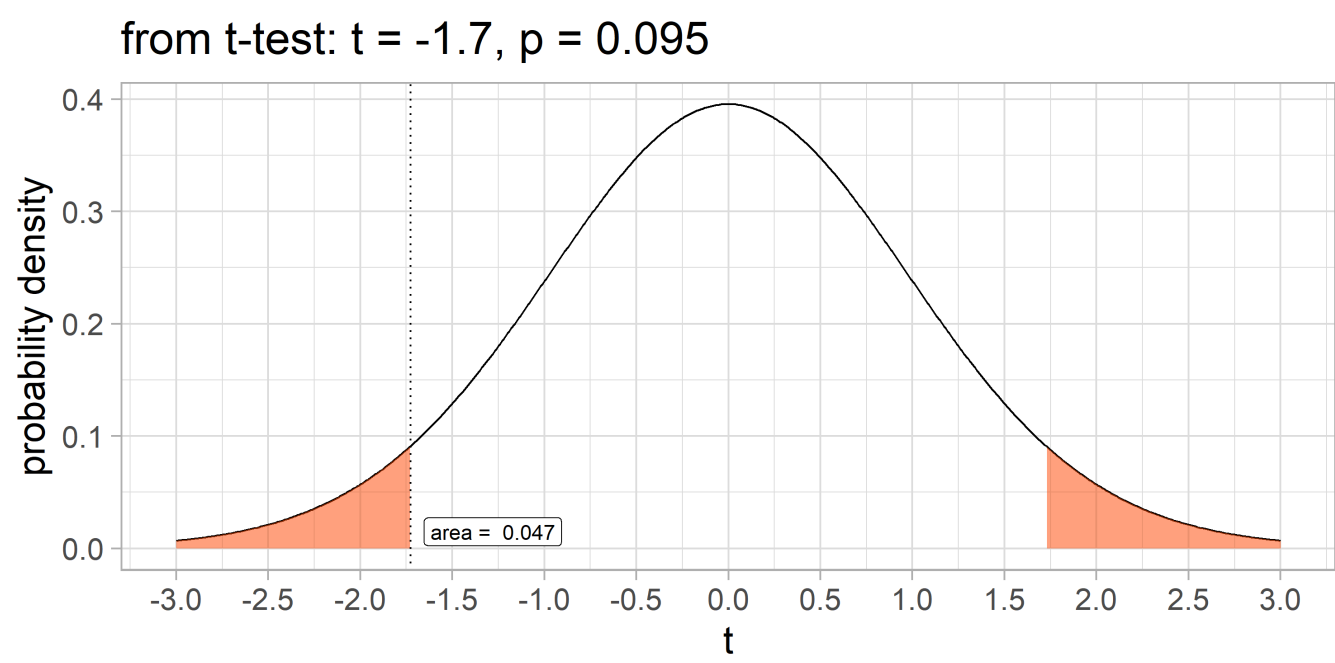
```
1 #Wilcoxon-Test
2 wilcox.test('Moisture Content (%)'~'Region',
3             data = rawdata)
```

Wilcoxon rank sum exact test

```
data: Moisture Content (%) by Region
W = 107, p-value = 0.7547
alternative hypothesis: true location shift is not equal to 0
```

28

From t to p



29

Model: Test 2 / multiple variables

```
1 test_gauss <- compare2numvars(data = rawdata,
2                               dep_vars = gaussvars$names,
3                               indep_var = 'Region',
4                               gaussian = TRUE,
5                               round_p = 5)
6 test_gauss |> flextable()|>
7   theme_zebra(even_body = 'aquamarine',odd_body = 'antiquewhite')
```

Variable	desc_all	Region A	Region D	p
Weight Of Empty Cup	4.1 ± 0.1	4.1 ± 0.0	4.2 ± 0.0	0.00001
Weight Of Cup + Sample	9.2 ± 0.2	9.2 ± 0.0	9.2 ± 0.3	0.74227
Weight Of Cup + Sample After Drying	4.6 ± 0.1	4.6 ± 0.2	4.6 ± 0.1	0.42213
Weight Of Sample Before Drying	5.1 ± 0.2	5.2 ± 0.1	5.0 ± 0.3	0.12830
Weight Of Sample After Drying	0.47 ± 0.14	0.50 ± 0.16	0.42 ± 0.07	0.04937
Moisture Content (%)	91 ± 3	90 ± 3	92 ± 1	0.09465
Dry Content (%)	9.2 ± 2.7	9.7 ± 3.1	8.3 ± 1.3	0.09465

30

```

1 test_ord <- compare2numvars(data = rawdata,
2                             dep_vars = ordvars$names,
3                             indep_var = 'Region',
4                             gaussian = FALSE, round_desc = 3)
5 test_ord |> flextable() |>
6   theme_zebra(even_body = 'aquamarine', odd_body = 'antiquewhite')

```

Variable	desc_all	Region A	Region D	p
Weight Of Cup + Sample	9.25 (9.22/9.29)	9.24 (9.22/9.25)	9.29 (9.27/9.31)	0.003
Weight Of Sample Before Drying	5.14 (5.09/5.18)	5.14 (5.13/5.20)	5.09 (5.06/5.13)	0.001

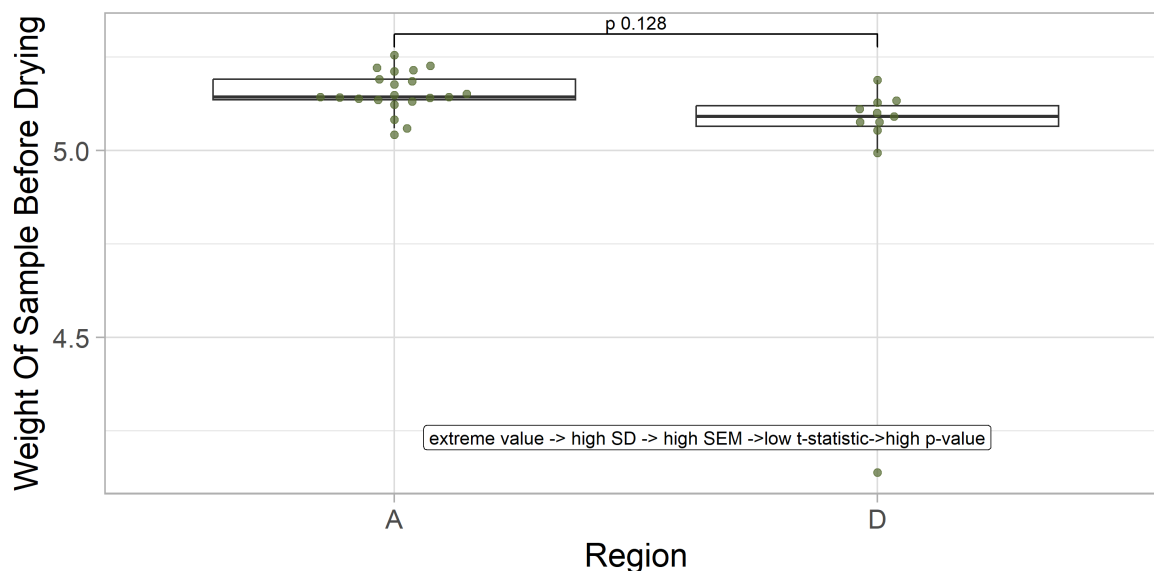
31

Show results

```

1 ggplot(rawdata, aes(x = `Region`, y = `Weight Of Sample Before Drying`))+
2   geom_boxplot(outlier.alpha = 0)+
3   geom_beeswarm(alpha=.7, size=2, cex = 2, color="darkolivegreen")+
4   annotate(geom = 'label', x=2, y=4.2,
5           label='extreme value -> high SD -> high SEM -> low t-statistic->high
6           hjust=0.8, vjust=0)+
7   geom_signif(comparisons = list(c(1,2)),
8             annotations = paste('p', formatP(t_out$p.value)))

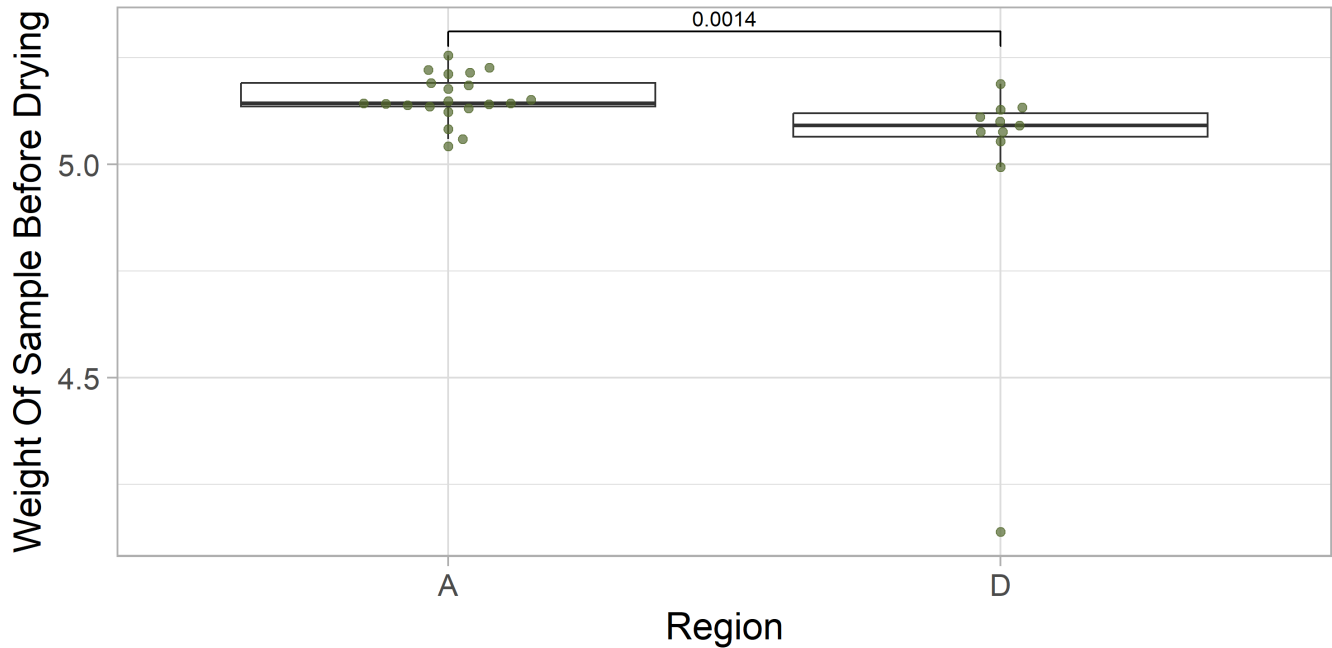
```



32

Re-thinking test decision?

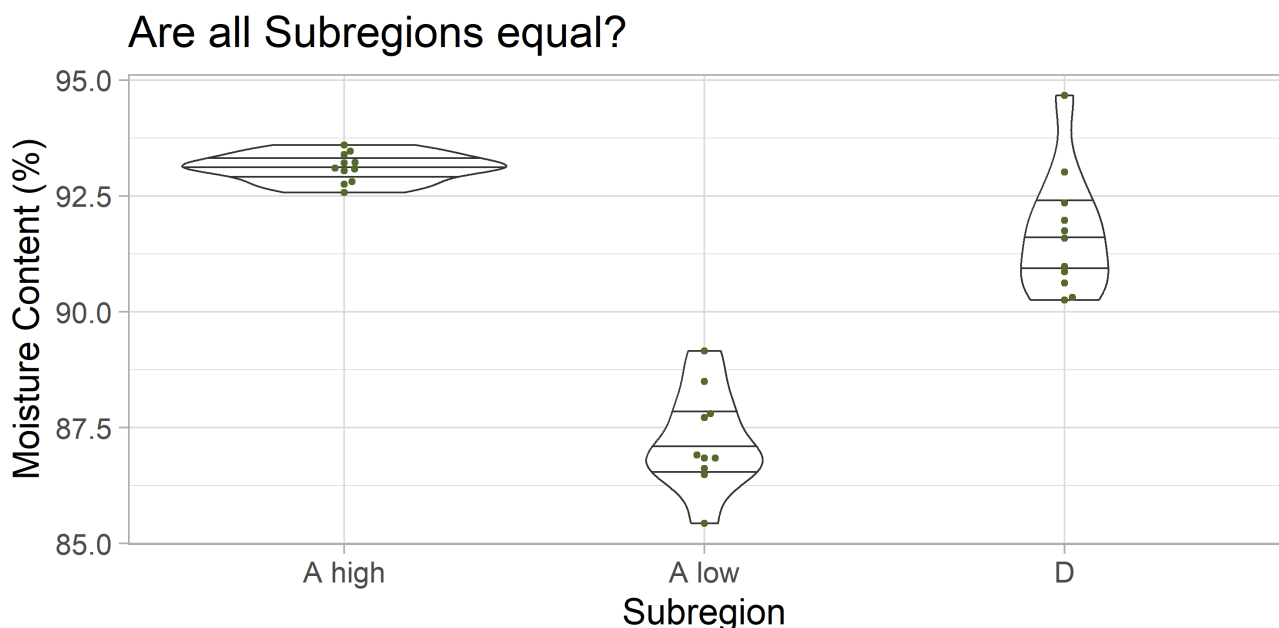
```
1 ggplot(rawdata, aes(x = `Region`, y = `Weight Of Sample Before Drying`))+  
2   geom_boxplot(outlier.alpha = 0)+  
3   geom_beeswarm(alpha=.7, size=2, cex = 2, color="darkolivegreen")+  
4   geom_signif(comparisons = list(c(1,2)), test = wilcox.test)
```



33

Model: linear models 1 / univariable

```
1 plottmp <- ggplot(rawdata, aes(Subregion, `Moisture Content (%)`))+  
2   geom_violin(draw_quantiles = c(.25, .5, .75))+  
3   geom_beeswarm(color="darkolivegreen")+  
4  
5   ggtitle('Are all Subregions equal?')  
6   print(plottmp)
```



34

ANOVA: build model

```
1 rawdata |> group_by(Subregion) |>
2   summarize(MeanMoisture=mean(`Moisture Content (%)`) |> roundR(4)) |>
3   pivot_wider(names_from = Subregion, values_from = MeanMoisture) |>
4   rename_with(~paste('Mean moisture %\n', .x)) |> flextable() |>
5   theme_zebra(even_body = 'aquamarine', odd_body = 'antiquewhite')
```

Mean moisture % A high	Mean moisture % A low	Mean moisture % D
93.11	87.23	91.67

```
1 lm1<- lm(`Moisture Content (%)`~Subregion, data=rawdata)
2 lm1
```

Call:

```
lm(formula = `Moisture Content (%)` ~ Subregion, data = rawdata)
```

Coefficients:

(Intercept)	SubregionA low	SubregionD
93.112	-5.879	-1.443

35

ANOVA: get p-values

```
1 anova(lm1) |> broom::tidy() |> flextable() |>
2   theme_zebra(even_body = 'aquamarine', odd_body = 'antiquewhite')
```

term	df	sumsq	meansq	statistic	p.value
Subregion	2	194.34239	97.1711969	97.76477	0.00000000000001292127
Residuals	29	28.82393	0.9939285		

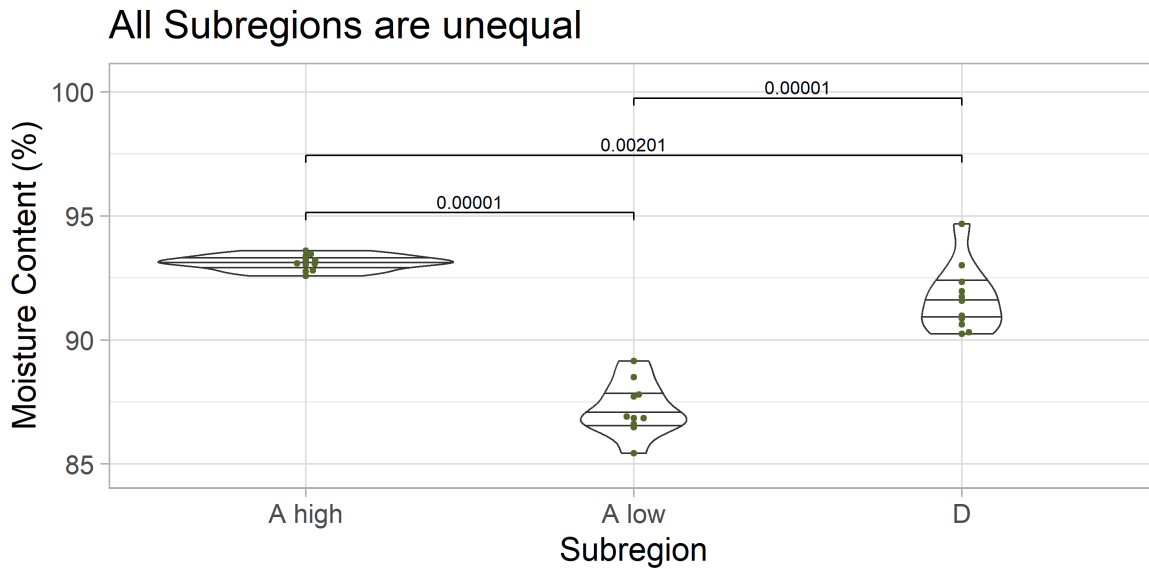
```
1 #post-hoc
2 (posthoc_out <- pairwise.t.test(x = rawdata$`Moisture Content (%)`,
3   g = rawdata$Subregion,
4   p.adjust.method = 'fdr'))$p.value |>
5   formatP(ndigits = 5)
```

	A high	A low
A low	"0.00001"	"NA"
D	"0.00201"	"0.00001"

36

Visualize ANOVA

```
1 ggplot(rawdata,aes(Subregion,`Moisture Content (%)`))+
2   geom_violin(draw_quantiles = c(.25,.5,.75))+
3   geom_beeswarm(color="darkolivegreen")+
4   geom_signif(comparisons = list(c(1,2),c(1,3),c(2,3)),
5               annotations = c(posthoc_out[,1], posthoc_out[2,2]),
6               step_increase = .25)+
7   scale_y_continuous(expand = expansion(mult = .1))+
8   ggtitle('All Subregions are unequal')
```



37

Analyze more than 1 outcome

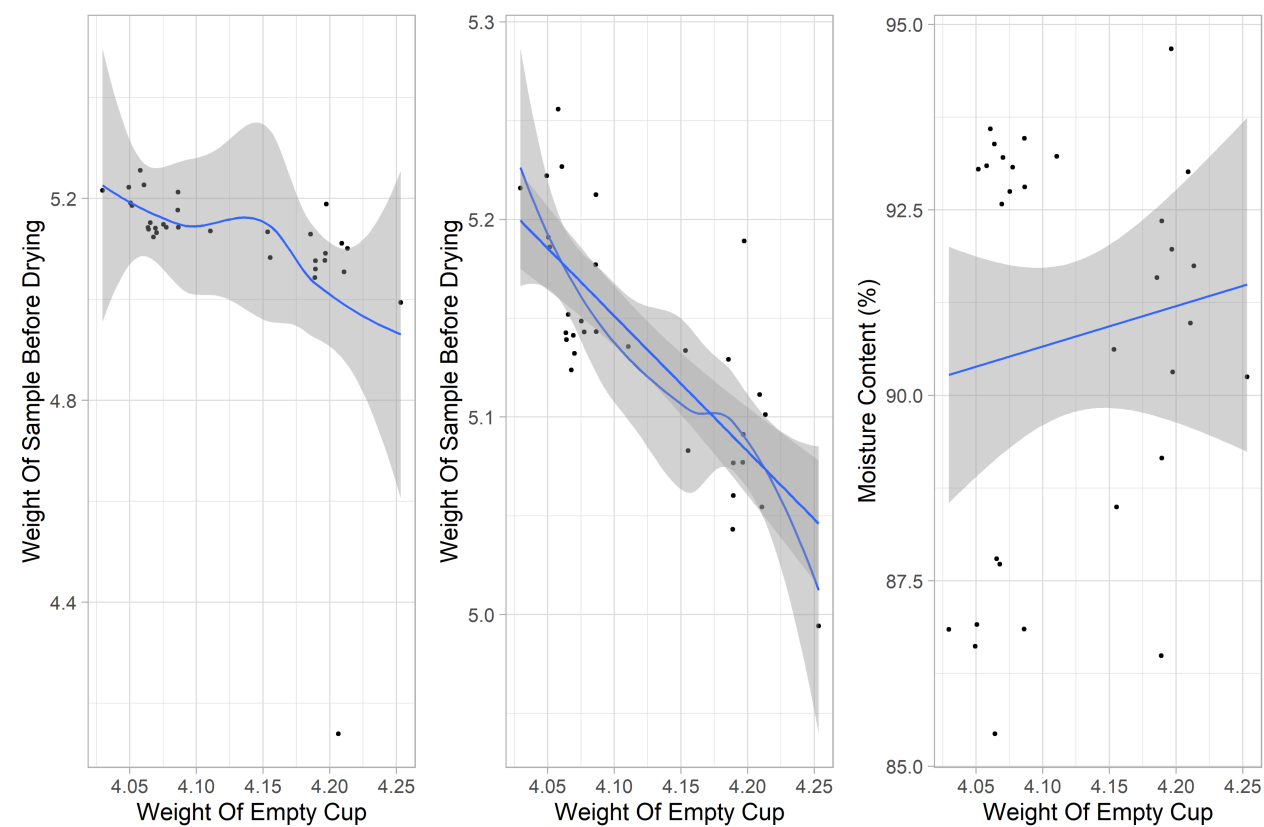
```
1 test_out <- compare_n_numvars(.data=rawdata,
2                               dep_vars=gaussvars$names,
3                               indep_var='Subregion',
4                               gaussian=TRUE)
5 test_out$results |>
6   select(Variable,contains("fn"),multivar_p) |>
7   rename_with(~str_remove(.x," fn")) |>
8   flextable() |>
9   theme_zebra(even_body = 'aquamarine',odd_body = 'antiquewhite') |>
10  add_footer_lines(
11    values='b indicates difference from group 2, c indicates difference from g
```

Variable	Subregion A high	Subregion A low	Subregion D	multivar_p
Weight Of Empty Cup	4.1 ± 0.0 c	4.1 ± 0.1 c	4.2 ± 0.0	0.001
Weight Of Cup + Sample	9.2 ± 0.0	9.2 ± 0.0	9.2 ± 0.3	0.894
Weight Of Cup + Sample After Drying	4.4 ± 0.0 bc	4.8 ± 0.1 c	4.6 ± 0.1	0.001
Weight Of Sample Before Drying	5.2 ± 0.0	5.1 ± 0.1	5.0 ± 0.3	0.095
Weight Of Sample After Drying	0.36 ± 0.02 bc	0.66 ± 0.06 c	0.42 ± 0.07	0.001
Moisture Content (%)	93 ± 0 bc	87 ± 1 c	92 ± 1	0.001
Dry Content (%)	7 ± 0 bc	13 ± 1 c	8 ± 1	0.001

b indicates difference from group 2, c indicates difference from group 3

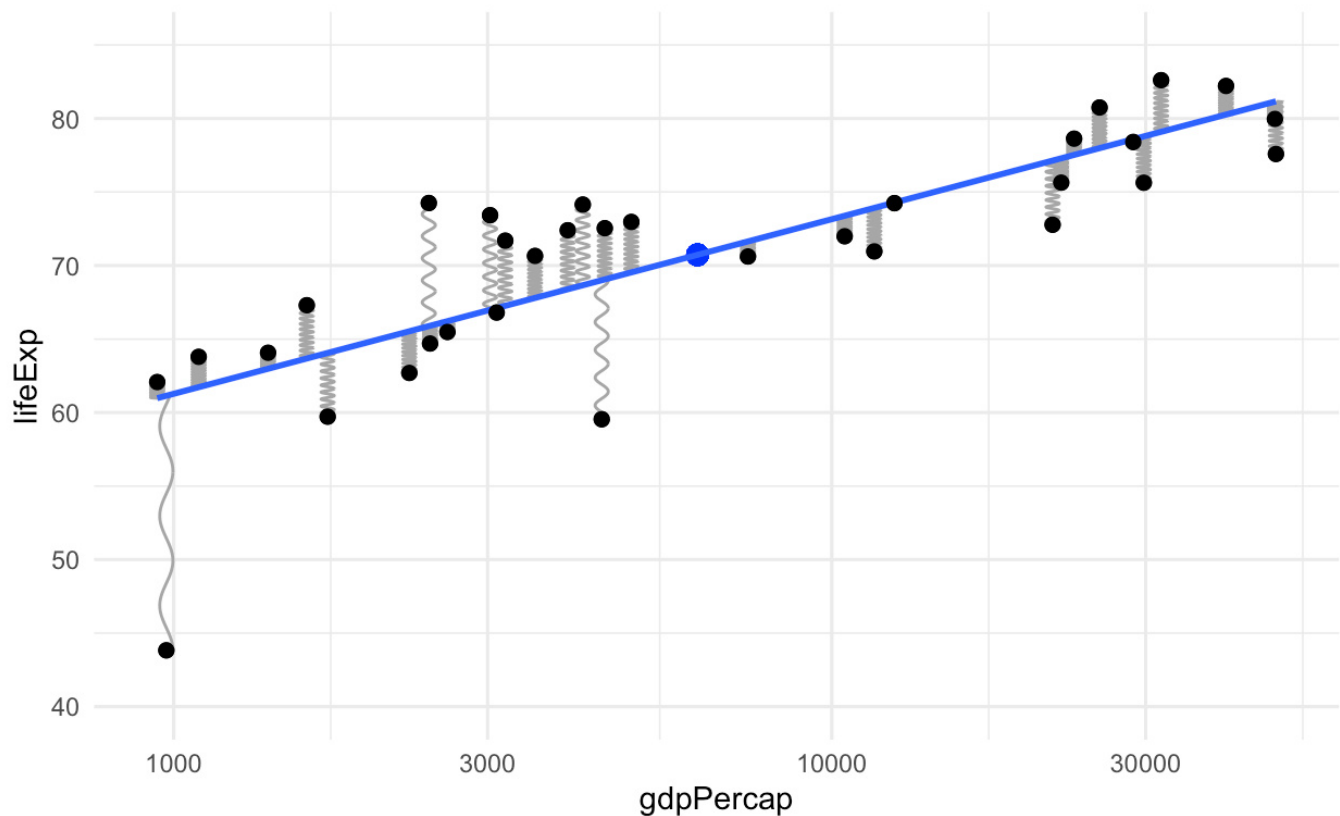
38

Regression: Scatterplot



39

Regression: Underlying *mechanics*



40

Regression: Statistics

```
1 lm_out0 <- lm(`Weight Of Sample Before Drying` ~ `Weight Of Empty Cup`,
2               data=rawdata)
3 lm_out0
```

Call:

```
lm(formula = `Weight Of Sample Before Drying` ~ `Weight Of Empty Cup`,
    data = rawdata)
```

Coefficients:

```
(Intercept)  `Weight Of Empty Cup`
      10.169           -1.228
```

```
1 # filtering outlier
2 lm_out <- lm(`Weight Of Sample Before Drying` ~ `Weight Of Empty Cup`,
3             data=rawdata |> filter(`Weight Of Sample Before Drying`>4.5))
4 lm_out
```

Call:

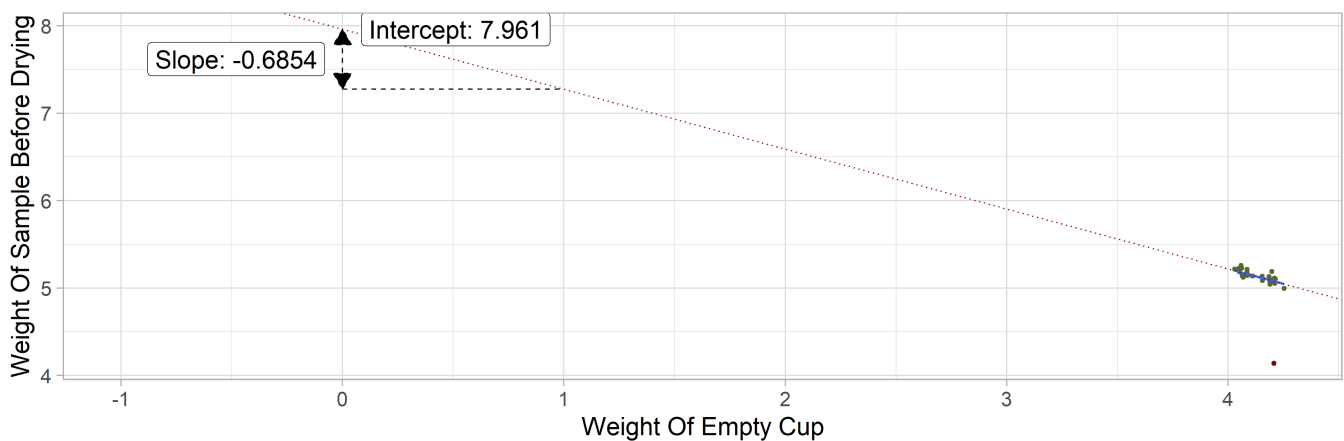
```
lm(formula = `Weight Of Sample Before Drying` ~ `Weight Of Empty Cup`,
    data = filter(rawdata, `Weight Of Sample Before Drying` >
4.5))
```

Coefficients:

```
(Intercept)  `Weight Of Empty Cup`
      7.9612           -0.6854
```

41

Regression: Visualize



42

Regression: Significance

```
1 anova(lm_out) |> broom::tidy()
```

```
# A tibble: 2 × 6
```

	term <chr>	df <int>	sumsq <dbl>	meansq <dbl>	statistic <dbl>	p.value <dbl>
1	`Weight Of Empty Cup`	1	0.0638	0.0638	42.1	0.000000425
2	Residuals	29	0.0440	0.00152	NA	NA

```
1 model_parameters(lm_out)
```

Parameter	Coefficient	SE	95% CI	t(29)	p
(Intercept)	7.96	0.44	[7.07, 8.85]	18.28	< .001
Weight Of Empty Cup	-0.69	0.11	[-0.90, -0.47]	-6.49	< .001

43

Report

- *RMarkdown and quarto are powerful tools to create reports and presentations*
- Export figures: `ggsave()` / `png()` / `pdf()`
- Export tables: `write_xlsx()`
- *Package flextable provides nice features for table formatting*

44

Flextable example

```
1 test_ord |> select(-desc_all) |> rename_with(~str_remove(., 'Code Of ')) |>
2 flextable() |>
3 theme_zebra(even_body = 'aquamarine', odd_body = 'antiquewhite') |>
4 italic(~p<=0.05, j = 1) |> bg(~p<=0.05, j = 4, bg = 'yellow') |>
5 set_caption('Treatment effects, measures following a normal distribution') |
6 add_footer_lines('Significance level is set at 0.05') |>
7 fontsize(size = 12, part = 'footer')
```

Variable	Region A	Region D	p
<i>Weight Of Cup + Sample</i>	9.24 (9.22/9.25)	9.29 (9.27/9.31)	0.003
<i>Weight Of Sample Before Drying</i>	5.14 (5.13/5.20)	5.09 (5.06/5.13)	0.001

Significance level is set at 0.05

45

Useful tools along the way

- Pick columns / rows: `select()` / `pull()` / `filter()` / `slice()`
- Change format of tibble wide <=> long (e.g. for repeated measures): `pivot_longer()`/`pivot_wider()`
- Regular expressions: `str_replace()` / `str_detect()` / `str_...`
- Merge text elements: `paste()` / `str_glue()`
- Apply functions: `purrr::map_XXX`

46