

431 Class 07

thomaselove.github.io/431

2020-09-15

Today's Data

NHANES data from 2011-12: a sample of 1000 adults.

- Fix a problem with the old `nh2` data I'd built in Classes 5 and 6
- Create `nh3` which solves this problem.

How should we explore these data before modeling?

- What can we learn about the center, spread, outliers, and shape of quantitative data?
- Are these blood pressure data well described by a Normal distribution?

Next Time

How might we look at Associations between our quantities?

- Scatterplots, Correlation, Linear Models, Smoothing

Loading our R Packages

```
library(NHANES)
library(janitor)
library(knitr)
library(broom)
library(magrittr)
library(patchwork)
library(tidyverse)

theme_set(theme_bw())
```

Creating the old_nh2 data set (from last week)

```
set.seed(20200908)

old_nh2 <- NHANES %>%
  filter(SurveyYr == "2011_12") %>%
  select(ID, SurveyYr, Age, Height, Weight, BMI, Pulse,
         SleepHrsNight, BPSysAve, BPDiaAve, Gender,
         PhysActive, SleepTrouble, Smoke100,
         Race1, HealthGen, Depressed) %>%
  rename(SleepHours = SleepHrsNight, Sex = Gender,
         SBP = BPSysAve, DBP = BPDiaAve) %>%
  filter(Age > 20 & Age < 80) %>%
  drop_na() %>%
  slice_sample(., n = 1000) %>%
  clean_names()
```

But, there's a nuisance...

```
old_nh2 %>% arrange(id) %>% select(id:weight) %>% head(7)
```

```
# A tibble: 7 x 5
```

	id	survey_yr	age	height	weight
	<int>	<fct>	<int>	<dbl>	<dbl>
1	62180	2011_12	35	179.	89
2	62205	2011_12	28	171.	84.8
3	62205	2011_12	28	171.	84.8
4	62208	2011_12	38	169.	63.2
5	62220	2011_12	31	167.	113.
6	62222	2011_12	32	179	80.1
7	62222	2011_12	32	179	80.1

There are duplicate records in NHANES

To make this sample representative of the US, some subjects appear in the sample multiple times. Suppose we want to look only at distinct rows.

```
nh_deduplicated <- NHANES %>%  
  filter(SurveyYr == "2011_12") %>%  
  select(ID, SurveyYr, Age, Height, Weight, BMI, Pulse,  
         SleepHrsNight, BPSysAve, BPDiaAve, Gender,  
         PhysActive, SleepTrouble, Smoke100,  
         Race1, HealthGen, Depressed) %>%  
  rename(SleepHours = SleepHrsNight, Sex = Gender,  
         SBP = BPSysAve, DBP = BPDiaAve) %>%  
  filter(Age > 20 & Age < 80) %>%  
  drop_na() %>%  
  distinct() # add this to avoid duplicate rows
```

```
dim(nh_deduplicated)
```

```
[1] 1727    17
```

New version of data, without duplications

```
set.seed(20200914)

nh3 <- nh_deduplicated %>%
  slice_sample(n = 1000) %>%
  clean_names()
```

The importance of clean names



Source: <https://github.com/allisonhorst/stats-illustrations>

Today's Questions

- ① How might we explore the blood pressure data to understand it better?
In particular, does a Normal model fit our systolic and diastolic blood pressures well?
- ② What is the nature of the association between systolic BP and diastolic BP in these NHANES subjects?

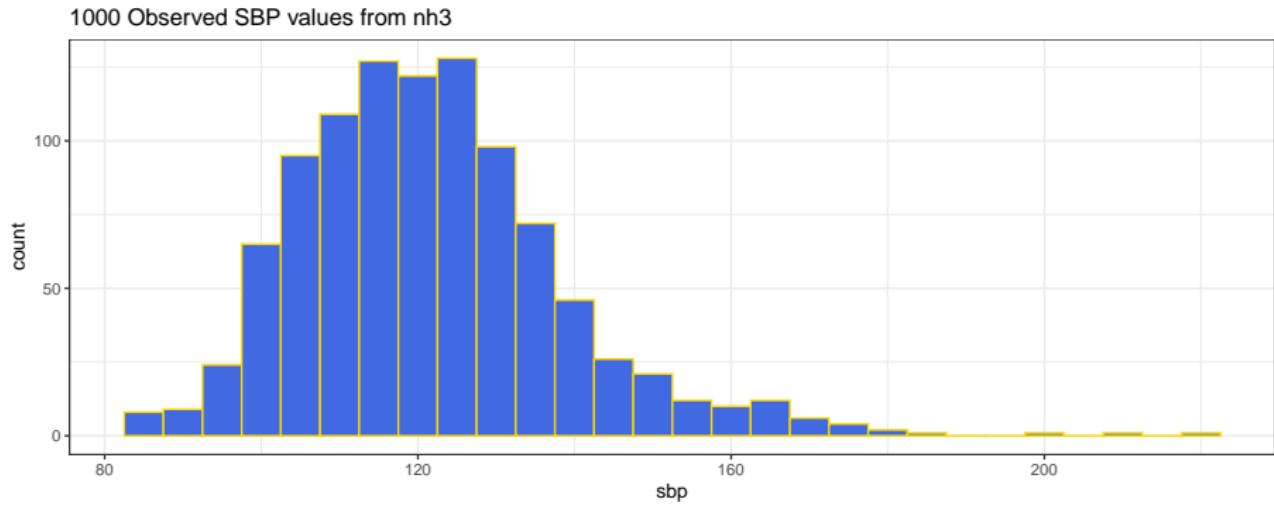
Today's Variables

Name	Description
id	Identifying code for each subject
sbp	Systolic Blood Pressure (mm Hg)
dbp	Diastolic Blood Pressure (mm Hg)

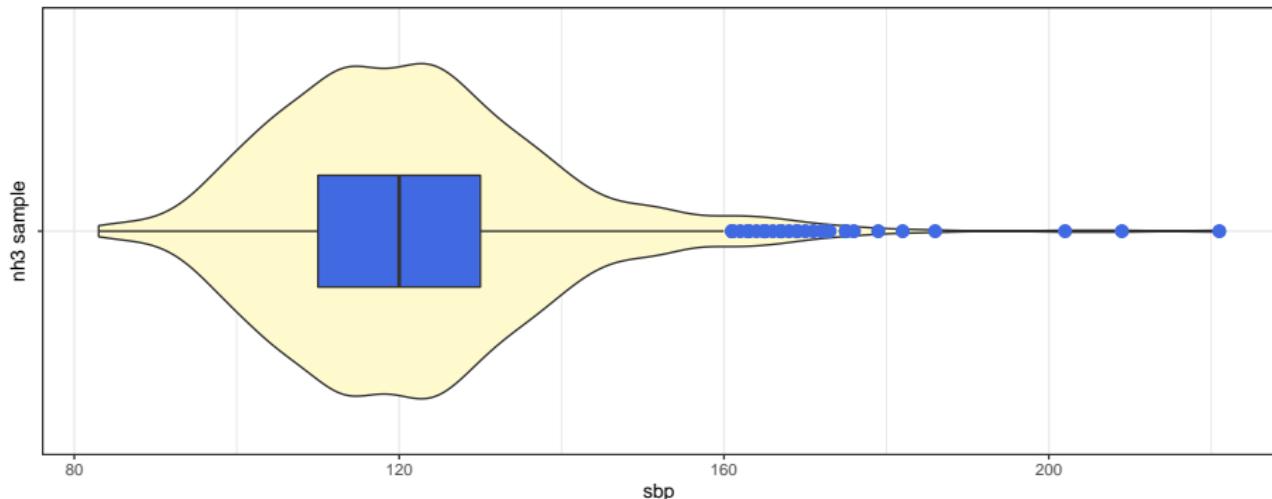
**Plotting the sbp data to learn about center,
spread, outliers, and shape**

Histogram of Systolic BP values from nh3

```
ggplot(data = nh3, aes(x = sbp)) +  
  geom_histogram(binwidth = 5,  
                 fill = "royalblue", col = "gold") +  
  labs(title = "1000 Observed SBP values from nh3")
```



Violin and Boxplot for nh3 SBP data



min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.7	17.1	1000	0

Can we describe these data as being
well-approximated by a Normal model?

What is a Normal Model?

By a Normal model, we mean that the data are assumed to be the result of selecting at random from a probability distribution called the Normal (or Gaussian) distribution, which is characterized by a bell-shaped curve.

- The Normal model is defined by establishing the values of two parameters: the mean and the standard deviation.

When is it helpful to assume our data follow a Normal model?

- When summarizing the data (especially if we want to interpret the mean and standard deviation)
- When creating inferences about populations from samples (as in a t test, or ANOVA)
- When creating regression models, it will often be important to make distributional assumptions about errors, for instance, that they follow a Normal model.

Does a Normal model fit our data “well enough”?

We evaluate whether a Normal model fits sufficiently well to our data on the basis of (in order of importance):

- ① Graphs (DTDP) are the most important tool we have
 - There are several types of graphs available that are designed to (among other things) help us identify clearly several of the potential problems with assuming Normality.
- ② Planned analyses after a Normal model decision is made
 - How serious the problems we see in graphs need to be before we worry about them changes substantially depending on how closely the later analyses we plan to do rely on the assumption of Normality.
- ③ Numerical Summaries are by far the least important even though they seem “easy-to-use” and “objective”.

Does a Normal model fit well for my data?

The least important approach (even though it is seemingly the most objective) is the calculation of various numerical summaries.

Semi-useful summaries help us understand whether they match up well with the expectations of a normal model:

- ① Assessing skewness with $skew_1$ (is the mean close to the median)?
- ② Assessing coverage probabilities:
 - In a Normal model, mean ± 1 standard deviation covers 68% of the data.
 - In a Normal model, mean ± 2 standard deviations covers 95% of the data.
 - In a Normal model, mean ± 3 standard deviations covers 99.7% of the data.

Quantifying skew with a simple $skew_1$ measure

$$skew_1 = \frac{mean - median}{standard\ deviation}$$

Interpreting $skew_1$ (for unimodal data)

- $skew_1 = 0$ if the mean and median are the same
- $skew_1 > 0.2$ indicates fairly substantial right skew
- $skew_1 < -0.2$ indicates fairly substantial left skew

Measuring skewness in the SBP values: nh3?

```
mosaic::favstats(~ sbp, data = nh3)
```

min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.725	17.09677	1000	0

```
nh3 %>% summarize(skew1 = (mean(sbp) - median(sbp))/sd(sbp))
```

```
# A tibble: 1 x 1
```

```
skew1
```

```
<dbl>
```

```
1 0.101
```

What does this suggest?

Empirical Rule for a Normal Model

If the data followed a Normal distribution, perfectly, then about:

- 68% of the data would fall within 1 standard deviation of the mean
- 95% of the data would fall within 2 standard deviations of the mean
- 99.7% of the data would fall within 3 standard deviations of the mean

Remember that, regardless of the distribution of the data:

- Half of the data will fall below the median, and half above it.
- Half of the data will fall in the Interquartile Range (IQR).

How many SBPs are within 1 SD of the mean?

```
nh3 %>%
  count(sbp > mean(sbp) - sd(sbp),
        sbp < mean(sbp) + sd(sbp)) %>%
  kable()
```

	sbp > mean(sbp) - sd(sbp)	sbp < mean(sbp) + sd(sbp)	n
FALSE		TRUE	134
TRUE		FALSE	131
TRUE		TRUE	735

How does this compare to the expectation under a Normal model?
Remember that there are 1000 observations in nh3.

SBP and the mean \pm 2 standard deviations rule?

The total sample size here is 1000.

```
nh3 %>%
  count(sbp > mean(sbp) - 2*sd(sbp),
        sbp < mean(sbp) + 2*sd(sbp)) %>%
  kable()
```

sbp > mean(sbp) - 2 * sd(sbp)	sbp < mean(sbp) + 2 * sd(sbp)	n
FALSE	TRUE	8
TRUE	FALSE	44
TRUE	TRUE	948

How does this compare to the expectation under a Normal model?

Hypothesis Testing to assess Normality

Don't. Graphical approaches are **far** better than hypothesis tests. . .

```
nh3 %$% shapiro.test(sbp)
```

Shapiro-Wilk normality test

```
data: sbp  
W = 0.95318, p-value < 2.2e-16
```

The very small p value (2.2×10^{-16} should be interpreted by you as meaning zero) indicates that the test finds some indications **against** adopting a Normal model for these data.

- Exciting, huh? But not actually all that useful, alas.

Why not test for Normality?

There are multiple hypothesis testing schemes (Kolmogorov-Smirnov, etc.) and each looks for one specific violation of a Normality assumption. None can capture the wide range of issues our brains can envision, and none by itself is great at its job.

- With any sort of reasonable sample size, the test is so poor at detecting non-normality compared to our eyes, that it finds problems we don't care about and ignores problems we do care about.
- And without a reasonable sample size, the test is essentially useless.

Whenever you *can* avoid hypothesis testing and instead actually plot the data, you should plot the data. Sometimes you can't plot (especially with really big data) but the test should be a last resort.

Can we simulate a Normal model for SBP?

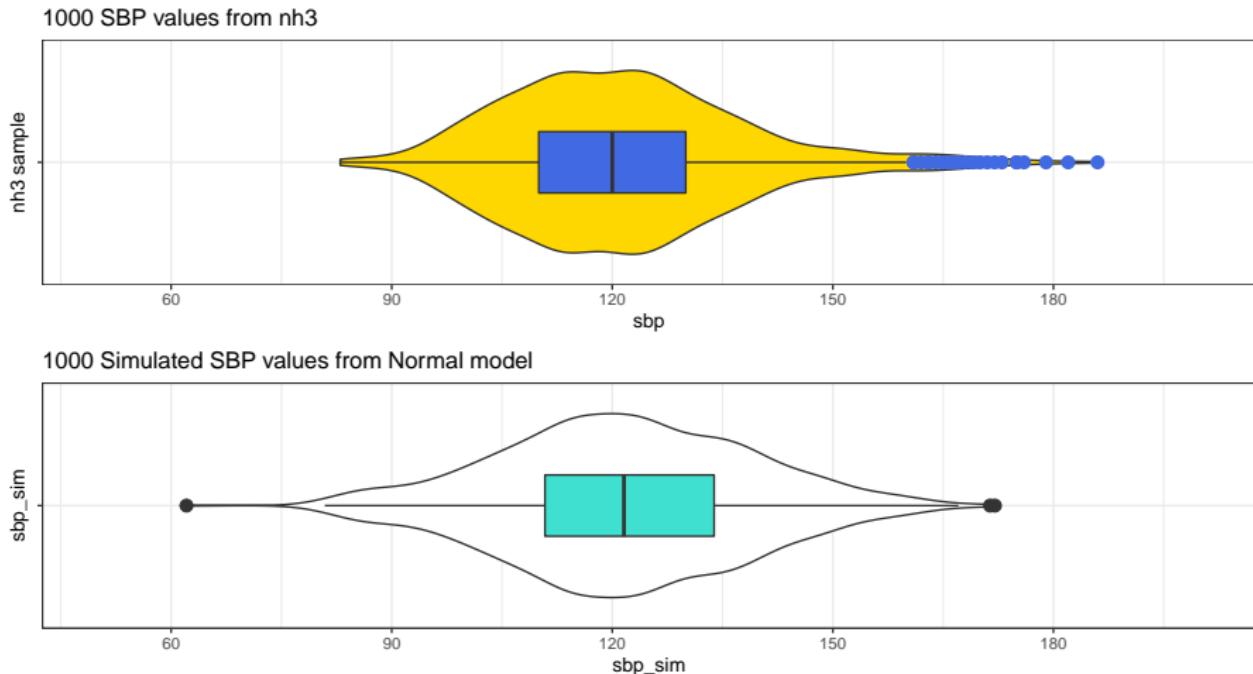
Simulate 1000 observations from a Normal distribution with the same mean and standard deviation as our nh3 systolic BP values...

```
set.seed(1234567)
sim1 <- tibble( sbp_sim =
                  rnorm(n = 1000, mean = 121.7, sd = 17.1))

mosaic::favstats(~ sbp_sim, data = sim1) %>%
  kable(digits = 1)
```

min	Q1	median	Q3	max	mean	sd	n	missing
62.1	110.8	121.6	133.9	172	121.6	17	1000	0

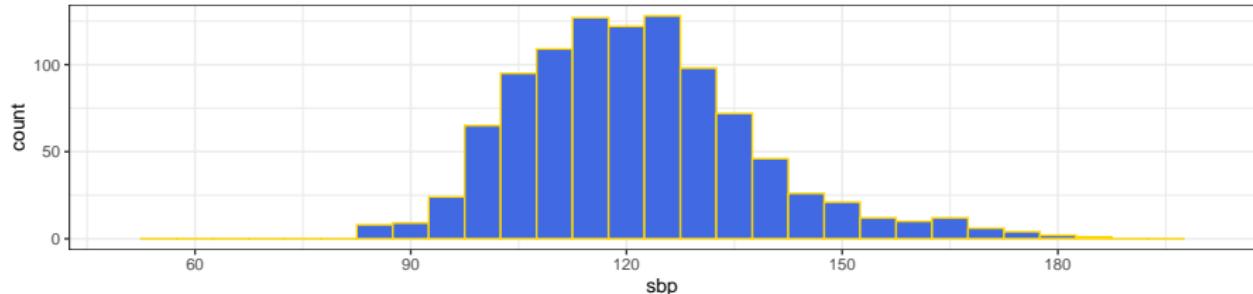
Comparing Boxplots of nh3 and simulated SBP



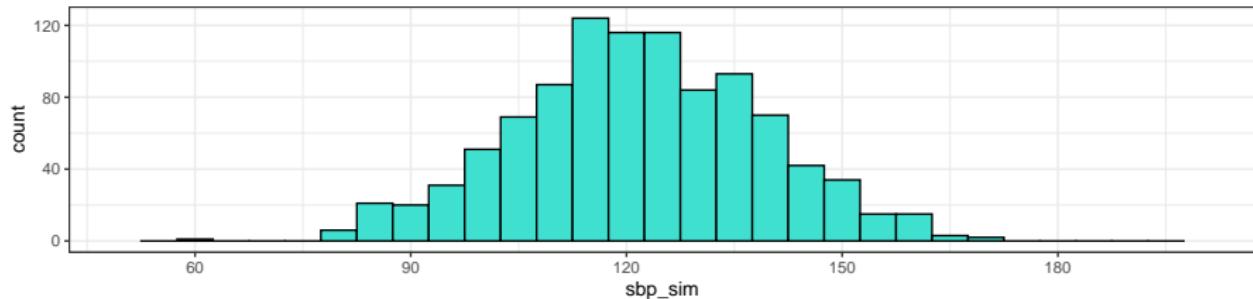
- Does a Normal model look appropriate for describing the nh3 SBP?

Comparing Histograms of nh3 and simulated SBP

1000 Observed SBP values from nh3 (sample mean = 121.7, sd = 17.1)

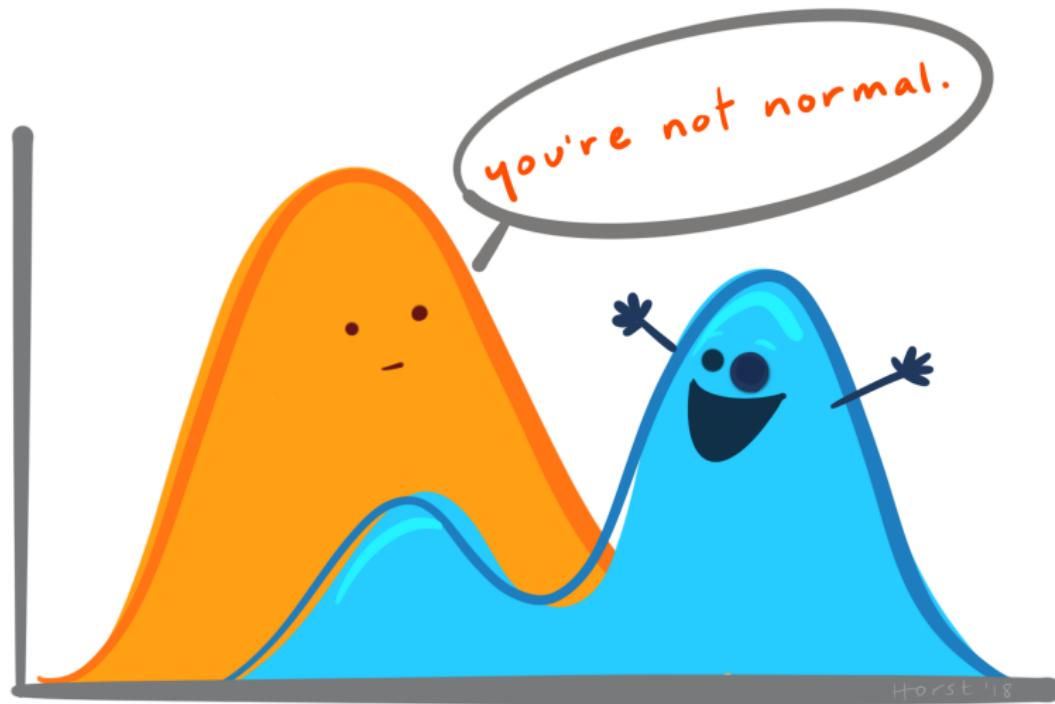


1000 Simulated Values from Normal model with mean = 121.7, sd = 17.1



- Does a Normal model look appropriate for describing the nh3 SBP?

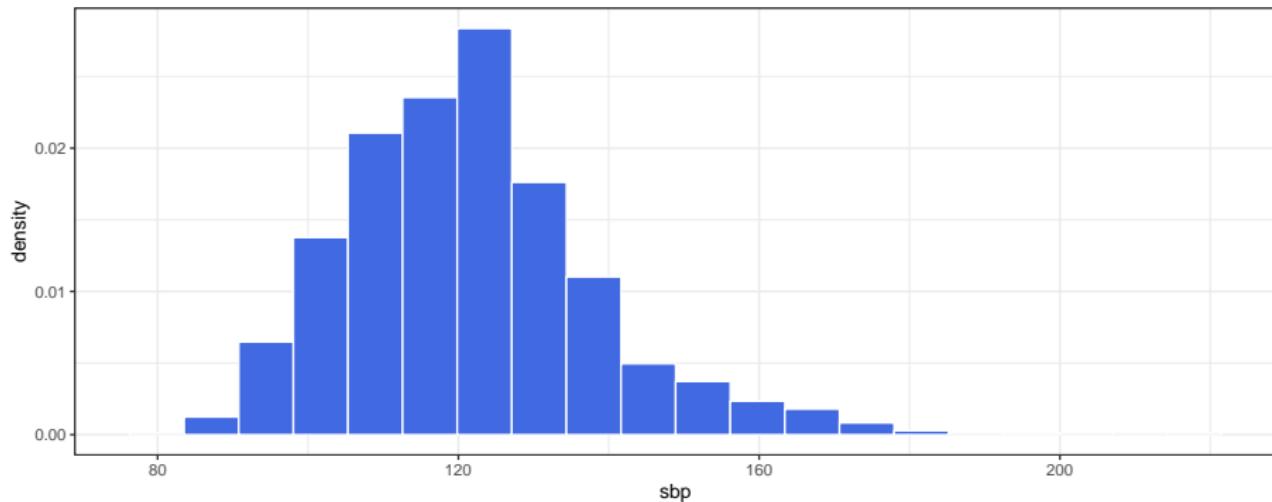
Graphs are our most important tool!



Rescale nh3 SBP histogram as density

Suppose we want to rescale the histogram counts so that the bar areas integrate to 1. This will let us overlay a Normal density onto the results.

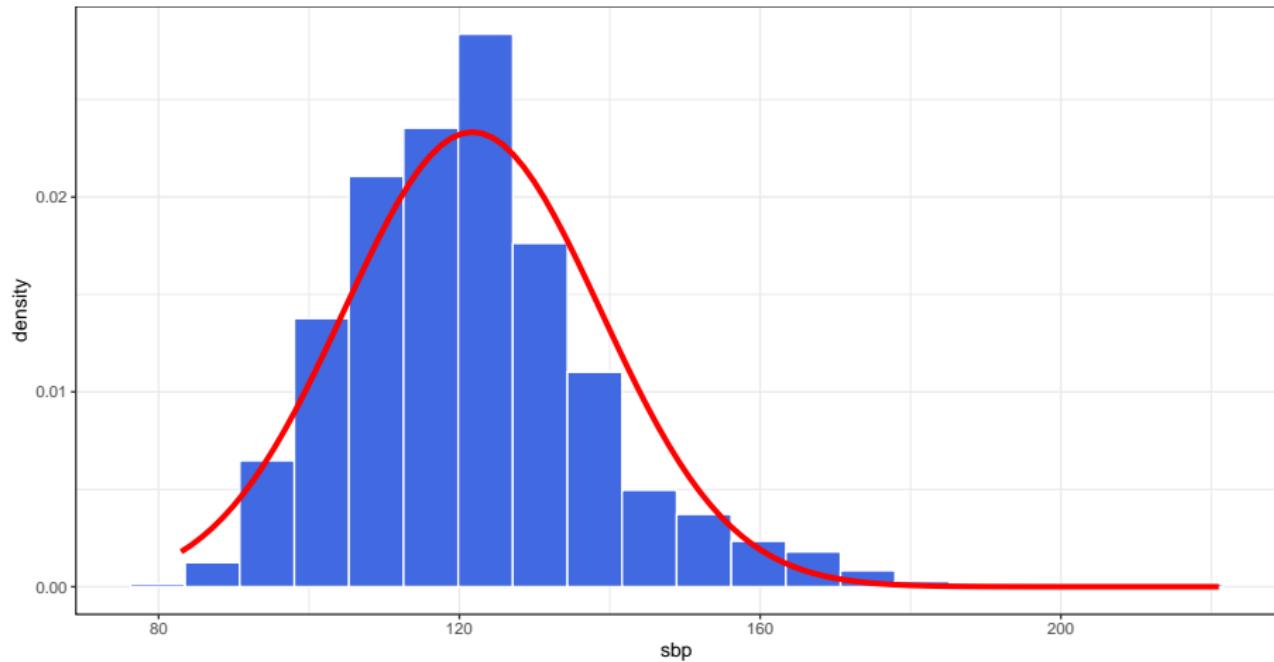
```
ggplot(nh3, aes(x = sbp)) +  
  geom_histogram(aes(y = stat(density)), bins = 20,  
                 fill = "royalblue", col = "white")
```



Density Function, with Normal superimposed

Now we can draw a Normal density curve on top of the rescaled histogram.

SBP density, with Normal model superimposed



Code for plotting Histogram as Density function

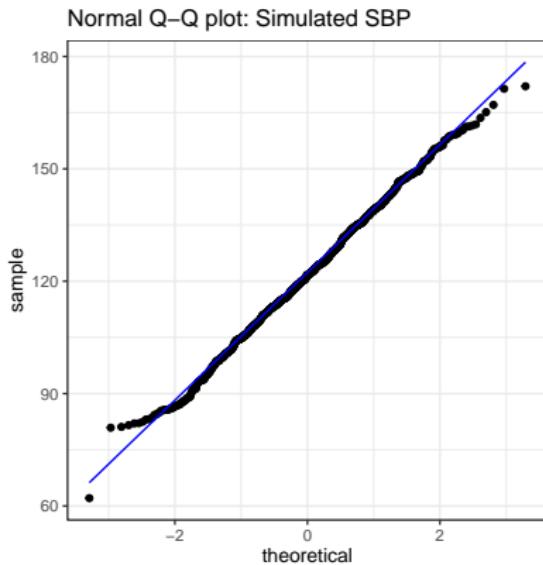
Including the superimposition of a Normal density on top of the histogram.

```
ggplot(nh3, aes(x = sbp)) +
  geom_histogram(aes(y = stat(density)),
                 bins = 20,
                 fill = "royalblue", col = "white") +
  stat_function(fun = dnorm,
                args = list(mean = mean(nh3$sbp),
                            sd = sd(nh3$sbp)),
                col = "red", lwd = 1.5) +
  labs(title = "SBP density, with Normal model superimposed")
```

Using a Normal Q-Q plot

Normal Q-Q plot of our simulated data

Remember that these are draws from a Normal distribution, so this is what a sample of 1000 Normally distributed data points should look like.



The Normal Q-Q Plot

Tool to help assess whether the distribution of a single sample is well-modeled by the Normal.

- Suppose we have N data points in our sample.
- Normal Q-Q plot will plot N points, on a scatterplot.
 - Y value is the data value
 - X value is the expected value for that point in a Normal distribution

Using the Normal distribution with the same mean and SD as our sample, R calculates what the minimum value is expected to be, given a sample of size N , then the next smallest value, and so forth all the way up until the maximum value.

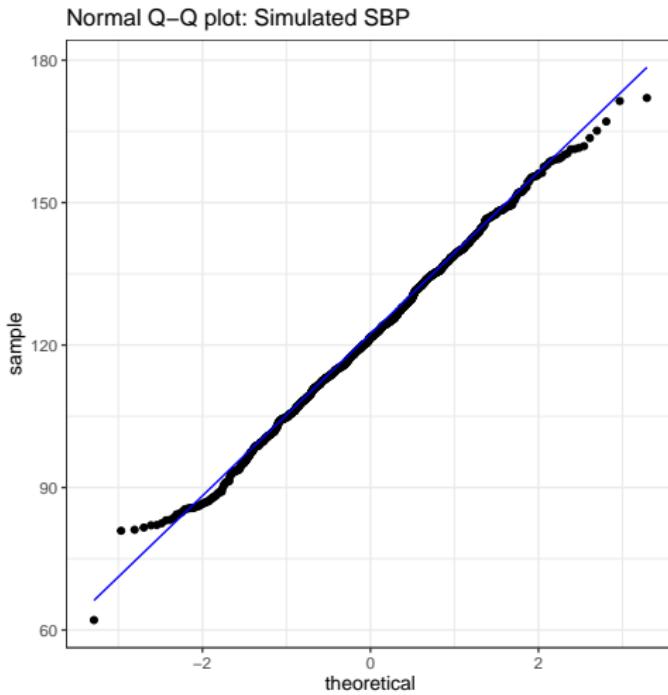
- X value in the Normal Q-Q plot is the value that a Normal distribution would take for that rank in the data set.
- We draw a line through $Y = X$, and points close to the line therefore match what we'd expect from a Normal distribution.

How do we create a Normal Q-Q plot?

For our simulated blood pressure data

```
ggplot(sim1, aes(sample = sbp_sim)) +  
  geom_qq() + # plot the points  
  geom_qq_line(col = "blue") + # plot the Y = X line  
  theme(aspect.ratio = 1) + # make the plot square  
  labs(title = "Normal Q-Q plot: Simulated SBP")
```

Result, again...



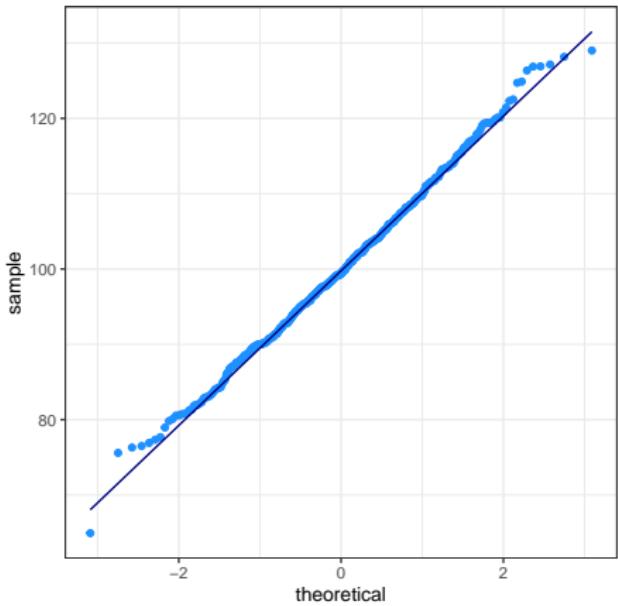
Interpreting the Normal Q-Q plot?

The Normal Q-Q plot can help us identify data as well approximated by a Normal distribution, or not, because of:

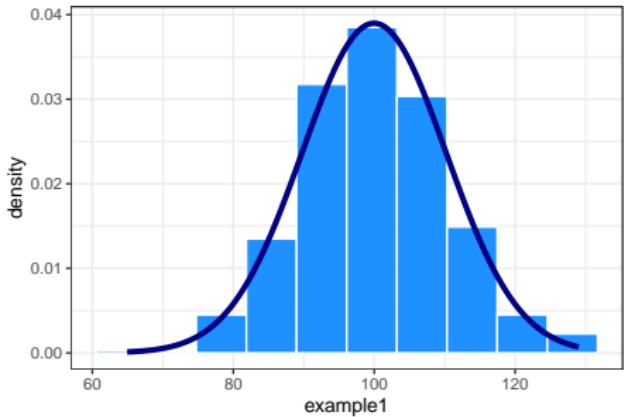
- skew (including distinguishing between right skew and left skew)
 - behavior in the tails (which could be heavy-tailed [more outliers than expected] or light-tailed)
- ① Normally distributed data are indicated by close adherence of the points to the diagonal reference line.
 - ② Skew is indicated by substantial curving (on both ends of the distribution) in the points away from the reference line (if both ends curve up, we have right skew; if both ends curve down, this indicates left skew)
 - ③ An abundance or dearth of outliers (as compared to the expectations of a Normal model) are indicated in the tails of the distribution by an “S” shape or reverse “S” shape in the points.

Example 1: Data from a Normal Distribution

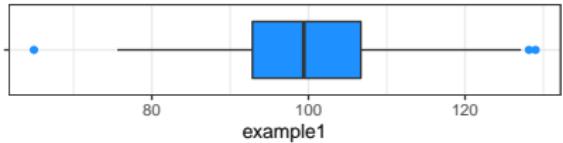
Normal Q-Q plot: Example 1



Density Function: Example 1



Boxplot: Example 1



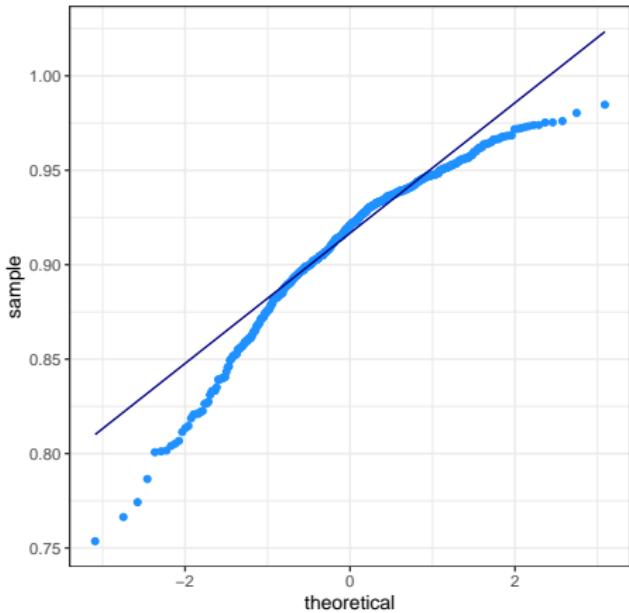
min	Q1	median	Q3	max	mean	sd	n	missing
64.9	92.8	99.4	106.7	129	100	10.2	500	0

Does a Normal model fit well for my data?

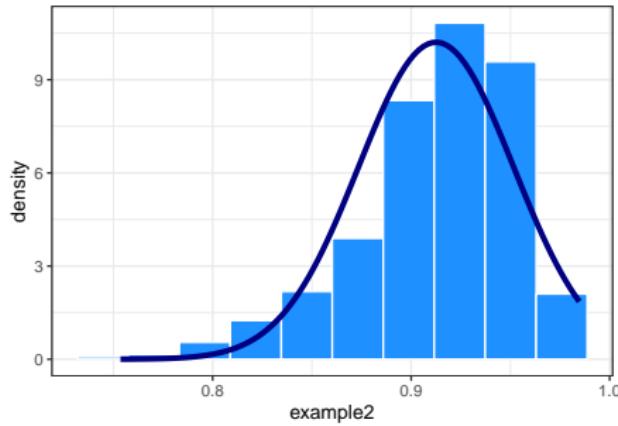
- ① Is a Normal Q-Q plot showing something close to a straight line, without clear signs of skew or indications of lots of outliers (heavy-tailedness)?
- ② Does a boxplot, violin plot and/or histogram also show a symmetric distribution, where both the number of outliers is modest, and the distance of those outliers from the mean is modest?
- ③ Do numerical measures match up with the expectations of a normal model?

Example 2: Data from a Left-Skewed Distribution

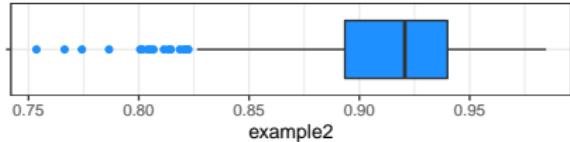
Normal Q-Q plot: Example 2



Density Function: Example 2



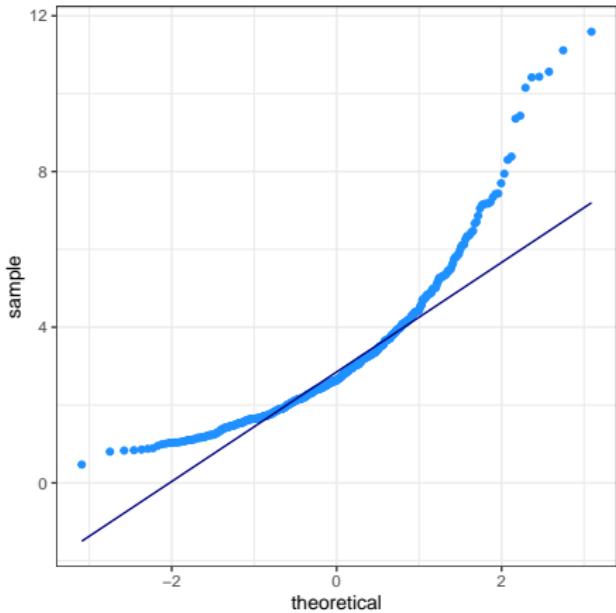
Boxplot: Example 2



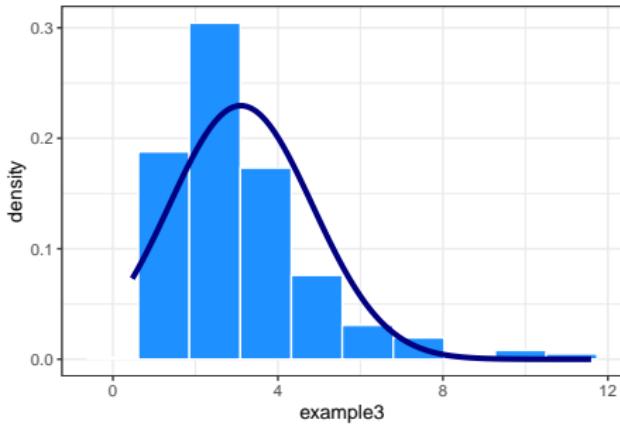
	min	Q1	median	Q3	max	mean	sd	n	missing
	0.8	0.9	0.9	0.9	1	0.9	0	500	0

Example 3: Data from a Right-Skewed Distribution

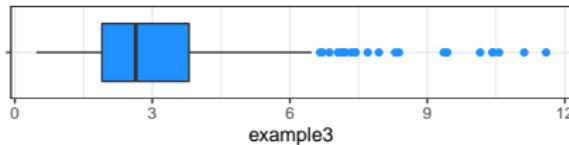
Normal Q-Q plot: Example 3



Density Function: Example 3



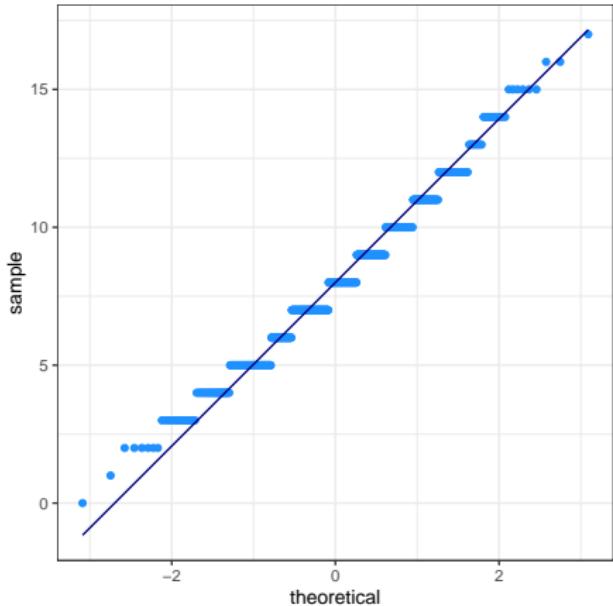
Boxplot: Example 3



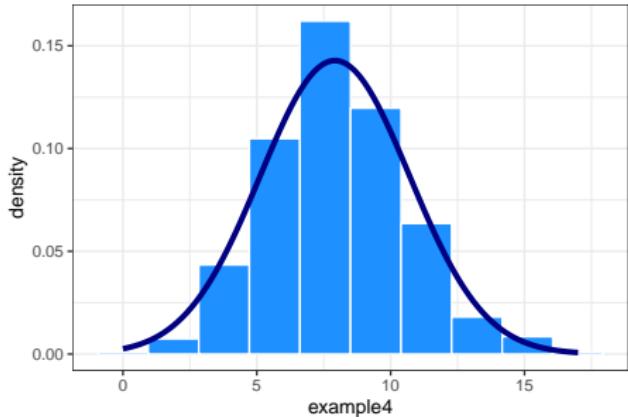
min	Q1	median	Q3	max	mean	sd	n	missing
0.5	1.9	2.6	3.8	11.6	3.1	1.7	500	0

Example 4: Discrete “Symmetric” Distribution

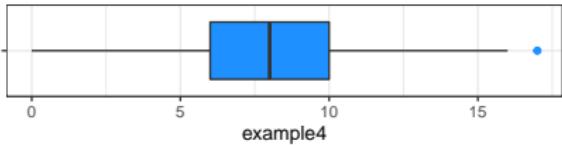
Normal Q-Q plot: Example 4



Density Function: Example 4



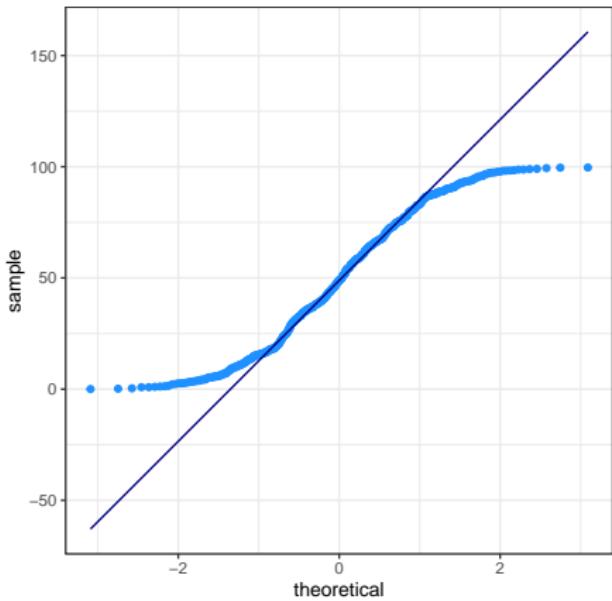
Boxplot: Example 4



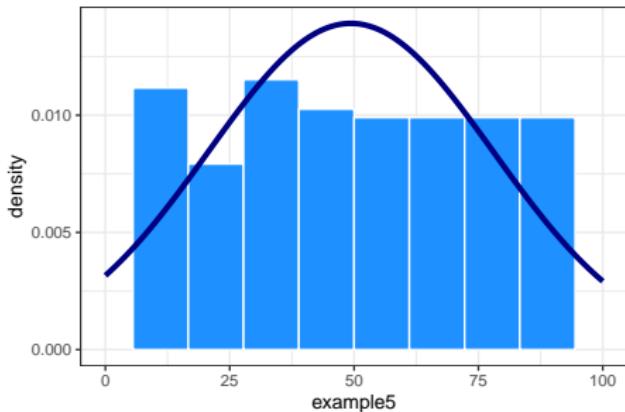
	min	Q1	median	Q3	max	mean	sd	n	missing
	0	6	8	10	17	7.9	2.8	500	0

Example 5: Data from a Uniform Distribution

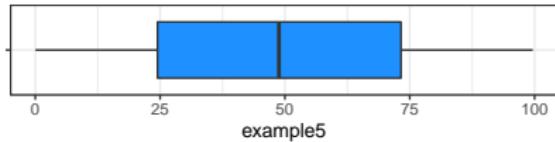
Normal Q-Q plot: Example 5



Density Function: Example 5



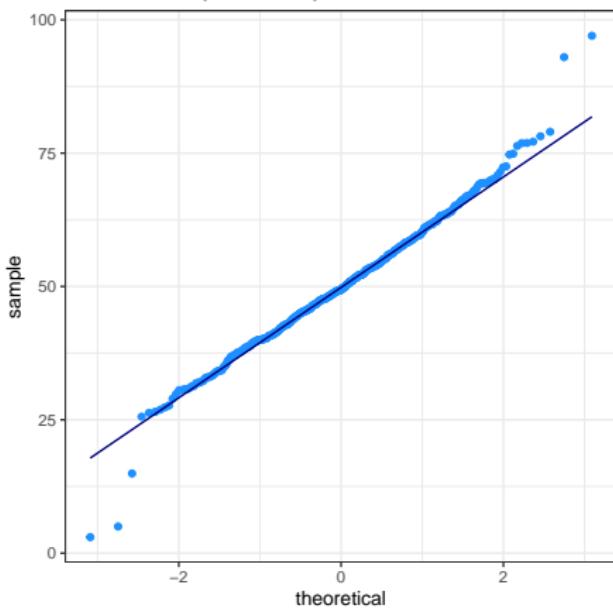
Boxplot: Example 5



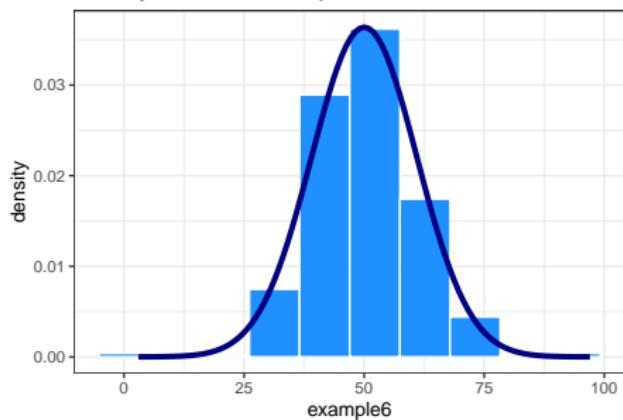
min	Q1	median	Q3	max	mean	sd	n	missing
0	24.5	48.8	73.2	99.6	49.3	28.7	500	0

Example 6: Symmetric data with outliers

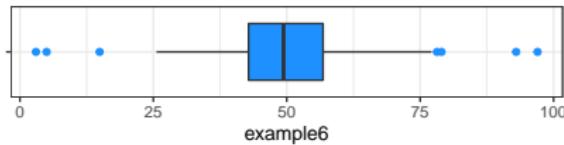
Normal Q-Q plot: Example 6



Density Function: Example 6

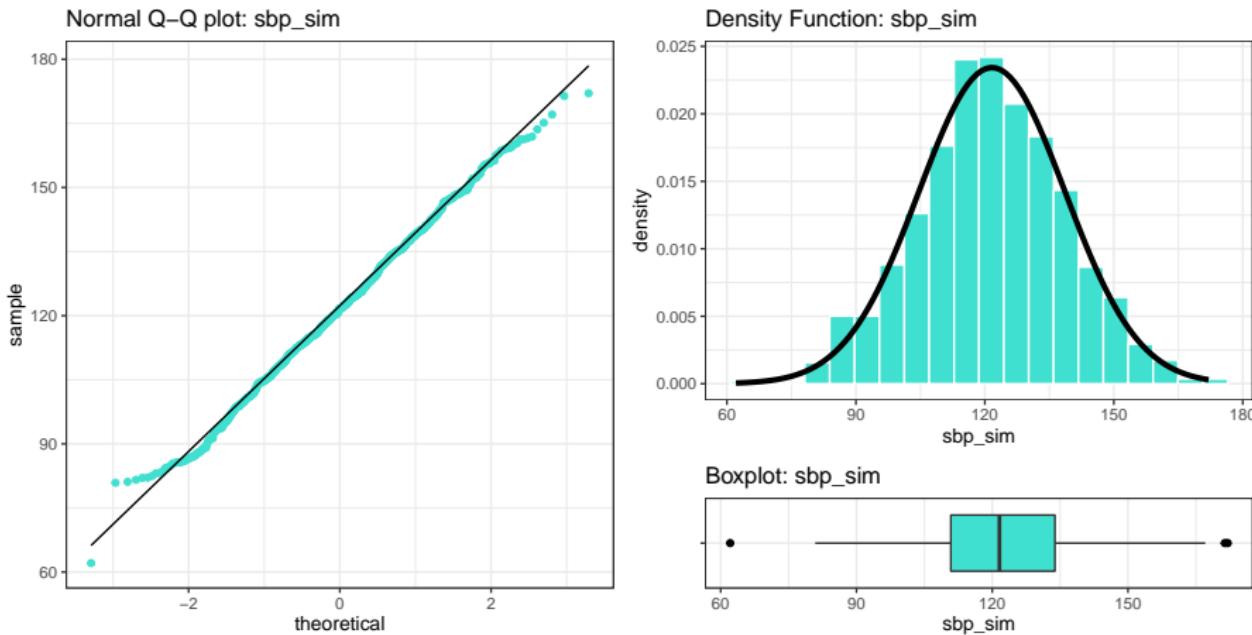


Boxplot: Example 6



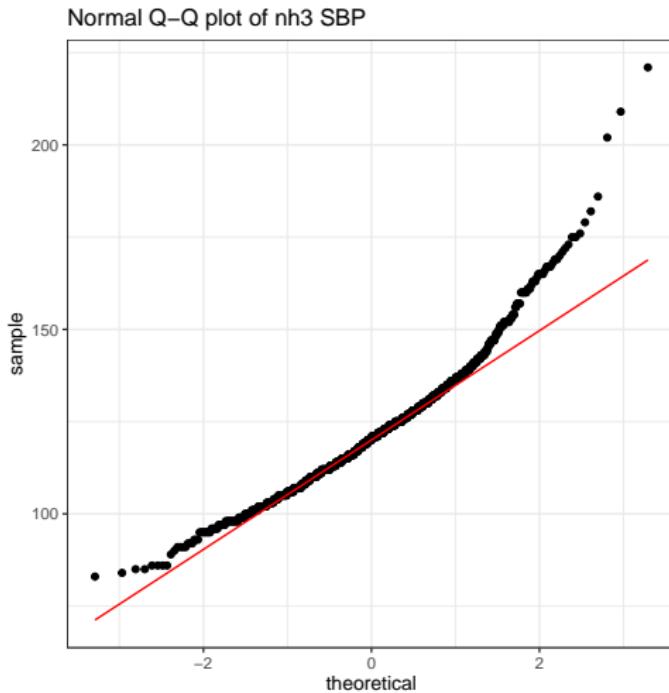
min	Q1	median	Q3	max	mean	sd	n	missing
3	42.8	49.4	56.8	97	50	11	500	0

Our 1000 simulated Systolic Blood Pressures



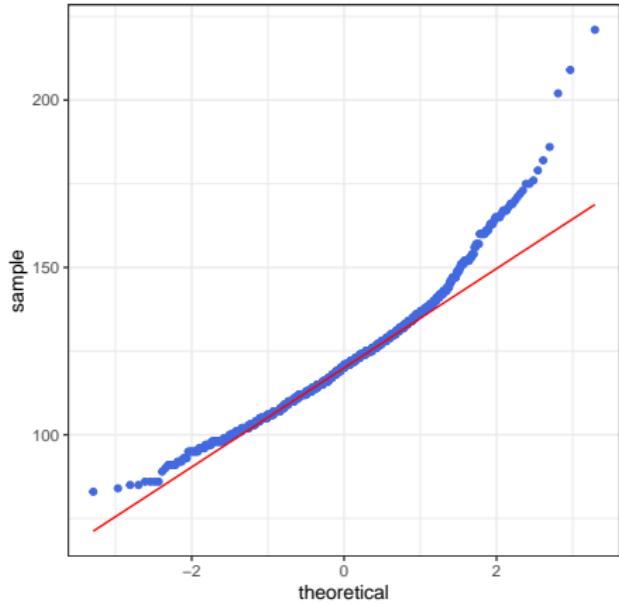
min	Q1	median	Q3	max	mean	sd	n	missing
62.1	110.8	121.6	133.9	172	121.6	17	1000	0

A Normal Q-Q Plot of the nh3 SBP data ($n = 1000$)

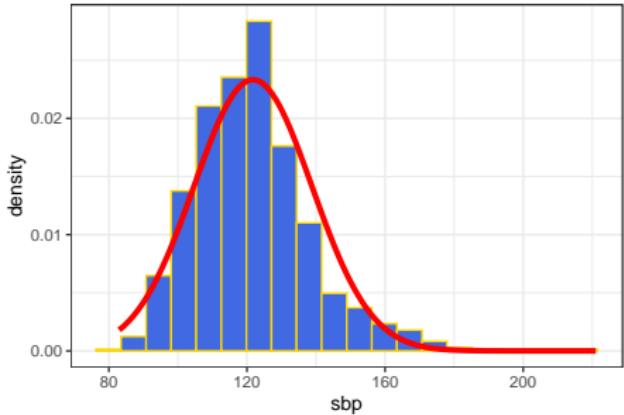


How do we build this slide?

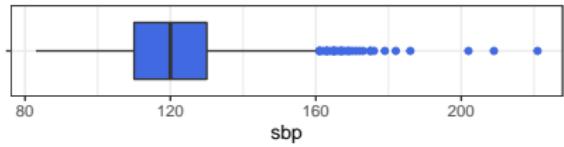
Normal Q-Q plot: nh3 SBP



Density Function: nh3 SBP



Boxplot: nh3 SBP



min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.7	17.1	1000	0

Code for sbp in nh3 (First of Three Plots)

```
p1 <- ggplot(nh3, aes(sample = sbp)) +
  geom_qq(col = "royalblue") +
  geom_qq_line(col = "red") +
  theme(aspect.ratio = 1) +
  labs(title = "Normal Q-Q plot: nh3 SBP")
```

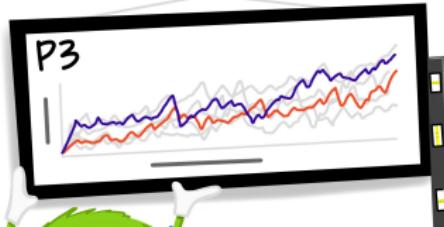
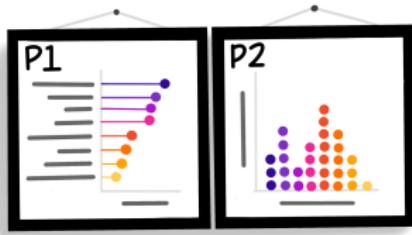
Code for sbp in nh3 (Second of Three Plots)

```
p2 <- ggplot(nh3, aes(x = sbp)) +
  geom_histogram(aes(y = stat(density)),
                 bins = 20,
                 fill = "royalblue", col = "gold") +
  stat_function(fun = dnorm,
                args = list(mean = mean(nh3$sbp),
                            sd = sd(nh3$sbp)),
                col = "red", lwd = 1.5) +
  labs(title = "Density Function: nh3 SBP")
```

Code for sbp in nh3 (Third of Three Plots)

```
p3 <- ggplot(nh3, aes(x = sbp, y = "")) +  
  geom_boxplot(fill = "royalblue",  
               outlier.color = "royalblue") +  
  labs(title = "Boxplot: nh3 SBP", y = "")
```

Putting the plots together...



Using patchwork

```
p1 + (p2 / p3 + plot_layout(heights = c(4,1)))
```

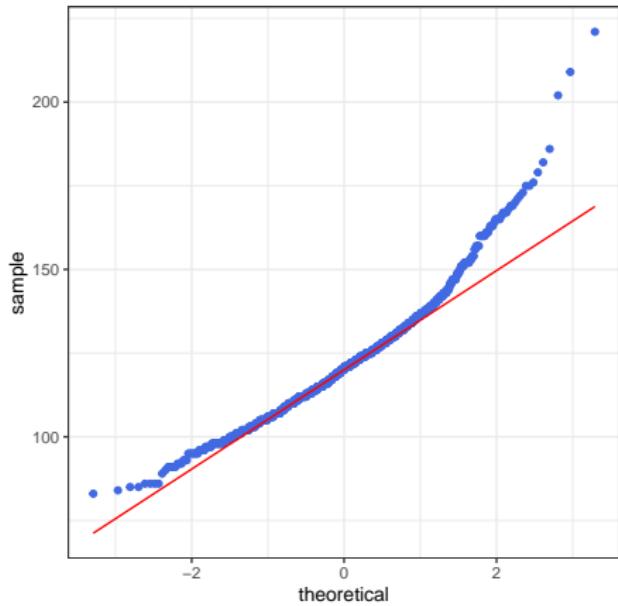
Also added...

```
mosaic::favstats(~ sbp, data = nh3) %>% kable(digits = 1)
```

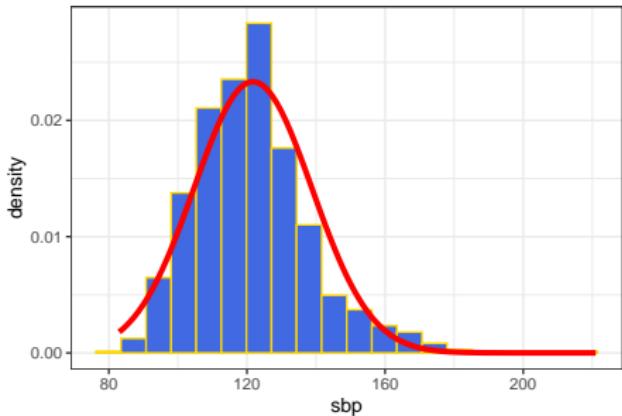
min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.7	17.1	1000	0

Result: 1000 observed Systolic BP values

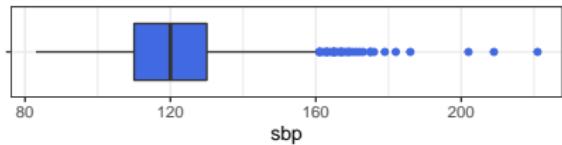
Normal Q-Q plot: nh3 SBP



Density Function: nh3 SBP



Boxplot: nh3 SBP



min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.7	17.1	1000	0

More Extensive Numerical Summaries?

We could try

```
nh3 %>% Hmisc::describe(sbp)
```

- but that will throw an error message, specifically Error in
describe.data.frame(., sbp) : object 'sbp' not found.
What is wrong? How can we fix that?

More Extensive Numerical Summaries?

We could try

```
nh3 %>% Hmisc::describe(sbp)
```

- but that will throw an error message, specifically Error in
describe.data.frame(., sbp) : object 'sbp' not found.
What is wrong? How can we fix that?
- We could drop the pipe and use \$ notation, so
`Hmisc::describe(nh3$sbp)`

More Extensive Numerical Summaries?

We could try

```
nh3 %>% Hmisc::describe(sbp)
```

- but that will throw an error message, specifically Error in
describe.data.frame(., sbp) : object 'sbp' not found.
What is wrong? How can we fix that?
- We could drop the pipe and use \$ notation, so
`Hmisc::describe(nh3$sbp)`
- Another option is to change the pipe (to the %\$% pipe available in the
`magrittr` package): `nh3 %$% Hmisc::describe(sbp)`

What do these summaries tell us?

```
nh3 %$% Hmisc::describe(sbp)
```

sbp

	n	missing	distinct	Info	Mean	Gmd
1000		0	93	1	121.7	18.52
.05		.10	.25	.50	.75	.90
98		102	110	120	130	142
.95						
152						

lowest : 83 84 85 86 89, highest: 182 186 202 209 221

- Gmd = Gini's mean difference (a robust measure of spread) = mean absolute difference between any pairs of observations. Larger Gmd indicates more spread.
- Info = a measure of relative information describing how “continuous” the data are. Higher Info indicates fewer ties.

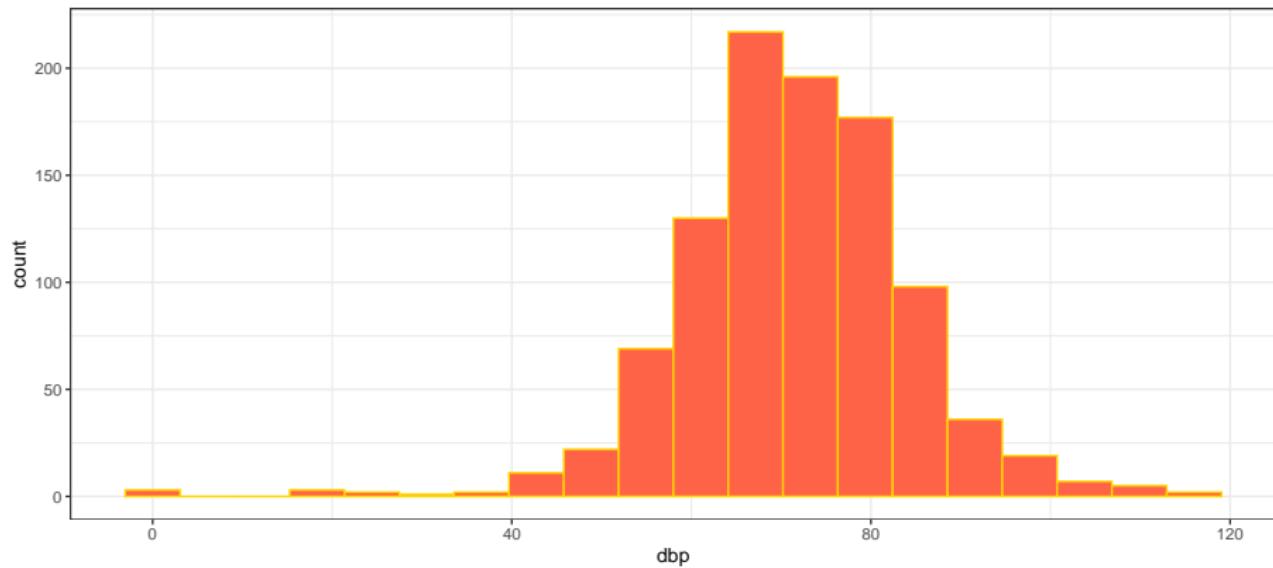
What Summaries to Report

It is usually helpful to focus on the shape, center and spread of a distribution. Bock, Velleman and DeVeaux provide some useful advice:

- If the data are skewed, report the median and IQR (or the three middle quantiles). You may want to include the mean and standard deviation, but you should point out why the mean and median differ. The fact that the mean and median do not agree is a sign that the distribution may be skewed. A histogram will help you make that point.
- If the data are symmetric, report the mean and standard deviation, and possibly the median and IQR as well.
- If there are clear outliers and you are reporting the mean and standard deviation, report them with the outliers present and with the outliers removed. The differences may be revealing. The median and IQR are not likely to be seriously affected by outliers.

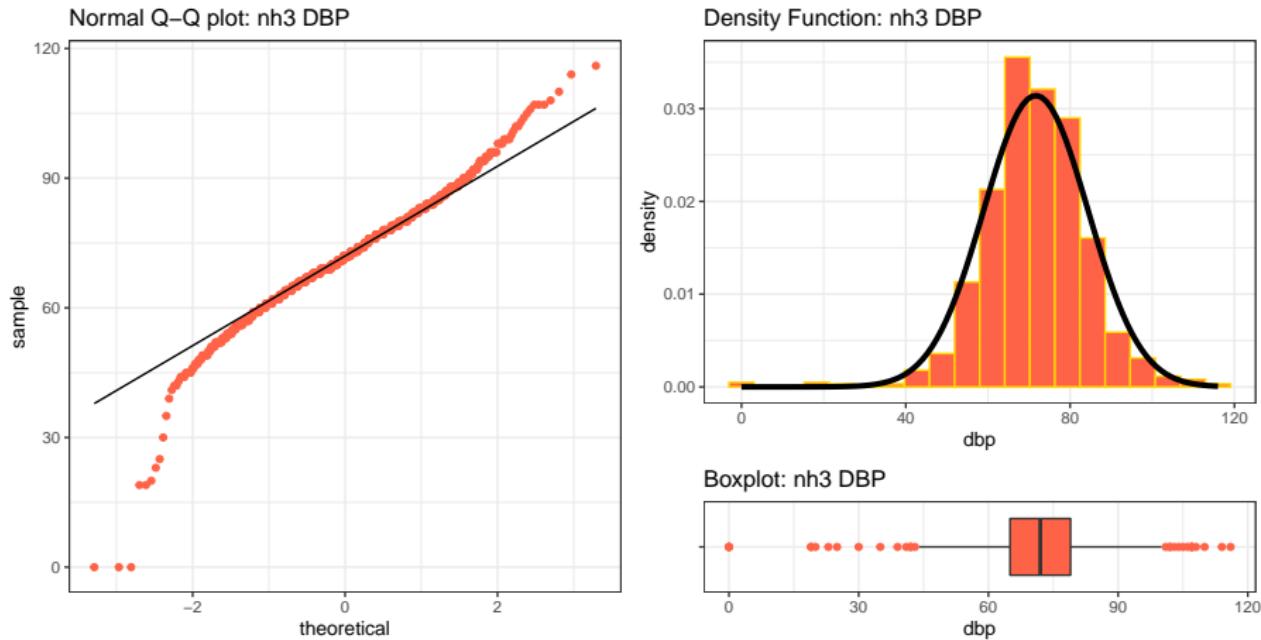
OK, what about Diastolic Blood Pressure?

```
ggplot(data = nh3, aes(x = dbp)) +  
  geom_histogram(bins = 20, fill = "tomato", col = "gold")
```



- We can generate the set of plots we've been using...

DBP in nh3: Center/Spread/Outliers/Shape?



	min	Q1	median	Q3	max	mean	sd	n	missing
0	65	72	79	116	71.7	12.7	1000	0	

Does a Normal model fit well for my data?

- ① Is a Normal Q-Q plot showing something close to a straight line, without clear signs of skew or indications of lots of outliers (heavy-tailedness)?
- ② Does a boxplot, violin plot and/or histogram also show a symmetric distribution, where both the number of outliers is modest, and the distance of those outliers from the mean is modest?
- ③ Do numerical measures match up with the expectations of a normal model?

Hmisc::describe for dbp?

```
nh3 %$% Hmisc::describe(dbp)
```

dbp

	n	missing	distinct	Info	Mean	Gmd
1000		0	78	0.999	71.66	13.58
.05		.10	.25	.50	.75	.90
52		57	65	72	79	86
.95						
91						

lowest : 0 19 20 23 25, highest: 107 108 110 114 116

What is a plausible diastolic blood pressure?

Stem-and-Leaf of dbp values?

`stem(nh3$dbp)`

The decimal point is 1 digit(s) to the right of the |

Who are those people with tiny dbp values?

```
nh3 %>%  
  filter(dbp < 40) %>%  
  select(id, sbp, dbp)
```

```
# A tibble: 11 x 3  
      id    sbp    dbp  
   <int> <int> <int>  
1 71598     86    30  
2 68528    133    39  
3 64298    135     0  
4 64616    111    25  
5 65298    126    35  
6 62649    122    23  
7 70664    152     0  
8 69237    120     0  
9 68561    119    19  
10 68908   129    20
```

Let's reset.

```
nh3_new <- nh3 %>%  
  filter(dbp > 39)
```

```
nrow(nh3)
```

```
[1] 1000
```

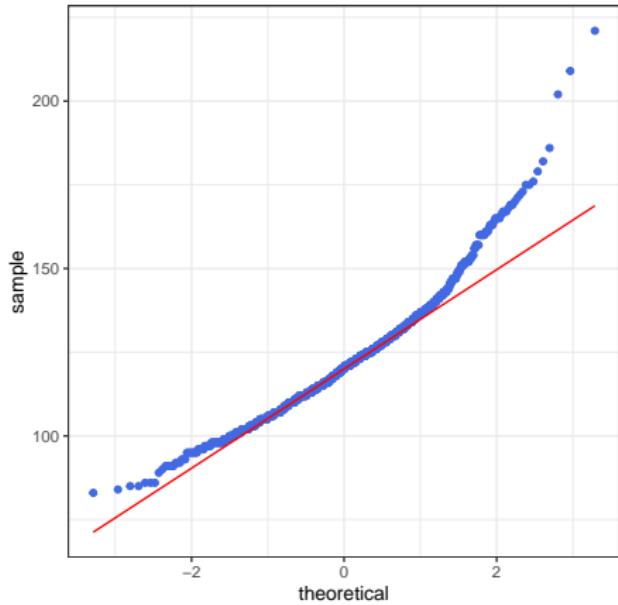
```
nrow(nh3_new)
```

```
[1] 989
```

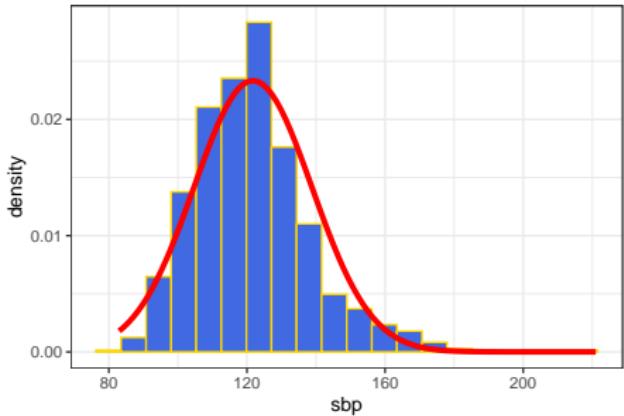
We'll work with nh3_new for the rest of today.

nh3_new: Systolic Blood Pressure

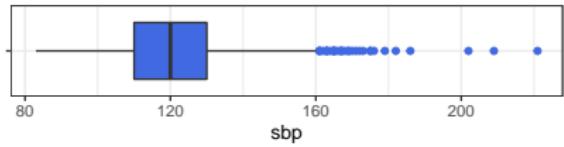
Normal Q-Q plot: nh3 SBP



Density Function: nh3 SBP



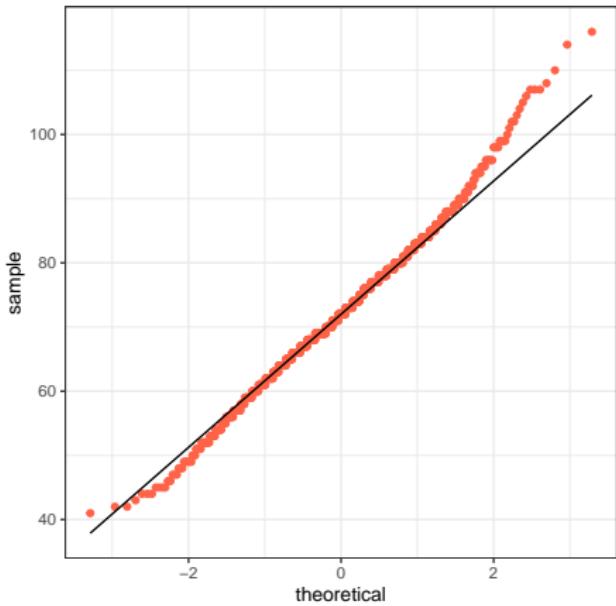
Boxplot: nh3 SBP



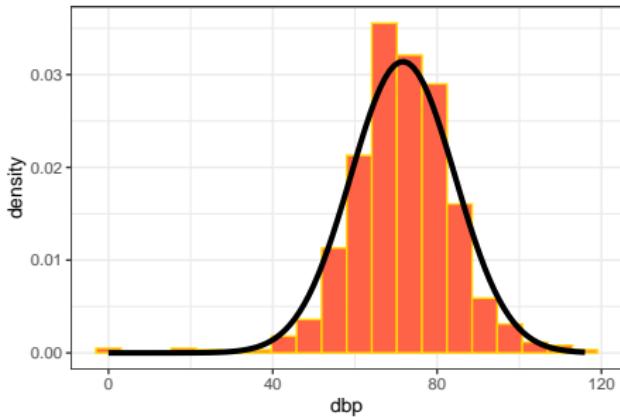
min	Q1	median	Q3	max	mean	sd	n	missing
83	110	120	130	221	121.7	17.1	1000	0

nh3_new: Diastolic Blood Pressure

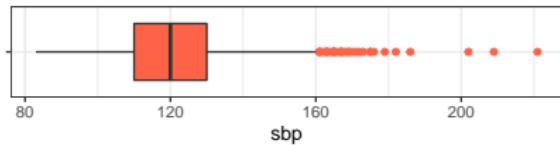
Normal Q-Q plot: nh3 DBP



Density Function: nh3 DBP



Boxplot: nh3 DBP



min	Q1	median	Q3	max	mean	sd	n	missing
0	65	72	79	116	71.7	12.7	1000	0

Summing Up: Does a Normal Model fit well?

If a Normal model fits our data well, then we should see the following graphical indications:

- ① A histogram that is symmetric and bell-shaped.
- ② A boxplot where the box is symmetric around the median, as are the whiskers, without a serious outlier problem.
- ③ A normal Q-Q plot that essentially falls on a straight line.

As for numerical summaries, we'd like to see

- ④ The mean and median within 0.2 standard deviation of each other.
- ⑤ No real evidence of too many outlier candidates (more than 5% starts to get us concerned about a Normal model)
- ⑥ No real evidence of individual outliers outside the reasonable range for the size of our data (we might expect about 3 observations in 1000 to fall more than 3 standard deviations away from the mean.)

NEXT TIME: How can we describe the relationship between SBP and DBP?

Scatterplot to study the SBP-DBP association

```
ggplot(nh3_new, aes(x = dbp, y = sbp)) +  
  geom_point(col = "purple") +  
  theme(aspect.ratio = 1) # make the plot square for slide
```

