

Assignment 1

Erling Tennøy Nordtvedt

Oleg Karpov

2024-09-16

Exercise 1 - R syntax & data structures

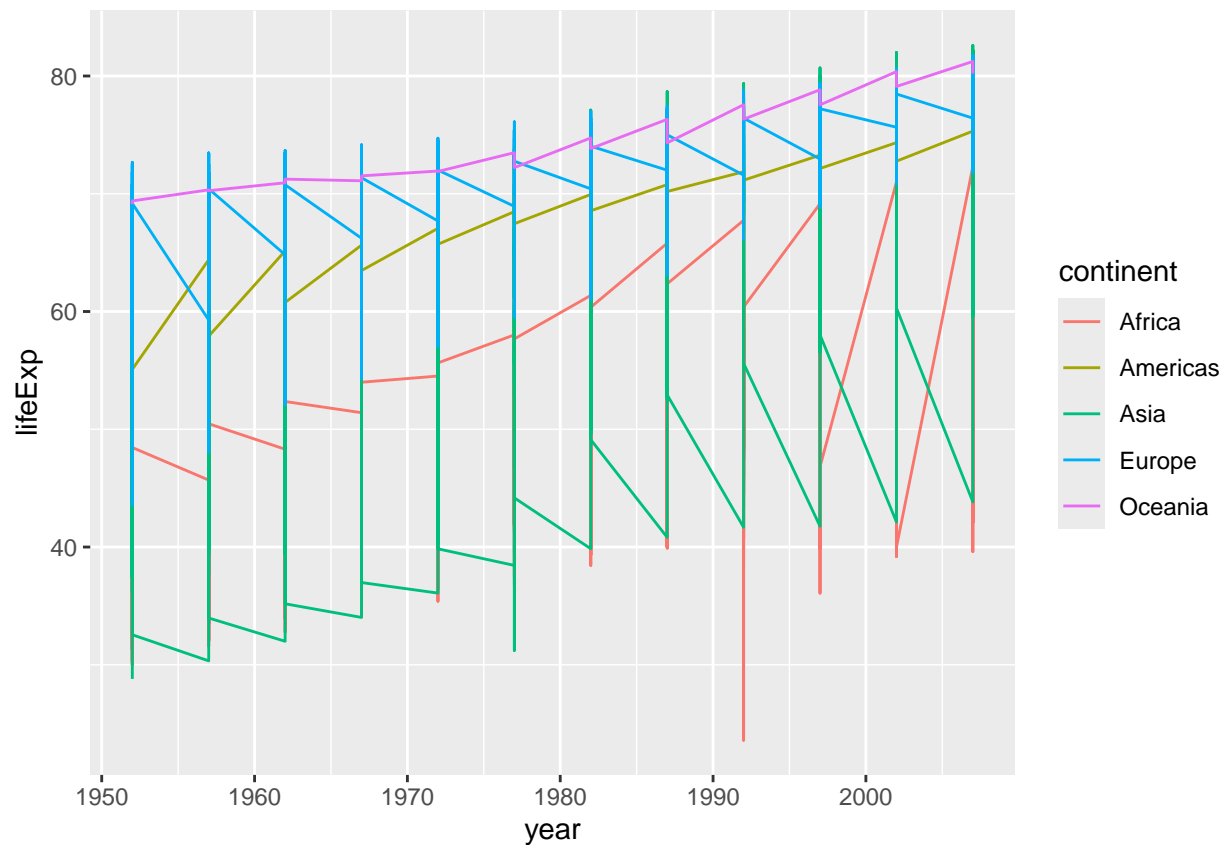
(a)

```
gapminder <- read.csv("gapminder.csv")  
  
summary(gapminder)
```

```
##           X           country           continent           year  
## Min.      : 1.0      Length:1704      Length:1704      Min.      :1952  
## 1st Qu.: 426.8      Class :character  Class :character  1st Qu.:1966  
## Median : 852.5      Mode  :character  Mode  :character  Median :1980  
## Mean    : 852.5  
## 3rd Qu.:1278.2  
## Max.    :1704.0  
##      lifeExp      pop      gdpPercap  
## Min.      :23.60   Min.      :6.001e+04   Min.      : 241.2  
## 1st Qu.:48.20     1st Qu.:2.794e+06     1st Qu.: 1202.1  
## Median :60.71     Median :7.024e+06     Median : 3531.8  
## Mean    :59.47     Mean    :2.960e+07     Mean    : 7215.3  
## 3rd Qu.:70.85     3rd Qu.:1.959e+07     3rd Qu.: 9325.5  
## Max.    :82.60     Max.    :1.319e+09     Max.    :113523.1
```

(b)

```
gapminder %>%  
  ggplot(aes(x=year, y=lifeExp, colour=continent)) +  
    #geom_bar(position='dodge', stat='identity')  
    geom_line()
```



(c)

```
knitr::kable(
  gapminder %>%
    group_by(continent, year) %>%
    summarise_at(vars(lifeExp), list(Min = min, Med = median, Mean = mean, Max = max, Sd = sd)) %>%
    data.frame()
)
```

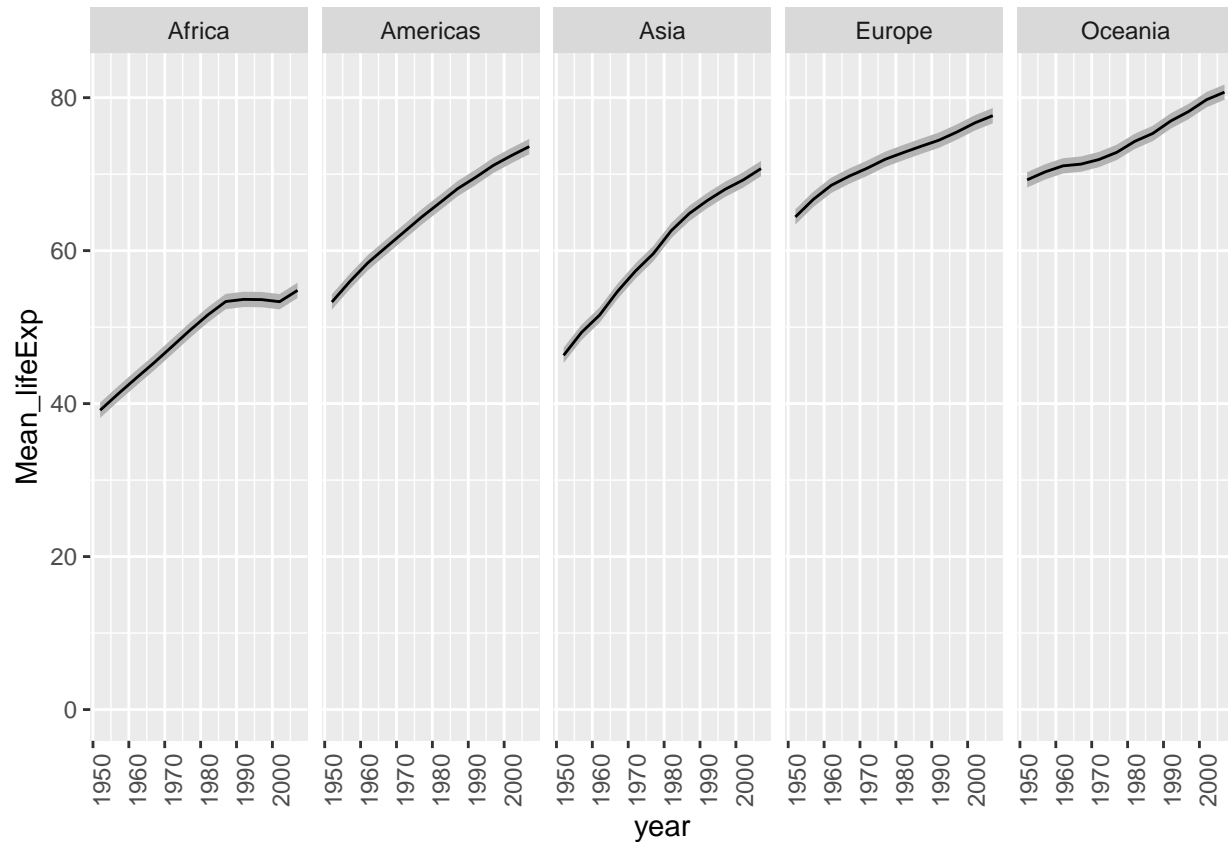
continent	year	Min	Med	Mean	Max	Sd
Africa	1952	30.000	38.8330	39.13550	52.724	5.1515814
Africa	1957	31.570	40.5925	41.26635	58.089	5.6201229
Africa	1962	32.767	42.6305	43.31944	60.246	5.8753639
Africa	1967	34.113	44.6985	45.33454	61.557	6.0826726
Africa	1972	35.400	47.0315	47.45094	64.274	6.4162583
Africa	1977	36.788	49.2725	49.58042	67.064	6.8081974
Africa	1982	38.445	50.7560	51.59287	69.885	7.3759401
Africa	1987	39.906	51.6395	53.34479	71.913	7.8640891
Africa	1992	23.599	52.4290	53.62958	73.615	9.4610710
Africa	1997	36.087	52.7590	53.59827	74.772	9.1033866
Africa	2002	39.193	51.2355	53.32523	75.744	9.5864959
Africa	2007	39.613	52.9265	54.80604	76.442	9.6307807
Americas	1952	37.579	54.7450	53.27984	68.750	9.3260819
Americas	1957	40.696	56.0740	55.96028	69.960	9.0331923
Americas	1962	43.428	58.2990	58.39876	71.300	8.5035437

continent	year	Min	Med	Mean	Max	Sd
Americas	1967	45.032	60.5230	60.41092	72.130	7.9091710
Americas	1972	46.714	63.4410	62.39492	72.880	7.3230168
Americas	1977	49.923	66.3530	64.39156	74.210	7.0694956
Americas	1982	51.461	67.4050	66.22884	75.760	6.7208338
Americas	1987	53.636	69.4980	68.09072	76.860	5.8019288
Americas	1992	55.089	69.8620	69.56836	77.950	5.1671038
Americas	1997	56.671	72.1460	71.15048	78.610	4.8875839
Americas	2002	58.137	72.0470	72.42204	79.770	4.7997055
Americas	2007	60.916	72.8990	73.60812	80.653	4.4409476
Asia	1952	28.801	44.8690	46.31439	65.390	9.2917507
Asia	1957	30.332	48.2840	49.31854	67.840	9.6354286
Asia	1962	31.997	49.3250	51.56322	69.390	9.8206319
Asia	1967	34.020	53.6550	54.66364	71.430	9.6509646
Asia	1972	36.088	56.9500	57.31927	73.420	9.7227000
Asia	1977	31.220	60.7650	59.61056	75.380	10.0221970
Asia	1982	39.854	63.7390	62.61794	77.110	8.5352214
Asia	1987	40.822	66.2950	64.85118	78.670	8.2037919
Asia	1992	41.674	68.6900	66.53721	79.360	8.0755490
Asia	1997	41.763	70.2650	68.02052	80.690	8.0911706
Asia	2002	42.129	71.0280	69.23388	82.000	8.3745954
Asia	2007	43.828	72.3960	70.72848	82.603	7.9637245
Europe	1952	43.585	65.9000	64.40850	72.670	6.3610883
Europe	1957	48.079	67.6500	66.70307	73.470	5.2958054
Europe	1962	52.098	69.5250	68.53923	73.680	4.3024996
Europe	1967	54.336	70.6100	69.73760	74.160	3.7997285
Europe	1972	57.005	70.8850	70.77503	74.720	3.2405764
Europe	1977	59.507	72.3350	71.93777	76.110	3.1210300
Europe	1982	61.036	73.4900	72.80640	76.990	3.2182603
Europe	1987	63.108	74.8150	73.64217	77.410	3.1696803
Europe	1992	66.146	75.4510	74.44010	78.770	3.2097811
Europe	1997	68.835	76.1160	75.50517	79.390	3.1046766
Europe	2002	70.845	77.5365	76.70060	80.620	2.9221796
Europe	2007	71.777	78.6085	77.64860	81.757	2.9798127
Oceania	1952	69.120	69.2550	69.25500	69.390	0.1909188
Oceania	1957	70.260	70.2950	70.29500	70.330	0.0494975
Oceania	1962	70.930	71.0850	71.08500	71.240	0.2192031
Oceania	1967	71.100	71.3100	71.31000	71.520	0.2969848
Oceania	1972	71.890	71.9100	71.91000	71.930	0.0282843
Oceania	1977	72.220	72.8550	72.85500	73.490	0.8980256
Oceania	1982	73.840	74.2900	74.29000	74.740	0.6363961
Oceania	1987	74.320	75.3200	75.32000	76.320	1.4142136
Oceania	1992	76.330	76.9450	76.94500	77.560	0.8697413
Oceania	1997	77.550	78.1900	78.19000	78.830	0.9050967
Oceania	2002	79.110	79.7400	79.74000	80.370	0.8909545
Oceania	2007	80.204	80.7195	80.71950	81.235	0.7290271

(d)

```
gapminder %>%
  group_by(continent, year) %>%
  summarise(Mean_lifeExp= mean(lifeExp, na.rm = T), .groups = 'drop') %>%
```

```
ggplot(aes(x=year, y=Mean_lifeExp)) +
  geom_ribbon(aes(ymin= Mean_lifeExp - 1, ymax = Mean_lifeExp + 1), fill = "grey70") +
  geom_line() +
  facet_grid(.~continent) +
  theme(axis.text.x = element_text(angle=90)) +
  ylim(0, NA)
```



Exercise 2 - Elementary data analysis and model training

(a)

```
weatherHistory <- read.csv("weatherHistory.csv")
head(weatherHistory)
```

```
##           Formatted.Date      Summary Precip.Type Temperature..C.
## 1 2006-04-01 00:00:00.000 +0200 Partly Cloudy      rain      9.472222
## 2 2006-04-01 01:00:00.000 +0200 Partly Cloudy      rain      9.355556
## 3 2006-04-01 02:00:00.000 +0200 Mostly Cloudy     rain      9.377778
## 4 2006-04-01 03:00:00.000 +0200 Partly Cloudy      rain      8.288889
## 5 2006-04-01 04:00:00.000 +0200 Mostly Cloudy     rain      8.755556
## 6 2006-04-01 05:00:00.000 +0200 Partly Cloudy      rain      9.222222
## Apparent.Temperature..C. Humidity Wind.Speed..km.h. Wind.Bearing..degrees.
## 1           7.388889      0.89      14.1197           251
## 2           7.227778      0.86      14.2646           259
## 3           9.377778      0.89       3.9284           204
## 4           5.944444      0.83      14.1036           269
```

```
## 5          6.977778      0.83          11.0446          259
## 6          7.111111      0.85          13.9587          258
##  Visibility..km. Loud.Cover Pressure..millibars.
## 1          15.8263          0          1015.13
## 2          15.8263          0          1015.63
## 3          14.9569          0          1015.94
## 4          15.8263          0          1016.41
## 5          15.8263          0          1016.51
## 6          14.9569          0          1016.66
##                      Daily.Summary
## 1 Partly cloudy throughout the day.
## 2 Partly cloudy throughout the day.
## 3 Partly cloudy throughout the day.
## 4 Partly cloudy throughout the day.
## 5 Partly cloudy throughout the day.
## 6 Partly cloudy throughout the day.
```

Qualitative nominal

- Summary
- Precip.Type
- Daily.Summary

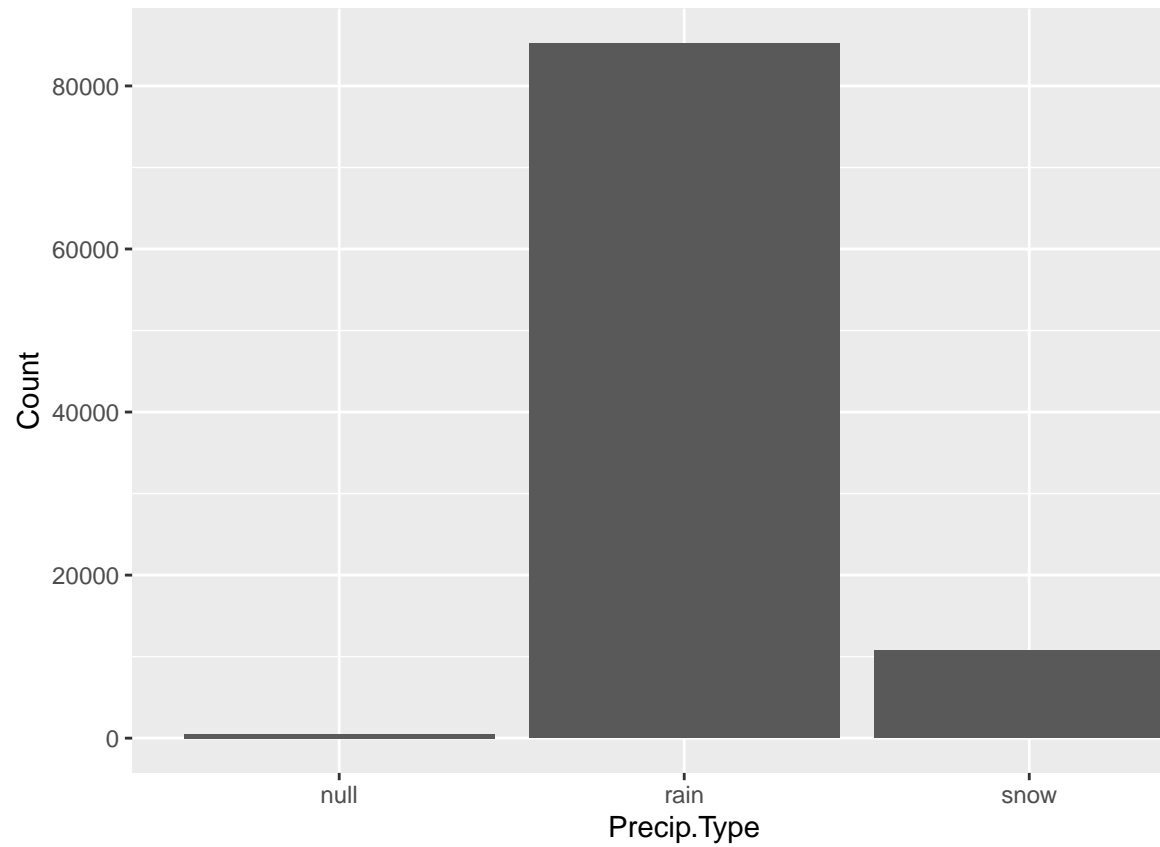
Quantitative Continuous:

- Temperature..C.
- Apparent.Temperature..C.
- Humidity
- Wind.Speed..km.h.
- Visibility..km.
- Wind.Bearing..degrees (Reason: Not ranked)

Quantitative Discrete:

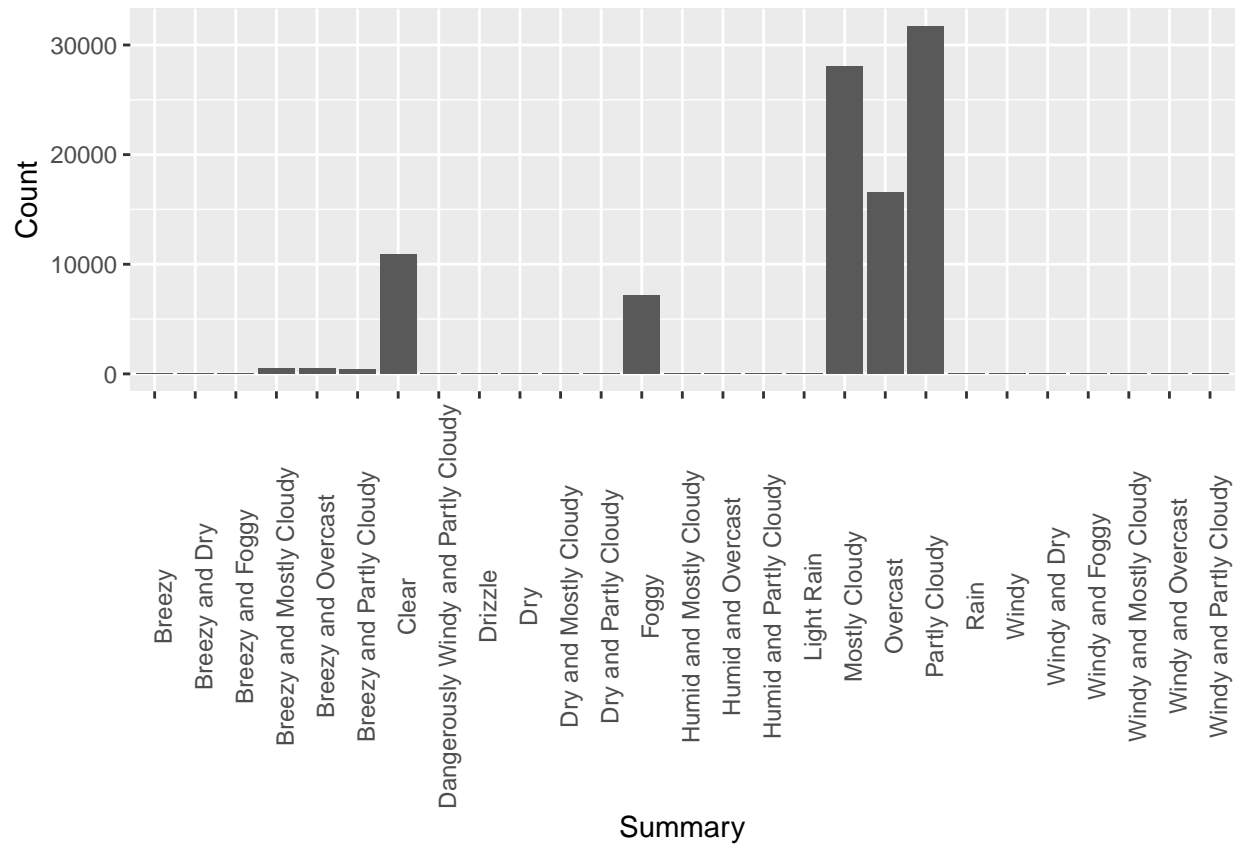
- Formatted.Date
- Loud.Cover

```
weatherHistory %>%
  group_by(Precip.Type) %>%
  summarize(Count = n()) %>%
  ggplot(aes(x=Precip.Type, y=Count)) +
  geom_bar(stat='identity', position='dodge')
```

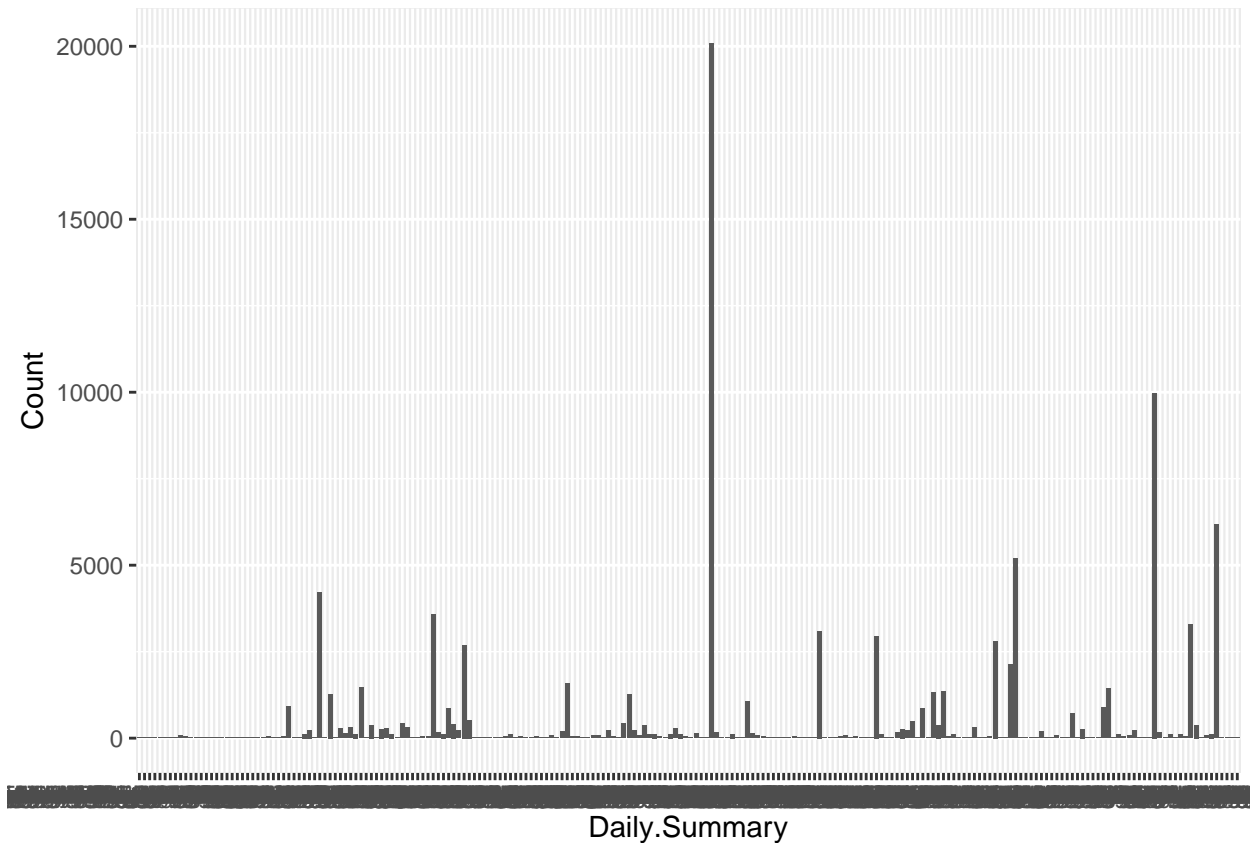


Qualitative nominal

```
weatherHistory %>%  
  group_by(Summary) %>%  
  summarize(Count = n()) %>%  
  ggplot(aes(x=Summary, y=Count)) +  
  geom_bar(stat='identity', position='dodge') +  
  theme(axis.text.x = element_text(angle=90))
```



```
weatherHistory %>%
  group_by(Daily.Summary) %>%
  summarize(Count = n()) %>%
  ggplot(aes(x=Daily.Summary, y=Count)) +
  geom_bar(stat='identity', position='dodge')
```

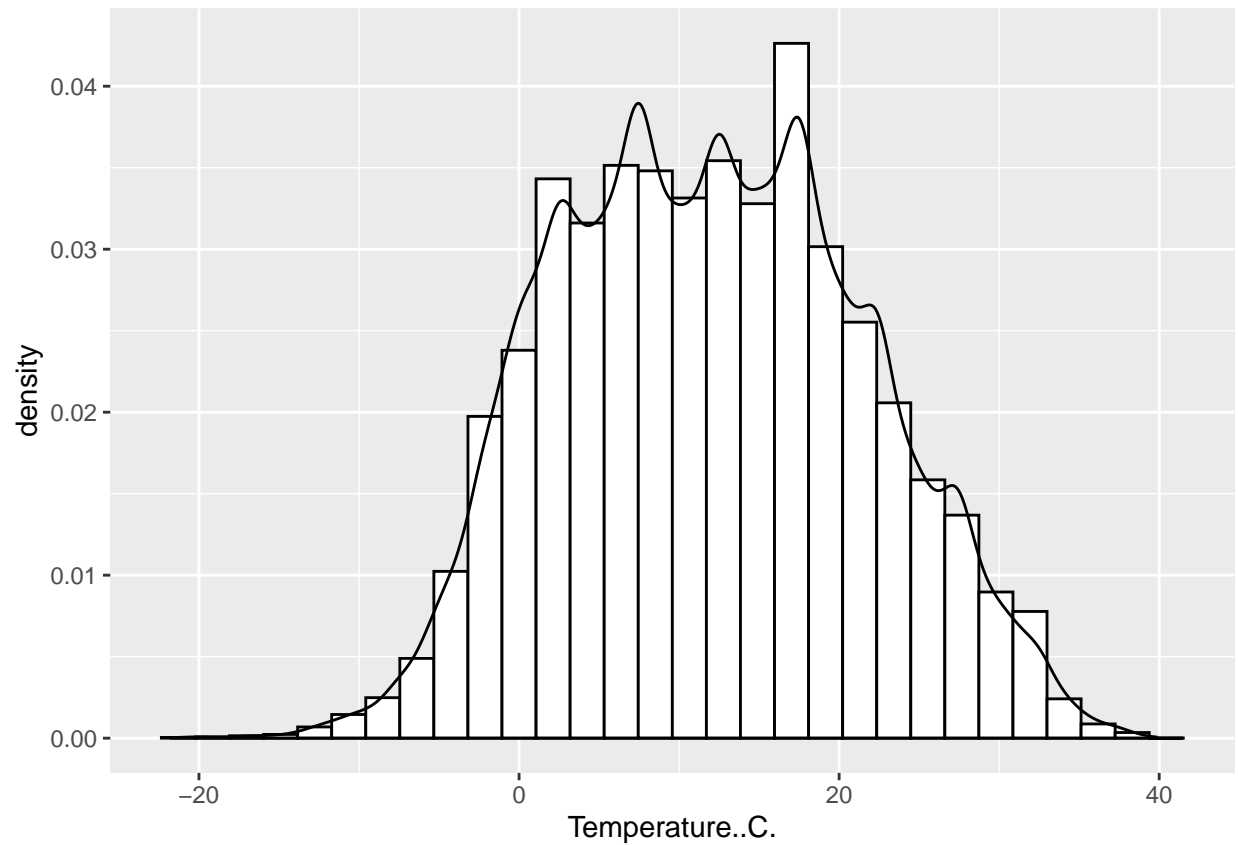


```
weatherHistory %>%
  ggplot(aes(Temperature..C.)) +
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +
  stat_density(kernel = "gaussian", fill = NA, colour = "black")
```

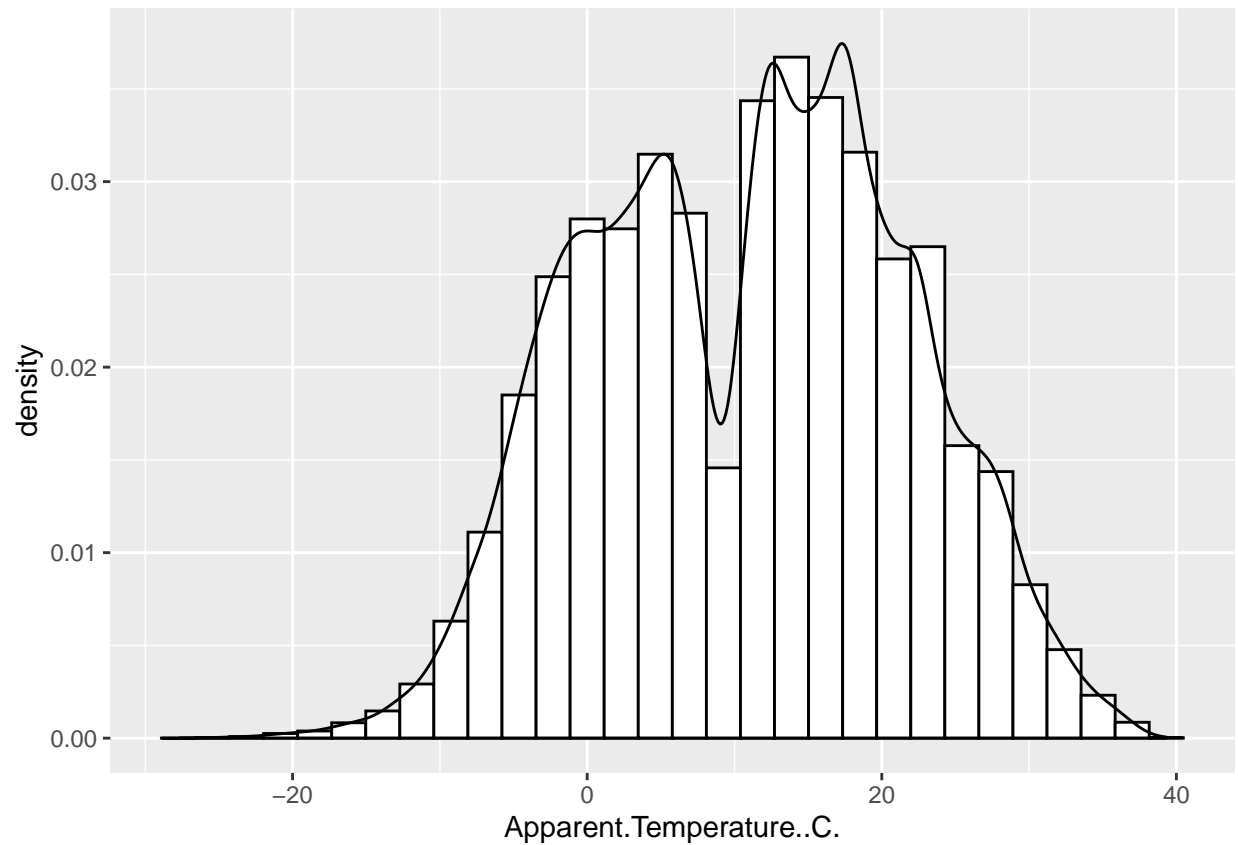
Discrete nominal

```
## Warning: The dot-dot notation (`..density..`) was deprecated in ggplot2 3.4.0.
## i Please use `after_stat(density)` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

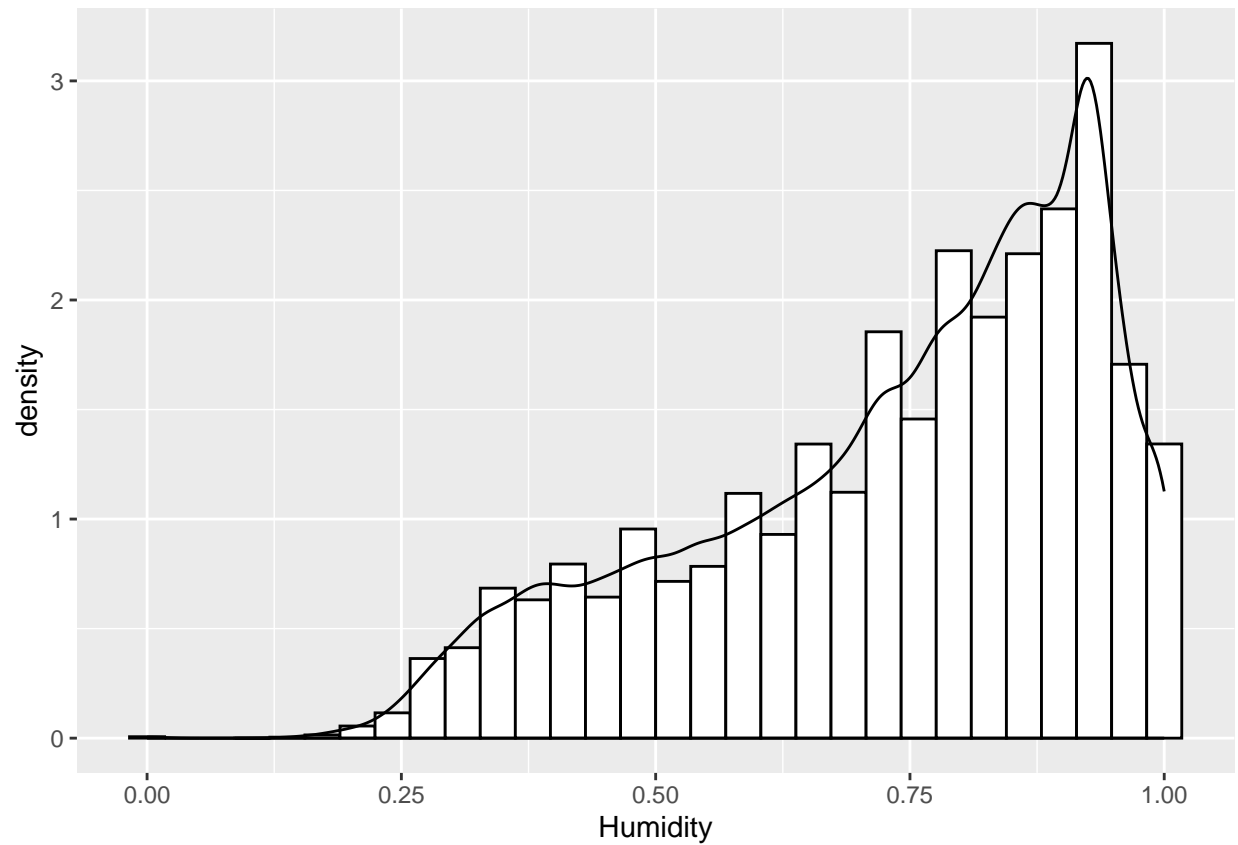
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

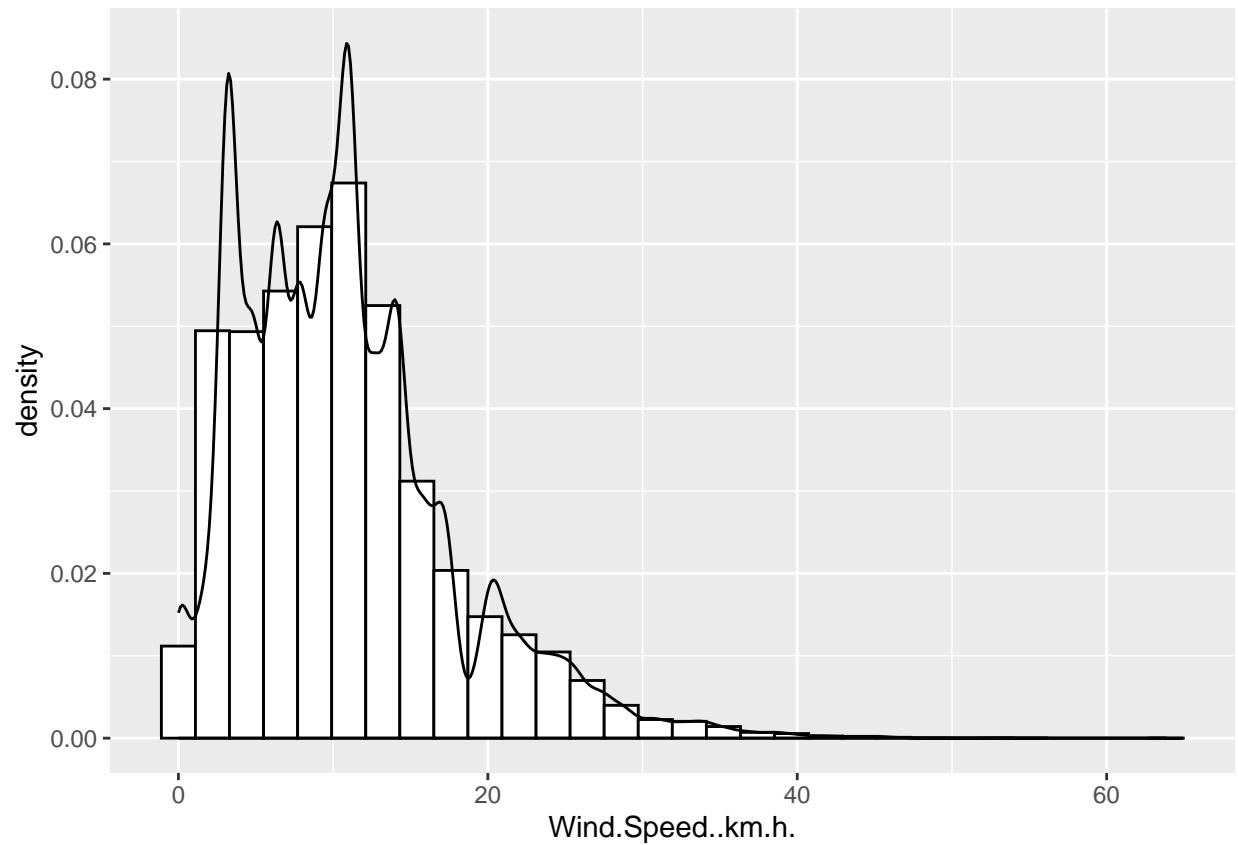
```
weatherHistory %>%  
  ggplot(aes(Apparent.Temperature..C.)) +  
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +  
  stat_density(kernel = "gaussian", fill = NA, colour = "black")  
  
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
weatherHistory %>%  
  ggplot(aes(Humidity)) +  
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +  
  stat_density(kernel = "gaussian", fill = NA, colour = "black")  
  
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

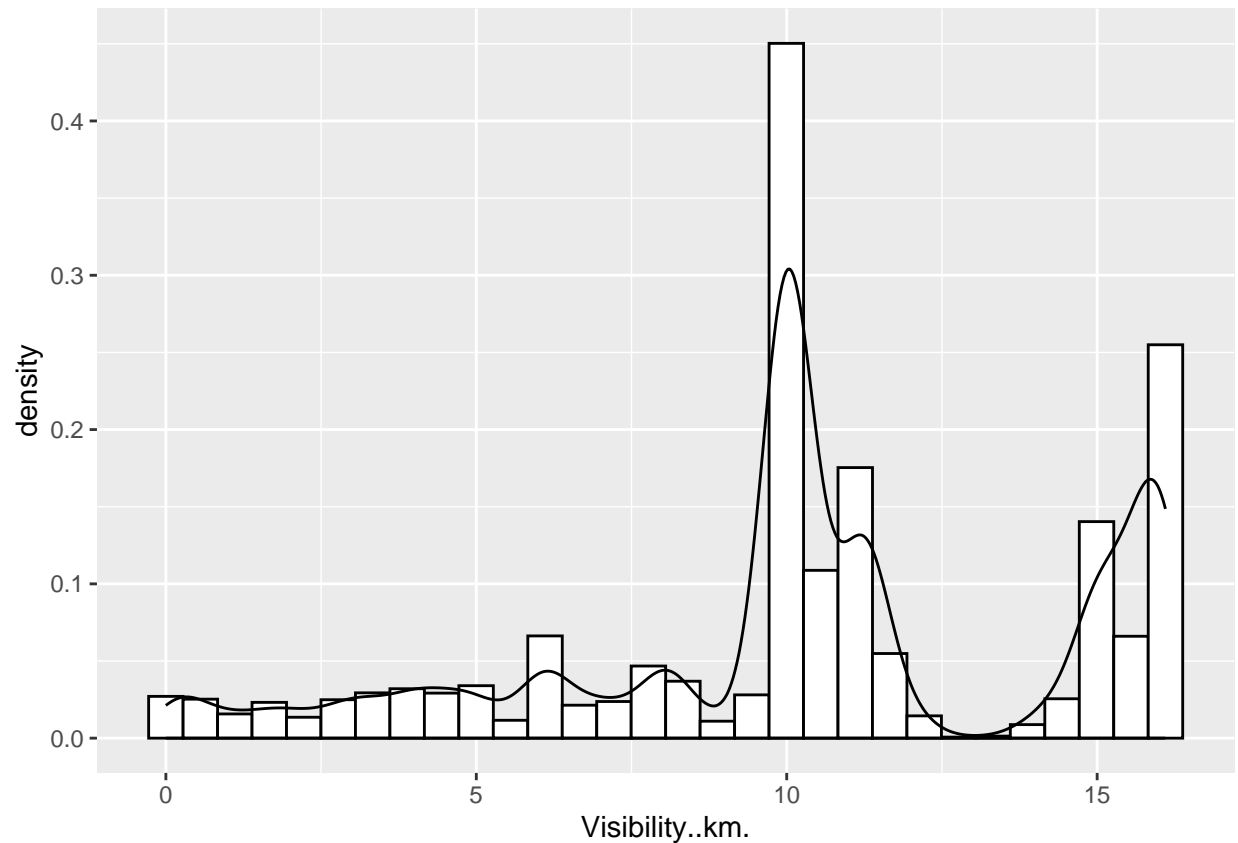


```
weatherHistory %>%  
  ggplot(aes(Wind.Speed..km.h.)) +  
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +  
  stat_density(kernel = "gaussian", fill = NA, colour = "black")  
  
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
weatherHistory %>%
  ggplot(aes(Visibility..km.)) +
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +
  stat_density(kernel = "gaussian", fill = NA, colour = "black")

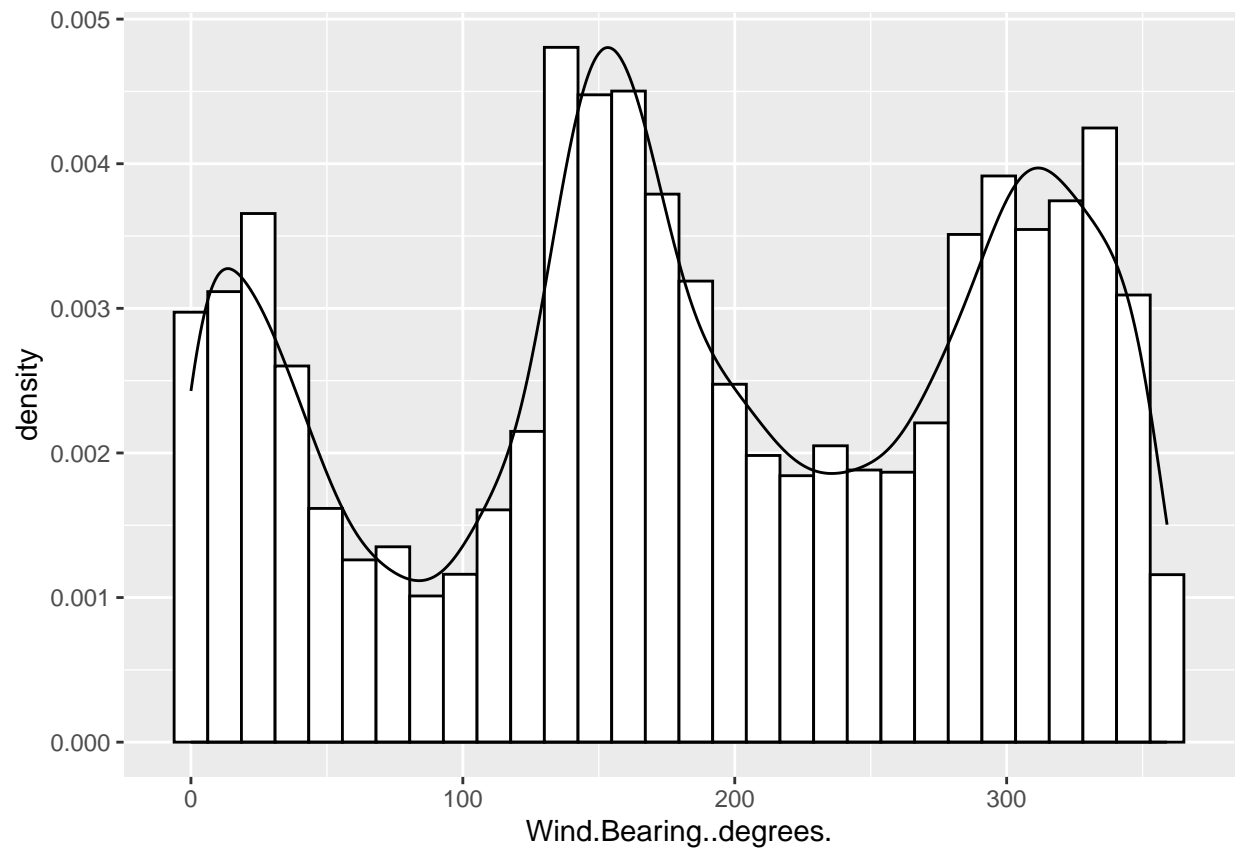
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



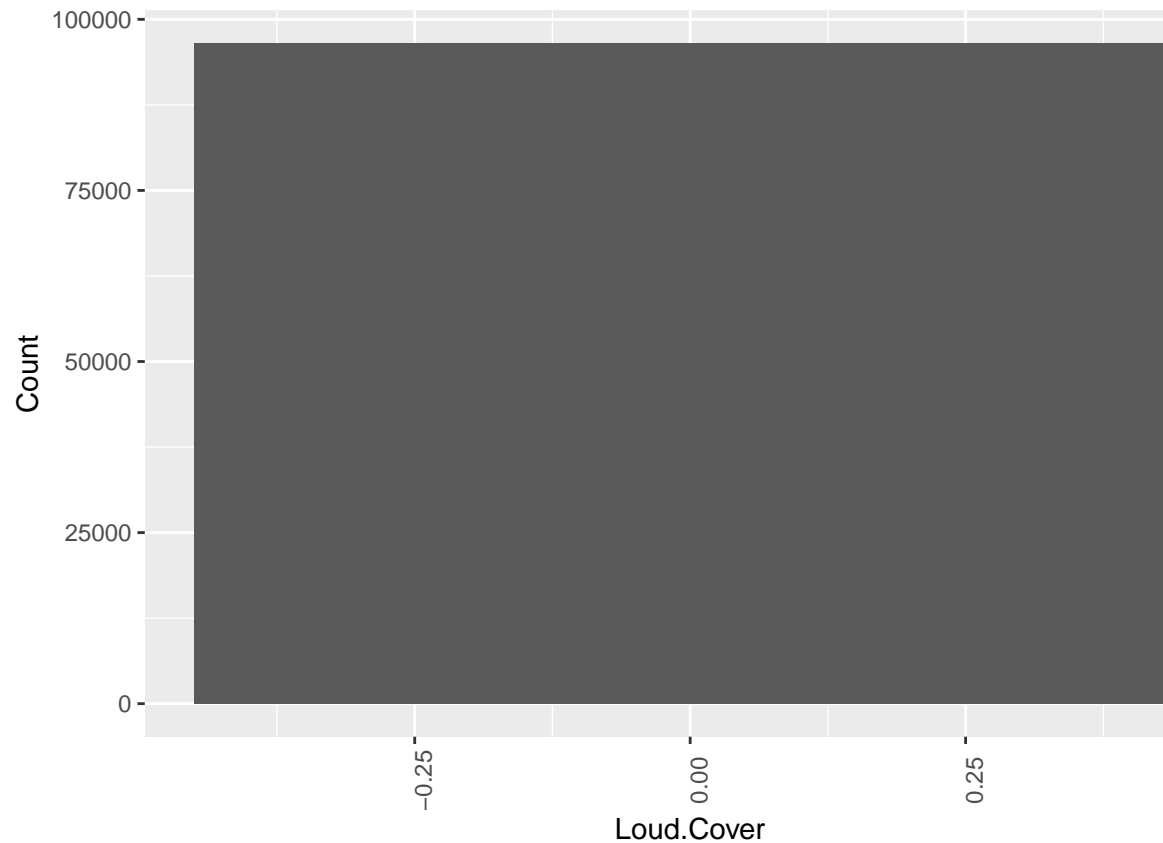
```
# weatherHistory %>%
#   group_by(Wind.Bearing..degrees.) %>%
#   summarize(Count = n()) %>%
#   ggplot(aes(x=Wind.Bearing..degrees., y=Count)) +
#   geom_bar(stat='identity', position='dodge') +
#   theme(axis.text.x = element_text(angle=90))

weatherHistory %>%
  ggplot(aes(Wind.Bearing..degrees.)) +
  geom_histogram(aes(y = ..density..), fill = "white", color="black") +
  stat_density(kernel = "gaussian", fill = NA, colour = "black")

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
weatherHistory %>%  
  group_by(Loud.Cover) %>%  
  summarize(Count = n()) %>%  
  ggplot(aes(x=Loud.Cover, y=Count)) +  
  geom_bar(stat='identity', position='dodge') +  
  theme(axis.text.x = element_text(angle=90))
```



Quantative discrete

(b)

First removing all columns that seem irrelevant, reasoning:

- Formatted.Date : When encoded it will be equal to row label (1, 2, 3, ...) which tells nothing
- Loud.Cover : All values are 0, therefore tells nothing
- Daily.Summary : Too big to onehotencode effectivly

Then remove all rows with NA, do this after removing irrelevant columns so data is not lost to having NA in the removed columns

```
library(caret)
```

```
## Loading required package: lattice
```

```
library(tidyr)
```

```
library(dplyr)
```

```
weatherHistory <- weatherHistory %>% select(-c("Formatted.Date", "Daily.Summary", "Loud.Cover"))
```

```
weatherHistory <- na.omit(weatherHistory) # Remove all NA
```

```
head(weatherHistory)
```

```
##      Summary Precip.Type Temperature..C. Apparent.Temperature..C. Humidity
## 1 Partly Cloudy      rain      9.472222      7.388889      0.89
## 2 Partly Cloudy      rain      9.355556      7.227778      0.86
## 3 Mostly Cloudy      rain      9.377778      9.377778      0.89
## 4 Partly Cloudy      rain      8.288889      5.944444      0.83
## 5 Mostly Cloudy      rain      8.755556      6.977778      0.83
## 6 Partly Cloudy      rain      9.222222      7.111111      0.85
```

```
## Wind.Speed..km.h. Wind.Bearing..degrees. Visibility..km. Pressure..millibars.
## 1      14.1197      251      15.8263      1015.13
## 2      14.2646      259      15.8263      1015.63
## 3       3.9284      204      14.9569      1015.94
## 4      14.1036      269      15.8263      1016.41
## 5      11.0446      259      15.8263      1016.51
## 6      13.9587      258      14.9569      1016.66

num_wH <- weatherHistory %>%
  select(-c("Summary", "Precip.Type"))
num_stand_wH <- as.data.frame(sapply(num_wH, function(x) ((x-mean(x))/sd(x))))

qualitative_wH <- weatherHistory %>%
  select(c("Summary", "Precip.Type")) #Omitted "Formatted.Date", "Daily.Summary"

# PT <- factor(qualitative_wH$Precip.Type)
# PT <- as.data.frame(model.matrix(~ Precip.Type - 1, PT))
#
# FD <- factor(qualitative_wH$Summary)
# FD <- as.data.frame(model.matrix(~ f - Summary - 1, FD))

q1 <- table(1:nrow(weatherHistory), weatherHistory$Precip.Type) # as.data.frame.matrix(
q2 <- table(1:nrow(weatherHistory), weatherHistory$Summary)
q <- as.data.frame.matrix(cbind(q1, q2))
#head(merge(PT, FD))

#oh_weatherHistory <- dummyVars("~ .", data = qualitative_wH)
#oh_weatherHistory <- data.frame(predict(oh_weatherHistory, newdata = qualitative_wH))
#head(num_stand_wH)
cleaned_wH <- cbind(num_stand_wH, q)
head(cleaned_wH)

## Temperature..C. Apparent.Temperature..C. Humidity Wind.Speed..km.h.
## 1      -0.2575977      -0.3240338 0.7934663      0.47863251
## 2      -0.2698121      -0.3390953 0.6399922      0.49959129
## 3      -0.2674856      -0.1381015 0.7934663      -0.99546821
## 4      -0.3814869      -0.4590684 0.4865181      0.47630376
## 5      -0.3326292      -0.3624667 0.4865181      0.03384067
## 6      -0.2837715      -0.3500020 0.5888342      0.45534498
## Wind.Bearing..degrees. Visibility..km. Pressure..millibars. null rain snow
## 1      0.5912529      1.306969      0.1016847      0      1      0
## 2      0.6657523      1.306969      0.1059593      0      1      0
## 3      0.1535690      1.099580      0.1086095      0      1      0
## 4      0.7588766      1.306969      0.1126276      0      1      0
## 5      0.6657523      1.306969      0.1134826      0      1      0
## 6      0.6564399      1.099580      0.1147649      0      1      0
## Breezy Breezy and Dry Breezy and Foggy Breezy and Mostly Cloudy
## 1      0      0      0      0
## 2      0      0      0      0
## 3      0      0      0      0
## 4      0      0      0      0
## 5      0      0      0      0
## 6      0      0      0      0
## Breezy and Overcast Breezy and Partly Cloudy Clear
## 1      0      0      0
```



```

## 2          0          0  0
## 3          0          0  0
## 4          0          0  0
## 5          0          0  0
## 6          0          0  0
##  Dangerously Windy and Partly Cloudy Drizzle Dry Dry and Mostly Cloudy
## 1          0          0  0
## 2          0          0  0
## 3          0          0  0
## 4          0          0  0
## 5          0          0  0
## 6          0          0  0
##  Dry and Partly Cloudy Foggy Humid and Mostly Cloudy Humid and Overcast
## 1          0  0          0          0
## 2          0  0          0          0
## 3          0  0          0          0
## 4          0  0          0          0
## 5          0  0          0          0
## 6          0  0          0          0
##  Humid and Partly Cloudy Light Rain Mostly Cloudy Overcast Partly Cloudy Rain
## 1          0          0          0  0          1  0
## 2          0          0          0  0          1  0
## 3          0          0          1  0          0  0
## 4          0          0          0  0          1  0
## 5          0          0          1  0          0  0
## 6          0          0          0  0          1  0
##  Windy Windy and Dry Windy and Foggy Windy and Mostly Cloudy
## 1  0          0          0          0
## 2  0          0          0          0
## 3  0          0          0          0
## 4  0          0          0          0
## 5  0          0          0          0
## 6  0          0          0          0
##  Windy and Overcast Windy and Partly Cloudy
## 1          0          0
## 2          0          0
## 3          0          0
## 4          0          0
## 5          0          0
## 6          0          0

```

```

sample <- sample(c(T, F), nrow(cleaned_wH), replace=T, prob=c(0.75, 0.25))
test_wH <- cleaned_wH[!sample,]
train_wH <- cleaned_wH[sample,]

```

(c)

Reason for choosen variables:

- Tempratrue (C) : Baseline that gets moved
- Humidity : Feels a lot hotter when its more humid, harder to sweat
- Wind speed : Wind makes skin feel colder
- Pressure : Pressure changes based on if it may rain or not, feels different
- Rain/Snow : If it rains the air feels colder

```
wH_mod <- train_wH %>%
  lm(Apparent.Temperature..C. ~ rain + snow + Pressure..millibars. + Humidity + Temperature..C. + Wind..km.h.)

summary.aov(wH_mod)
```

```
##              Df Sum Sq Mean Sq  F value Pr(>F)
## rain          1  22292    22292  2.216e+06 <2e-16 ***
## snow          1    750      750  7.457e+04 <2e-16 ***
## Pressure..millibars. 1      1      1  7.374e+01 <2e-16 ***
## Humidity       1  16920    16920  1.682e+06 <2e-16 ***
## Temperature..C.  1  30885    30885  3.070e+06 <2e-16 ***
## Wind.Speed..km.h.  1    254      254  2.528e+04 <2e-16 ***
## Residuals      71999      724        0
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(wH_mod)
```

```
##
## Call:
## lm(formula = Apparent.Temperature..C. ~ rain + snow + Pressure..millibars. +
##      Humidity + Temperature..C. + Wind.Speed..km.h., data = .)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.39827 -0.06843 -0.00989  0.06134  0.46461
##
## Coefficients:
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    0.0405563   0.0052229     7.765 8.27e-15 ***
## rain          -0.0360736   0.0052403    -6.884 5.87e-12 ***
## snow          -0.0752042   0.0053886   -13.956 < 2e-16 ***
## Pressure..millibars. 0.0022826   0.0003748     6.090 1.14e-09 ***
## Humidity       0.0160131   0.0005182    30.899 < 2e-16 ***
## Temperature..C.  0.9967281   0.0005958  1672.983 < 2e-16 ***
## Wind.Speed..km.h. -0.0628036   0.0003950  -159.011 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1003 on 71999 degrees of freedom
## Multiple R-squared:  0.9899, Adjusted R-squared:  0.9899
## F-statistic: 1.178e+06 on 6 and 71999 DF,  p-value: < 2.2e-16
```

RMSE

$$\text{RMSE}(y, \hat{y}) = \sqrt{\frac{\sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2}{N}}$$

Exercise 3 - Linear Regression and Diagnostic Plots

Exercise 4 - correlation and partial correlation