

Data cleaning and exploration

Jody Daniel

2020-11-03

Contents

Traits and Plant Growth Rates	1
Table of Contents	1
Extracting Principle Components for Environmental Traits	2
Including Plots	3
Estimate Growth Rate	5

Traits and Plant Growth Rates

We have two datasets: 1) basal area each year from 2006 to 2011 and 2) plant and environmental traits for each tree. Below, I aim to examine the data - assess data types, check for missingness, normality, outliers. Depending on the degree of missingness, I will impute the data as we would like to keep as many observations as possible.

Table of Contents

- Initial Assessment
- Imputation
- Estimate Growth Rate
- Conclusion

```
library(corrplot)
library(fastDummies)
library(RColorBrewer)
library(factoextra)
library(ggplot2)
require(ggrepel)
library(knitr)
library(kableExtra)
library(tidyverse)
library(dplyr)
library(here)
library(skimr)
library(reshape2)
```

```

library(tidymodels)
library(qdapTools)
library(rsample)
library(corr)
library(broom)
library(vegan)
library(extrafont)
library(viridis)
library(car)
library(mice)
library(sjmisc)
library(skimr)
library(RVAideMemoire)
source(here("scripts/archive/1. functions.R"))
theme_set(theme_special())

```

Extracting Principle Components for Environmental Traits

```

rgr_raw <- read.csv(here("data/RGR.csv"))
msh_raw <- read.csv(here("data/MSH.csv"))
# what do these data look like?
kable(skim(rgr_raw), "latex", booktabs = T) %>%
  kable_styling(latex_options="scale_down")

```

skim_type	skim_variable	n_missing	complete_rate	character.min	character.max	character.empty	character.n_unique	character.whitespace	numeric.mean	numeric.sd	numeric.p0	numeric.p25	numeric.p50	numeric.p75	numeric.p100	numeric.hist
character	SampleID	0	1.000000	6	10	0	395	0	NA	NA	NA	NA	NA	NA	NA	NA
numeric	BA.0.2006	9	0.9772152	NA	NA	NA	NA	NA	6.843472	6.298380	0.000000	2.520196	5.154124	9.998959	50.22676	<U+2587><U+2582><U+2581><U+2581><U+2581>
numeric	BA.0.2007	10	0.9746835	NA	NA	NA	NA	NA	7.784849	6.662033	0.000000	2.997015	5.626886	11.34191	51.51465	<U+2587><U+2583><U+2581><U+2581><U+2581>
numeric	BA.0.2008	10	0.9746835	NA	NA	NA	NA	NA	8.872749	7.052922	0.0302196	3.730206	7.189076	12.68398	52.24834	<U+2587><U+2583><U+2581><U+2581><U+2581>
numeric	BA.0.2009	10	0.9746835	NA	NA	NA	NA	NA	10.115077	7.547538	0.0758724	4.226576	8.733359	14.51579	52.78817	<U+2587><U+2585><U+2581><U+2581><U+2581>
numeric	BA.0.2010	10	0.9746835	NA	NA	NA	NA	NA	11.366633	8.163326	0.1405362	5.171318	9.901115	15.69798	57.18753	<U+2587><U+2585><U+2581><U+2581><U+2581>
numeric	BA.0.2011	9	0.9772152	NA	NA	NA	NA	NA	12.494634	8.797904	0.1697704	5.679960	11.022433	16.68853	58.07456	<U+2587><U+2586><U+2582><U+2581><U+2581>

```

# skim(rgr_raw) - for markdown visualization
kable(skim(msh_raw), "latex", booktabs = T) %>%
  kable_styling(latex_options="scale_down")

```

skim_type	skim_variable	n_missing	complete_rate	character.min	character.max	character.empty	character.a_unique	character.whitespace	numeric.mean	numeric.sd	numeric.p0	numeric.p25	numeric.p50	numeric.p75	numeric.p100	numeric.hist
character	SampleID	0	1.000000	6	10	0	395	0	NA	NA	NA	NA	NA	NA	NA	NA
character	Species	0	1.000000	3	4	0	25	0	NA	NA	NA	NA	NA	NA	NA	NA
character	Site	0	1.000000	3	19	0	23	0	NA	NA	NA	NA	NA	NA	NA	NA
numeric	Height:DBH:Ratio	0	1.000000	NA	NA	NA	NA	149.737218	41.4776792	66.2606567	120.000000	144.6153846	173.8119626	292.000000	<U+2583><U+2587><U+2585><U+2582><U+2581>	
numeric	Estom	0	1.000000	NA	NA	NA	NA	4056.1447300	3332.2107114	360.4871410	1670.3688000	3158.5020190	5427.7383348	21146.010000	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	Erwig	0	1.000000	NA	NA	NA	NA	343.4024492	192.0000004	22.5902873	200.9902223	312.1629240	431.5603081	1124.2157920	<U+2586><U+2587><U+2582><U+2581><U+2581>	
numeric	Breaching:Distance	0	1.000000	NA	NA	NA	NA	23.9343212	21.2721000	5.1333333	11.0934783	16.1428571	27.3000000	100.000000	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	Twig:Distance	0	1.000000	NA	NA	NA	NA	5.5794038	1.8103700	3.0250000	4.4500000	5.1200000	6.2500000	15.050000	<U+2587><U+2583><U+2581><U+2581><U+2581>	
numeric	Twig:Wood:Density	0	1.000000	NA	NA	NA	NA	0.5018101	0.0772251	0.2607191	0.4558797	0.5115347	0.6759929	0.5538173	<U+2581><U+2583><U+2587><U+2587><U+2582><U+2581>	
numeric	Stem:Wood:Density	0	1.000000	NA	NA	NA	NA	0.6162200	0.1232942	0.2642821	0.5167250	0.6604948	0.7272644	0.8971172	<U+2581><U+2583><U+2587><U+2587><U+2587><U+2583>	
numeric	Leaf:Mass:Fraction	0	1.000000	NA	NA	NA	NA	1.2869830	0.6559224	0.2307848	0.6101039	1.1522226	1.6449689	5.5339309	<U+2586><U+2587><U+2583><U+2582><U+2581>	
numeric	Leaf:Area	0	1.000000	NA	NA	NA	NA	59.3640807	78.2130797	3.7144000	18.2802344	29.7845887	47.4472723	191.4577273	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	LMA	0	1.000000	NA	NA	NA	NA	52.5578845	24.2827043	17.4592553	32.9624497	45.3863990	67.1790945	134.2225368	<U+2587><U+2586><U+2583><U+2582><U+2581>	
numeric	LCC	0	1.000000	NA	NA	NA	NA	47.5027863	2.3657353	41.7288800	46.6458900	47.6190000	49.1278000	52.5320000	<U+2582><U+2583><U+2587><U+2586><U+2582>	
numeric	LNC	0	1.000000	NA	NA	NA	NA	2.6920081	0.5666528	1.4219000	2.3377700	2.6778000	3.0489000	4.1779000	<U+2581><U+2586><U+2587><U+2583><U+2581>	
numeric	LPC	0	1.000000	NA	NA	NA	NA	15.3871587	4.8429936	5.5412378	12.0663117	14.4697182	17.8304845	41.3029907	<U+2583><U+2587><U+2582><U+2581><U+2581>	
numeric	d12N	0	1.000000	NA	NA	NA	NA	2.2289654	2.8404845	-0.0300275	-1.1032924	-2.2207275	-4.7338487	16.1917920	<U+2583><U+2587><U+2581><U+2581><U+2581>	
numeric	t12	7	0.982275	NA	NA	NA	NA	0.0250138	0.0090180	0.0007625	0.0191504	0.0325503	0.0580334	0.0325503	<U+2587><U+2581><U+2586><U+2582><U+2581>	
numeric	K1	8	0.9797468	NA	NA	NA	NA	117.0343280	144.7195939	10.9835078	30.2909683	57.4040847	122.2114627	985.8300149	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	Kwig	8	0.9797468	NA	NA	NA	NA	1156.3121313	901.7567951	90.7286105	122.4739055	899.1533593	1331.7584990	7135.0907761	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	Huber:Value	4	0.9898734	NA	NA	NA	NA	0.1440426	0.0063790	0.0014640	0.0050178	0.0083976	0.0165608	5.9961641	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	X_Lam	8	0.9797468	NA	NA	NA	NA	0.1326261	0.0444133	0.0357374	0.1030985	0.2781733	0.6114926	0.7781733	<U+2582><U+2587><U+2581><U+2582><U+2581>	
numeric	VID	8	0.9797468	NA	NA	NA	NA	10521.0232322	9650.2092383	1988.4913796	3095.1447790	6765.9227750	10350.0596500	32249.4432300	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	X_Sapwood	4	0.9898734	NA	NA	NA	NA	0.7648435	0.2847693	0.0089190	0.7203337	0.9036817	0.9596712	0.9596712	<U+2581><U+2581><U+2587><U+2582><U+2581>	
numeric	d13C	0	1.000000	NA	NA	NA	NA	-30.2257774	1.4615324	-33.8295092	-31.5694677	-30.4214853	-29.38093104	-25.8229007	<U+2582><U+2587><U+2587><U+2583><U+2581>	
numeric	Biomass	11	0.9721519	NA	NA	NA	NA	6.1811851	7.9142907	0.0319446	1.8609618	3.8306772	7.6292229	66.0562218	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	pent.max.BioI	11	0.9721519	NA	NA	NA	NA	0.3116330	0.2892126	0.0063699	0.855031	0.2116297	0.4476306	1.0000000	<U+2587><U+2583><U+2582><U+2581><U+2582>	
numeric	Symmetric.Competition	9	0.9772152	NA	NA	NA	NA	2.7314834	3.8865711	0.0221145	0.0696180	1.3246754	2.2648763	27.0587752	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	Asymmetric.Competition	9	0.9772152	NA	NA	NA	NA	20.2829284	23.1118657	0.2162967	4.4161511	10.0840750	27.5335537	149.3885902	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	Soil.Humidity	0	1.000000	NA	NA	NA	NA	0.1547539	0.1578224	0.0036667	0.0610000	0.1053333	0.1858333	1.000000	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	pH	130	0.6708861	NA	NA	NA	NA	5.1269792	0.7340890	3.6100000	4.5700000	5.0900000	5.6900000	7.1900000	<U+2583><U+2587><U+2587><U+2587><U+2583><U+2581>	
numeric	Organic:C	130	0.6708861	NA	NA	NA	NA	11.0842800	10.0567790	0.6850316	3.2127755	7.4423036	14.3739874	44.2885339	<U+2587><U+2583><U+2581><U+2582><U+2581>	
numeric	N	130	0.6708861	NA	NA	NA	NA	0.4505730	0.4793418	0.0076733	0.1025588	0.2692148	0.5488671	2.1087924	<U+2587><U+2582><U+2581><U+2581><U+2581>	
numeric	P	130	0.6708861	NA	NA	NA	NA	32.2434435	30.6929773	2.1562818	12.028980	18.1015122	33.9273519	222.939748	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	K	130	0.6708861	NA	NA	NA	NA	138.2513033	95.1734560	24.2527806	47.2213774	112.7995236	189.7051184	531.1374880	<U+2587><U+2583><U+2581><U+2581><U+2581>	
numeric	C	130	0.6708861	NA	NA	NA	NA	2074.9988264	2312.2158636	116.9189371	505.7524772	1351.8075700	2564.6472360	14134.2452000	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	Mg	130	0.6708861	NA	NA	NA	NA	107.0131431	93.1819573	9.4122099	40.2245095	78.6918416	149.2823380	615.2652603	<U+2587><U+2583><U+2581><U+2581><U+2581>	
numeric	Soil.Depth	82	0.7321651	NA	NA	NA	NA	31.1948882	14.9059141	2.0000000	21.0000000	30.0000000	40.0000000	55.0000000	<U+2582><U+2586><U+2587><U+2583><U+2587>	
numeric	Shape	0	1.000000	NA	NA	NA	NA	8.6549367	10.1938941	0.0000000	0.0000000	0.0000000	3.0000000	39.0000000	<U+2587><U+2581><U+2582><U+2581><U+2581>	
numeric	North.Aspect	0	1.000000	NA	NA	NA	NA	-0.1690550	0.4854480	-1.0000000	-0.5145956	0.0000000	0.0000000	0.9902681	<U+2583><U+2581><U+2587><U+2583><U+2581>	
numeric	East.Aspect	0	1.000000	NA	NA	NA	NA	0.0592598	0.5074686	-0.9992808	0.0000000	0.0000000	0.0000000	0.0000000	<U+2582><U+2581><U+2587><U+2581><U+2582>	
numeric	Elevation	0	1.000000	NA	NA	NA	NA	142.4911392	76.4742868	40.0000000	81.5000000	146.0000000	200.0000000	350.0000000	<U+2586><U+2587><U+2585><U+2581><U+2581>	
numeric	Summer.Max	0	1.000000	NA	NA	NA	NA	27.1680759	1.6051104	24.1700000	26.1700000	26.6700000	27.7300000	32.1700000	<U+2582><U+2587><U+2582><U+2583><U+2581>	
numeric	Fall.Min	0	1.000000	NA	NA	NA	NA	-10.2153418	1.2669112	-14.6700000	-10.3800000	-10.5000000	-9.5000000	-8.0000000	<U+2581><U+2581><U+2587><U+2586><U+2585>	
numeric	Mean.Temp	0	1.000000	NA	NA	NA	NA	13.7220000	0.5813513	12.3100000	13.2200000	13.8300000	14.1000000	14.6700000	<U+2581><U+2583><U+2585><U+2587><U+2587><U+2585>	
numeric	X.Canopy.Opening	0	1.000000	NA	NA	NA	NA	9.1948065	7.9404999	1.0000000	4.7650000	6.4000000	10.5000000	57.0100000	<U+2587><U+2581><U+2581><U+2581><U+2581>	
numeric	Tree.Age	6	0.9848101	NA	NA	NA	NA	20.8228221	10.2883423	5.0000000	14.0000000	17.0000000	20.0000000	60.0000000	<U+2587><U+2585><U+2582><U+2581><U+2581>	
numeric	Tree.Height	0	1.000000	NA	NA	NA	NA	421.4204990	131.4865550	170.0000000	228.0000000	404.5094528	589.0000000	922.0000000	<U+2585><U+2587><U+2587><U+2581><U+2581>	
numeric	julian.date.2011	0	1.000000	NA	NA	NA	NA	194.3518987	121.570038	171.0000000	186.0000000	194.0000000	204.0000000	218.0000000	<U+2583><U+2587><U+2587><U+2583><U+2583>	
numeric	Breath.Height	0	1.000000	NA	NA	NA	NA	3.0430739	1.0307331	1.2000000	2.9800000	3.4500000	4.0000000	6.8000000	<U+2585><U+2587><U+2585><U+2582><U+2581>	
numeric	SHL	212	0.3873418	NA	NA	NA	NA	105.8022774	61.145386	14.1784286	64.2643965	92.0142346	143.7505535	288.1373258	<U+2587><U+2587><U+2585><U+2582><U+2581>	
numeric	Root.Wood.Density	196	0.5037975	NA	NA	NA	NA	0.5683471	0.0990170	0.2954545	0.5080836	0.5695900	0.6401440	0.8489686	<U+2581><U+2585><U+2587><U+2586><U+2581>	
numeric	Twig.Breaching.angle	100	0.7461574	NA	NA	NA	NA	53.7744974	16.9091939	0.0000000	46.8821125	54.9505000	64.5755950	102.8261590	<U+2581><U+2582><U+2587><U+2583><U+2581>	
numeric	Hmax	0	1.000000	NA	NA	NA	NA	20.2582278	9.2461396	5.0000000	12.0000000	25.0000000	25.0000000	35.0000000	<U+2586><U+2585><U+2582><U+2587><U+2583>	
numeric	Shade.Tolerance	37	0.9662291	NA	NA	NA	NA	2.7761732	1.1405385	1.2100000	1.5600000	2.5900000	3.5600000	4.7600000	<U+2587><U+2586><U+2585><U+2583><U+2585>	
numeric	Drought.Tolerance	37	0.9662291	NA	NA	NA	NA	2.4795331	0.7942697	1.5000000	1.7700000	2.3800000	2.9200000	4.0000000	<U+2587><U+2583><U+2586><U+2582><U+2582>	
numeric	WaterLogging.Tolerance	37	0.9662291	NA	NA	NA	NA	1.7887809	0.7117861	1.0000000	1.0700000	1.5600000	2.5000000	3.3700000	<U+2587><U+2	

Now, I will run the imputation on these data. Since I have removed the columns with a large amount of missing cases, the imputation should work. If not, I will set a more stringent condition.

```
# fit a lm and see if results are comparable between mice output and raw data
summary(with(data = MSH.IP.Final, exp = lm(Tree.Height ~ Soil.Depth + Biomass1 + Huber.Value + pH)))
```

```
##
## Call:
## lm(formula = Tree.Height ~ Soil.Depth + Biomass1 + Huber.Value +
##     pH)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -261.03  -97.07  -16.84   88.79  522.04
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  366.6112    54.4005   6.739 5.74e-11 ***
## Soil.Depth    0.8010     0.4709   1.701 0.089781 .
## Biomass1     11.7975     3.4656   3.404 0.000732 ***
## Huber.Value  -15.9981    10.9075  -1.467 0.143261
## pH           3.4535     9.8303   0.351 0.725540
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 129.2 on 390 degrees of freedom
## Multiple R-squared:  0.04325,    Adjusted R-squared:  0.03343
## F-statistic: 4.407 on 4 and 390 DF,  p-value: 0.001701
```

```
# fit a lm and see if results are comparable between mice output and raw data
summary(lm(Tree.Height ~ Soil.Depth + Biomass1 + Huber.Value + pH, data = msh_raw))
```

```
##
## Call:
## lm(formula = Tree.Height ~ Soil.Depth + Biomass1 + Huber.Value +
##     pH, data = msh_raw)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -275.02 -100.11  -25.28   87.24  527.31
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  376.5703    60.2157   6.254 1.67e-09 ***
## Soil.Depth    0.4151     0.5481   0.757  0.4495
## Biomass1     20.7473     4.6604   4.452 1.27e-05 ***
## Huber.Value  313.5464   135.3190   2.317  0.0213 *
## pH           1.9247    11.0954   0.173  0.8624
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 130.9 on 256 degrees of freedom
## (134 observations deleted due to missingness)
```

```
## Multiple R-squared:  0.09589,    Adjusted R-squared:  0.08176
## F-statistic: 6.788 on 4 and 256 DF,  p-value: 3.307e-05
```

The regression coefficients do vary slightly. But the direction and significance of these relationships are near the same.

```
# how do the combined and imputed datasets compare?
xyplot(MSH.IP.MICE, Soil.Depth ~ Soil.Humidity|.imp, pch = 20, cex = 1.4)
```

\includegraphics[]{{data-exploration-cleaning_files/figure-latex/unnamed-chunk-6-1}}

```
xyplot(MSH.IP.MICE, Soil.Depth ~ pH|.imp, pch = 20, cex = 1.4)
```

\includegraphics[]{{data-exploration-cleaning_files/figure-latex/unnamed-chunk-6-2}}

```
#The imputed data seems to match that of the raw data. Now, how do the combined imputed data compare to
merge_imputations(msh_raw, MSH.IP.MICE, summary = "dens")
```

\includegraphics[]{{data-exploration-cleaning_files/figure-latex/unnamed-chunk-6-3}}

There is strong overlap between the mean and merged values. As such, I can be confident that the merged data set is representative of individual imputations.

Estimate Growth Rate

There are three metrics we can use to measure growth rates. They are: * Basal Area Relative Growth Rate
* Biomass Relative Growth Rate

```
# makes more sense to name the rows by their sample ID as a unique identifier
MSH.IP.Final$SampleID <- rgr_raw[,1]

MSH.RGR_IMP <- merge(rgr_raw,MSH.IP.Final, by = "SampleID")
MSH.RGR <- merge(rgr_raw,MSH.50, by = "SampleID")
# now for growth rate parameters
BAI_GR <- with(data = MSH.RGR_IMP, exp = ((BA.0.2011 - BA.0.2006)/5))
BIO_Gain <- with(data = MSH.RGR_IMP,
  exp = ((BA.0.2011*Stem.Wood.Density*Tree.Height*3)-(BA.0.2006*Stem.Wood.Density*(Tree.Height/Tree.Age)))
BIO_2006 <- with(data = MSH.RGR_IMP,
  exp = (BA.0.2006*Stem.Wood.Density*(Tree.Height/Tree.Age*(Tree.Age-6))*3))
BIO_GR <- BIO_Gain/((BIO_2006+1)*5)
# wanna place species and site at the start of the data frame
site_species <- c("SampleID", "Site", "Species", "Porosity")
site_species_no <- which(colnames(MSH.50)%in%site_species)
site_BA <- c(colnames(rgr_raw)[-1], "SampleID")
site_BA_no <- which(colnames(MSH.RGR_IMP)%in%site_BA)

RGR_MSH_Final <- data.frame(MSH.50[,site_species_no], BAI_GR = BAI_GR,
  BIO_GR = BIO_GR,
  MSH.RGR_IMP[, -site_BA_no])

# now for growth rate parameters
```

```

BAI_GR_NA <- with(data = MSH.RGR, exp = ((BA.0.2011 - BA.0.2006)/5))
BIO_Gain_NA <- with(data = MSH.RGR, exp = ((BA.0.2011*Stem.Wood.Density*Tree.Height*3) - (BA.0.2006*Stem.Wood.Density*Tree.Height*3)))
BIO_2006_NA <- with(data = MSH.RGR, exp = (BA.0.2006*Stem.Wood.Density*(Tree.Height/Tree.Age*(Tree.Age-1))))
BIO_GR_NA <- BIO_Gain_NA/((BIO_2006_NA+1)*5)

RGR_MSH_NA_Final <- data.frame(MSH.50[,site_species_no], BAI_GR = BAI_GR_NA,
                               BIO_GR = BIO_Gain_NA,
                               MSH.50[,numeric_columns_msh])

write_csv(RGR_MSH_Final,here("data/RGR_MSH.csv"))
write_csv(RGR_MSH_NA_Final,here("data/RGR_MSH_NA.csv"))

```