

Š 6 W ñ Ž Û

Š 6 W ñ u Ž ' , ž ý v í ½ 6 9 | Ÿ % | Ā W ñ &

36 • Ô c

9 y | Ā

x † ĩ l G L V U F W H í ½ q 9 ' Á l † , ý 7 ö ì | Ÿ † Ÿ †
% : y ,

x l F R Q W Q X R X V ½ í ½ | Ô † L Q I L W H , 6 u 9 q l 6 9 ,
6 9 , % b U D @ J H Á ñ - v ö ½ H f í ½ Ô | â 9 ' Á l ;
A Á Ô - , 4 - Ÿ † p % : y ,

• N W ñ

í ½ l 6 9 P W L F P B E D O H X U T W D Q W Y S H Ž ' , h í ½ | ž ý v í ½
6 9 í ½ / b | Ñ ç Á Á Ô - á ° , 4 | Ì Ç • 9 | † , l
6 9 » • í ½ l 6 9 † È † ĩ l 6 9 ø í ½ l 6 9 | ' v 9 ' , ì | Ì È
ç 9 Ÿ † x
ñ ~ L Q W Y D O : . | , Ÿ ~ f | m
x m U R W ½ > † Ñ ç ² • ... , x ì | 9 W L P H V | > U í ½ l
Ö p ï î Û k ò ĩ , + 3 4

W f l % ĩ l 6 9 Q R Q H W U L F I L F I G D V R D M F H J O W W X L Y Ç l q Ÿ ' f b %
í ½ y 6 9 q , ö ² ý 7 ¼ P U Ÿ † † ,
, ¼ l Q R P L Q D O ' q Ÿ V f f , ý 7 ¼ P U /
; ð l R U Ç Q D O q ; ð , ý ¶ \$ % - ð Ø ø y ò • % o f °

• l f ð f ð ... • l Y f ð
Ö w c T î u j

T î

x m ' =

á ' =

W ĩ ' =

» 9 ¶ & Y V x » 9 ¶ &

LQ §'

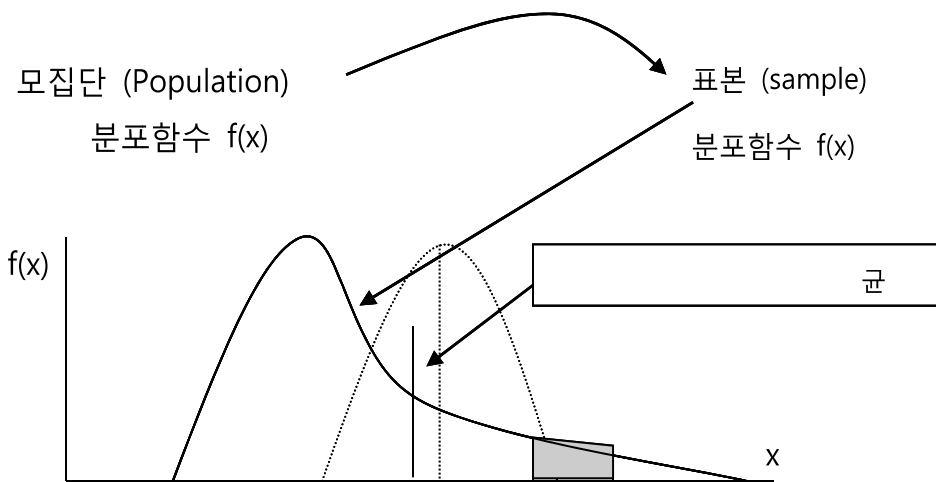
$y \in V^q$; ∂C $\check{Z}'q \text{ ñy } \gg /, \frac{1}{2} ? q \text{ ñ } ,$
 $y \in V^q$; $\partial C x_i \text{ } x_m \text{ } \emptyset z$; $A \beta \{$
 $y \in V^q$, $\emptyset \hat{O} \hat{O} \gg i \frac{1}{2} f \text{ } G I y \text{ } I S \text{ } \check{Y} v \text{ } \partial C$

$> \check{U} \text{ } > \acute{u}$

$\check{Z} 7 \text{ } 69 p > \check{S} 6 \text{ } W \text{ } \check{n} u \text{ } 69 | | \frac{1}{2} ? q$; $A \beta \text{ } \check{Y} \text{ } i A \text{ } @ \bullet$
 $\dots ; | \check{A} \text{ } A \beta \text{ } \text{ } y , \hat{O} p \hat{O} i \check{z} \text{ } \ddot{z} \text{ } A \beta u \text{ } 69 p > | \frac{1}{2} \ddot{z}$
 $Y i \hat{O} B \text{ } L Q G \check{H} S \check{H} Q \check{W} C \check{L} G \check{H} Q \check{W} E \check{D} \check{W} \check{O} \check{L} \check{E} \check{X} \check{W} \check{H} / \text{ } \text{ } y ,$
 $/ u \text{ } V \text{ } \check{z} \wedge , \hat{U} \text{ } | > \check{Y} \text{ } x \gg \check{x} \text{ } \check{V} \text{ } \backslash \text{ } P P \check{H} \check{W} \text{ } E \text{ } S \check{O} \check{G} \text{ } V \text{ } \check{H} \check{O} \text{ } / v \text{ } \text{ } -$
 $y , \text{ } \check{z} j \text{ } \acute{a} \text{ } \check{z} \text{ } i y \text{ } f \check{n} \check{z} \text{ } \check{z} \times \check{z} , \acute{a} \text{ } u \text{ } \hat{O} i \text{ } \acute{U} \hat{O} \text{ } \check{z} \hat{O} p \text{ } \hat{A} \beta \text{ } ,$

Yi Wæ!9 SURE DWE LGH OXQ FWLR

$\check{S} 6 \text{ } V \check{z} \wedge | | \frac{1}{2} ? u \text{ } Y i \check{u} f \text{ } ! 9 p , \% \text{ } \partial C y ,$
 $V f 3 6 \text{ } \hat{O} B \hat{O} S , \gg \sim U \text{ } \frac{1}{4} \hat{O} w \text{ } p \text{ } h 9 x v > \{ d \text{ } C \text{ } \hat{O} v \text{ } Y i 69 \text{ } X(w) \text{ } x$
 $V f T \check{a} 6 \text{ } Y i 69 x , Y i \check{u} f ! \text{ } 9 \text{ } S U R E D W E L G H O X Q F W L R Q Y S 6 9 X$
 $| | \dots x^* \text{ } p \text{ } \check{x} \text{ } Y i \text{ } p(x) v \text{ } 9 e \text{ } \hat{O} \text{ } I S ; \check{Y} I^2 \text{ } \check{z} ,$



$f \beta \text{ } \acute{U} u \text{ } \hat{O} B \check{z} \gg / , \mu \check{z} \text{ } \acute{x} \check{u} , \acute{x} \text{ } \hat{O} B , W \check{x} \gg / , W \check{x}^* \text{ } \check{x} ,$
 $\hat{O} ; f \beta \gg / , W \check{x} h \hat{O} ! 9 q \check{C} , \acute{x} , x v \text{ } \check{n} \text{ } 9 \text{ } \check{x} , \gg / \check{Z}' \text{ } \emptyset$

$\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $W \sim \chi^2_{n-1}$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $b = \frac{1}{2} p \hat{\sigma}^2 \sim \frac{1}{2} p \frac{1}{n-1} W s^2$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$

S. [T i V f T a n

$V \sim \chi^2_{n-1}$
 $Y_i \sim N(\mu, \sigma^2)$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$

$\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$

$\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$

$$\Pr\left(\left|\frac{\bar{x} - \mu}{s/\sqrt{n}}\right| > z_{\alpha/2}\right) = 1 - D(\bar{x}, z_{\alpha/2}, s/\sqrt{n}) \approx \alpha$$

$\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$
 $\hat{\mu} \sim t_{n-1}(\bar{x}, s^2/n)$
 $\hat{\sigma}^2 \sim \frac{1}{n-1} W s^2$

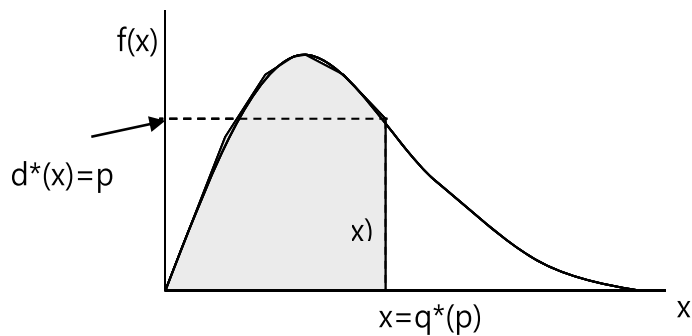
$B(p < 0.1)$ *n ü † Y i W æ q y d
 * n ü † W æ B(p < 0.1) p ð : % Y ò B v ´ . ò B á v ñ d
 ò B á v † H f G ñ ~ v ñ È ñ % f G ñ ~ † » á v æ ! È
 ‘ ¶ / d
 b , - ö v ! @ Ž B á v ñ È Ž f G ñ ~ ø
 » á v æ ! É u f G ñ ~ , 9 q ¶ s d
 ñ í ò B á Ž q † H ò B V † ^ Y i W æ 9 q y d
 ø [u v y d

* n ü † W æ B(p < 0.1) p ð : % Y i ò B v ´ . á v ñ È †
 - ö v ! @ Ž ò B á v ñ d
 ñ í ò B á Ž q † H ò B V † ^ Y i W æ 9 q y d
 ø [u v y

b, p Š - ö u Normal(10,5) p ð h d d

y À T ã 6 U V K U V K L E C T I K D W W I P C E R K Q P

6	
G [» 9	Y i ù f ! 9 Y i ... I [
S [S 9 »	W æ ! 9 ...) [
T S » 9	W æ ! 9 ...) S
U Q » 9	W æ ! 9 - n V † ^ Q Ž P - â ý



K W W S Z Z V W H M R G V P Q Q H W H P H L Q G H Q K M x y

¼ ê Ô E Q P V T Q W G O G P V

I R U 96 L Q ^ x `	p ½y ... f > 6 9 ... ‡ 6 X ´ ð ' x ' v @ h - ,
-----------------------	---

$L(\beta) = \sum_{i=1}^n \log f(x_i; \beta)$	$= \sum_{i=1}^n \log f(x_i; \beta)$
--	-------------------------------------

$Z(\beta) = \sum_{i=1}^n \log f(x_i; \beta)$	$= \sum_{i=1}^n \log f(x_i; \beta)$
--	-------------------------------------

6 € —

	SORWP[D \Q VXE' [LP F D \ØDE WSH' '	! 9
--	-------------------------------------	-----

'c • 6 X E

줄기 잎 그림 stem and leaf

+ 진단 내용

- 1) 분포의 개략적인 형태를 알 수 있다.
 - (1) 좌우 대칭인가? 아니면 skewed 되었는가?
 - (2) 봉우리(modal)는 하나인가? 아니면 여러 개인가?
- 2) 이상치의 존재 여부를 쉽게 파악할 수 있다.

+ 데이터

53 145
 43 621
 33 262
 45 208
 46 362
 55 424
 41 339
 55 736
 36 291
 45 58
 55 498
 50 643
 49 390
 47 332
 69 750
 51 368
 48 659
 62 234
 45 396
 37 300
 50 343
 50 536
 50 543
 58 217
 53 298
 57 1103
 53 406
 61 254
 47 862
 56 204

줄기	잎
0	5
1	4
2	00135699
3	03346699
4	029
5	34
6	245
7	35
8	6
9	
10	
11	1

(정렬)

줄기	잎
0	5
1	4
2	60931950
3	63936904
4	290
5	34
6	245
7	35
8	6
9	
10	
11	1

(정렬 없음)

+ 그리는 순서

- 자료를 크기 순으로 정리한다.
 자료의 수가 많을 때는 자료 정렬을 수작업 하기 어려움으로 이 단계는 무시해도 되지만 자료를 크기 순으로 정렬해 놓으면 plot 을 그리기 편리하다.
- 자료를 살펴 줄기와 잎을 결정한다.
 CEO 연봉 자료를 살펴보면 100 단위를 줄기로 하고 10 단위 이하를 잎으로 하여 plot 을 그리면 될 것이라는 것을 알 수 있다. 줄기 수는 히스토그램의 계급 구간 수에 해당되므로 8~12 정도가 적절하다. 적정 개수가 아닌 경우 줄기 수 조정에 대해서는 다음에 다루기로 한다.
- 한 열에 줄기(stem)를 먼저 그린다.



위에서 100 단위 이상을 줄기로 하기로 결정하였고 자료의 최소값이 58, 최대값이 1103 이므로 0 부터 11 까지 줄기를 한 열에 크기 순으로 적는다.

```

줄기|
-----
0
1
2
3
4
5
6
7
8
9
10
11
    
```

- 줄기(stem) 옆에 잎을 그린다.

잎을 그리는 방법은 간단하다. 줄기 바로 뒤의 숫자를 줄기 옆에 차례로 적으면 된다. CEO 연봉 자료는 잎이 두 자리이지만 앞에 것 하나만 적으면 된다. 굳이 반올림하는 수고를 할 필요는 없다. 줄기-잎 그림의 목적은 자료의 분포 형태와 이상치를 아는 것이 주된 목적이기 때문이다.



줄기-잎 그림

```
ds=read.csv("ceo.csv",header=T)
stem(ds$CEO.Salary)
```

- + 엑셀에서逗마가 있는 파일 형식으로 저장한 후 읽어 들인다.
- + ds\$변수명; 오브젝트 ds 내의 변수명 변수를 이용 지정

The decimal point is 3 digit(s) to the right of the |

```

0 | 178889
1 | 000012334446
2 | 0
3 | 3
    
```



줄기-잎 그리기

<http://lib.stat.cmu.edu/DASL/Stories/SingerHeights.html>

합창단원의 키에 대한 (단위: inch) 데이터이다.

- 1) 각 파트의 줄기-잎 그림을 그리고 해석하시오.
- 2) 키 데이터 전체에 대한 줄기-잎을 그리고 해석하시오.



+ Stem-leaf plot 해석하기

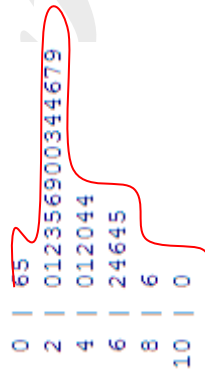
- > 자료의 분포 형태
stem-leaf plot 을 통하여 자료의 분포 형태를 알 수 있으므로 분포의 형태를 알 수 있다.
이는 히스토그램과 같은 역할이다.
- > 봉우리 (최빈값) 위치 및 개수 => 봉우리의 개수가 집단의 개수이다
- > 좌우 대칭 여부
- > 자료의 범위 및 분산
- > 이상치 존재 여부 및 위치

(히스토그램과 비교)

줄기-잎 그림을 90 도 회전하면 히스토그램 (이를 bar chart 라고도 함)이 된다. 히스토그램은 자료의 값의 정보가 상실되지만(실제 값은 알 수 없고 빈도만 바의 크기로 나타난다) stem-leaf plot 은 자료 값이 나타난다. 그러므로 히스토그램에 비해 더 많은 정보를 얻을 수 있다.

(1) 확률분포함수 추정

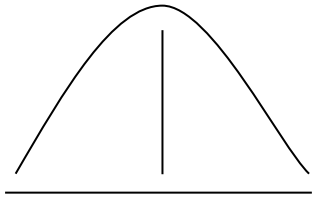
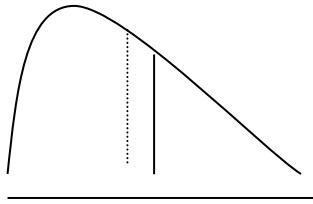
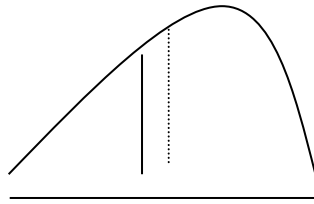
위의 예제처럼 stem-leaf plot 의 정점을 연결하면 확률분포함수를 얻게 된다. 아래 그림은 모집단 CEO 연봉의 확률밀도함수의 추정 형태이고 ($f(x)$) 면적은 1 이다.



(2) 대칭, 치우침 여부

symmetric (bell-shaped) 좌우 대칭, 종 모양	skewed to the right positively skewed 우로 치우침	skewed to the left negatively skewed 좌로 치우침
--	--	---



		
평균과 중앙값 일치	평균 > 중앙값	평균 < 중앙값
좌우대칭으로 만들려면... 자료 변환을 하면 된다.	$X^* = \sqrt{X} \rightarrow$ mild pos. $X^* = \log(X) \rightarrow$ pos. $X^* = -1/\sqrt{X} \rightarrow$ severe pos. $X^* = -1/X \rightarrow$ more severe	$X^* = X^2 \rightarrow$ mild neg. $X^* = X^3 \rightarrow$ extreme neg.

(정규성 검정) Anderson-Darling test for normality

```
library(nortest)
ad.test(ds$CEO.Salary)
```

```
> ad.test(ds$CEO.Salary)
```

Anderson-Darling normality test

```
data: ds$CEO.Salary
A = 1.1633, p-value = 0.003674
```

(연봉 모평균에 대한 95% 신뢰구간 구하기) conf.level=0.95

```
> t.test(ds$CEO.Salary)
```

One Sample t-test

```
data: ds$CEO.Salary
t = 8.3815, df = 19, p-value = 8.333e-08
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 896.6596 1493.5404
sample estimates:
mean of x
 1195.1
```

```
output=t.test(ds$CEO.Salary)
names(output)
output$conf.int
```

```
> output$conf.int
[1] 896.6596 1493.5404
```



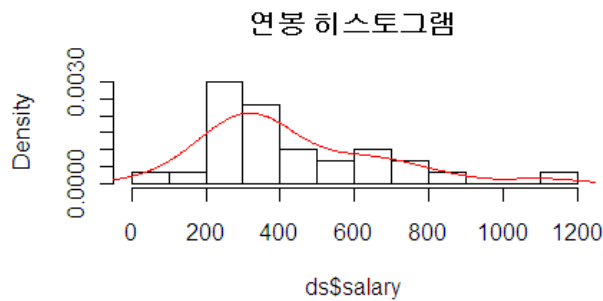


히스토그램 그리기

우로 치우침이 있으므로 제곱근 변환, 혹은 로그 변환 데이터 중 어느 변환이 더 좌우 대칭분포에 가까운지 알아보시오.

(3) 봉우리 위치 및 개수

히스토그램의 봉우리는 분포의 **최빈값**에 해당되는 부분으로 일반적으로 **최빈값**은 하나일 가능성이 가장 높다. 구간 설정에 따라 바로 옆의 구간이 동시에 **최빈값**이 되는 경향이 가끔 나타나기도 한다. 다음의 경우는 bi-modal 분포 함수라 하지는 않는다. 왜냐하면 구간을 조정하면 봉우리가 하나로 될 수 있기 때문이다. CEO 연봉은 단봉 형태를 갖는다.



단봉 uni-modal	다봉 bi-modal / multi-modal	
<p>봉우리가 2 개 이상인 의미는 모집단이 하나 이상일 가능성이 많다. 예를 들어 한남 대학생들 100 명의 몸무게를 조사하여 히스토그램을 그리면 bi-modal 형태가 될 가능성이 높다. 왜냐하면 여자와 남자 몸무게의 차이가 나기 때문에 그런 현상이 발생한다. 즉 측정 변수의 특성에 따라 모집단이 나누어진다. 용돈을 조사하여 히스토그램을 그려보면 아마 봉우리가 3-4 개일 가능성이 있다. 왜? 학년별 차이로 인하여... 이처럼 어떤 변수를 측정하느냐에 따라 같은 모집단이라도 봉우리의 개수가 다를 수 있다. 봉우리가 2 개 이상인 경우는 집단을 분리하여 추정 및 검정을 시행하는 것이 바람직하다. 그러나 집단에 대한 정보가 없다면 데이터를 분리하여 분석하는 것이 쉽지 않다.</p>		



(왜 좌우 대칭이어야 하나?)

- 1) 회귀 분석, 분산 분석 등 대부분의 통계 분석에서 종속변수는 정규분포를 따르고 있다는 가정을 한다. 만약 이것이 무너지면 t-검정, F-검정을 사용할 수 없다. → 3 학년 수업에서 배우기를...
- 2) 대표본 표본 크기 n 의 크기? : 자료 분석의 목적은 그래프 정리(bar chart, pie chart)나 숫자적 정리(평균, 표준편차)에서 끝나는 것이 아니라 이 정보를 가지고 모수(예:모집단의 평균)를 추정하거나 그에 대한 가설을 검정하게 된다. CEO 30 명의 연봉 자료를 이용하여 전체 CEO 의 연봉에 대해 알고 싶은 것이다.

통계 소프트웨어에서 출력되는 p-값은 two-sided(양측 검정) 가설 검정 시 값을 출력한다. 그러므로 위의 경우 대립 가설을 $H_a: \mu \neq 350$ (양측 검정) 설정하면 p-값이 0.0821 로 0.05 보다 크므로 귀무가설을 기각할 수 없으나 대립 가설을 $H_a: \mu > 350$ (단측 검정) 설정하면 p-값이 0.04105 이므로 0.05 보다 적어 귀무가설을 기각하고 연봉은 높아졌다고 결론 지을 수 있다. 그러므로 양측 검정 결과 귀무가설이 기각되면 같은 유의수준에서 단측 검정 결과도 귀무가설을 기각한다.

(4) 범위와 흩어진 정도

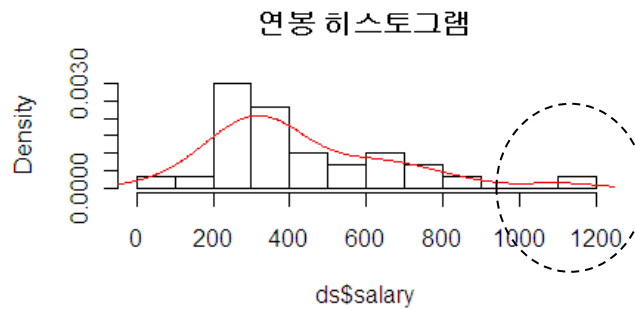
분포의 형태를 알 수 있으므로 자료의 범위(range=최대값-최소값)와 흩어진(spread) 정도를 알 수 있다.



(5) 이상치(outlier) 발견

다른 관측치에 비해 매우 크거나 적은 관측치를 이상치(outlier)라 한다. 이런 이상치는 히스토그램에서 쉽게 발견될 수 있다. 히스토그램이나 stem-leaf plot 의 경우 다른 관측치와 멀리 떨어져 있으면 이를 이상치라 한다. CEO 연봉 자료에서 이상치는 연봉이 1103(백만)인 사람이다. 물론 이 값이 이상치인지는 검정 통계량을 이용하여 (Box-plot 이나 검정 방법을 이용하여 검정해야 하지만 우선 쉽게 찾을 수 있다는 장점이 있다. CEO 연봉의 경우 다른 CEO 에 비해 연봉을 이상적으로 높게 받는 CEO (이를 이상치라 함) 가 있음을 알 수 있다.





이상치가 발견되면 그 해결책은

이상치인 관측치의 원 자료를 확인, 입력 오류인지 살펴본다. 오류가 있으면 정정한다.

이상치의 대상인 개체를 조사해 문제가 있는 개체이면 자료에서 제외한다.

예를 들면 1103(백만) 연봉을 받는 사람을 조사하였더니 외국인 전문 사장이었다. 국내 CEO 연봉으로 간주하기 어렵다면 제외

여전이 유효한 데이터이면 자료 변환을 통하여 이 문제를 해결하게 된다. 변수변환 (자료 변환)을 통하여 이상치 문제가 해결되면 이는 치우침의 한 부분이 된다.

+ 줄기 수 조정

일반적으로 자료의 분포 형태를 잘 파악하기 위해서는 줄기의 수가 8-10 개 정도 되어야 한다고 한다. 연봉 데이터 예제에서 본 것처럼 줄기 수는 변수 측정치의 범위에 의해 결정된다. 그러므로 줄기의 수를 조정하여 적절한 줄기-잎 그림을 그려야 한다.

> 줄기 수가 너무 많으면 (squeezed stems)

줄기를 일정한 수만큼 합치는 방법을 생각하면 된다. 만약 줄기가 1-20 까지 있다면 1-2, 3-4, 5-6, ..., 19-20 을 각각 줄기로 하면 줄기 수가 20 개에서 10 개로 줄어든다. 이처럼 줄기 수에 따라 2 배, 3 배, 4 배씩 줄이면 된다.

> 줄기 수가 너무 적으면 (stretched stems)

줄기를 2 등분(double stem) 혹은 5 등분(five-line stem)하여 사용한다.

(예) double stem: 1* (1.0~1.4), 1. (1.5~1.9)

(예) five-line stem: 1* (.0, .1), 1_t (.2, .3), 1_f (.4, .5), 1_s (.6, .7), 1. (.8, .9)



자료	줄기	잎
11	1*	1
12	1t	2
15	1f	5
17	1s	7
18	1.	88
19	2*	0011
20	2t	
20	2f	4
21	2s	8
21	2.	
24		
28		

적정 줄기 수에 관한 공식

- > Sturges formula $L = [1 + \log_2 n]$ (예) $n=30 \rightarrow L=5$
- > Velleman formula $L = [2\sqrt{n}]$ (예) $n=30 \rightarrow L=10$
- > Dixon-Kronmal formula $L \leq [10\log_{10} n]$ (예) $n=30 \rightarrow L=14$

그러나 위의 공식에 의해 줄기 수(L)를 결정하면 자료 값에 따라 줄기를 결정하기 어렵고 분포 형태를 제대로 알기 어려운 문제가 있어 이 공식들은 사용되지는 않는다. [x]의 의미는 x보다 크지 않는 최대 정수 값을 의미한다. $[2.9]=2$ / $[3.1]=3$



R 활용

```
hist(ds$CEO.Salary, nclass=10, freq=F, main="Histogram")
lines(density(ds$CEO.Salary), col="blue")
```

- nclass 옵션은 구간의 개수를 결정한다.
- freq 옵션은 빈도 대신 상대빈도 (확률)을 y-축으로 사용하라는 옵션
- 함수 lines()는 확률밀도함수를 그리라는 옵션



히스토그램 그리기

<http://lib.stat.cmu.edu/DASL/Stories/SingerHeights.html>

합창단원의 키에 대한 (단위: inch) 데이터이다.

- 1) 각 파트별 히스토그램을 그리고 해석하시오. (확률밀도함수도 그리시오)
- 2) 키 전체에 대한 히스토그램을 그리고 해석하시오.



상자 수염 그림 box whisker plot

Stem and leaf(줄기-잎) plot 은 자료의 분포의 형태(좌우 대칭, 단봉) 파악과 이상치를 발견할 수 있는 도구이다. 그러나 S-L plot 만 가지고는 정확한 중앙 위치, 자료의 사분위 값의 위치 (25%, 50%, 75% 부분의 값들), 이상치라 시각적으로 판단되는 관측치가 정말 이상치인지 알 수 진단해야 한다.

Box-whisker (상자 수염) 그림은 plot 중앙값 (때로 평균까지) [상자 안의 실선], 사분위 위치, [상자 끝단], 자료의 최대값, 최소값(whisker), 이상치 존재 여부(bullet)를 그려 놓은 상자 형태의 그림이다.

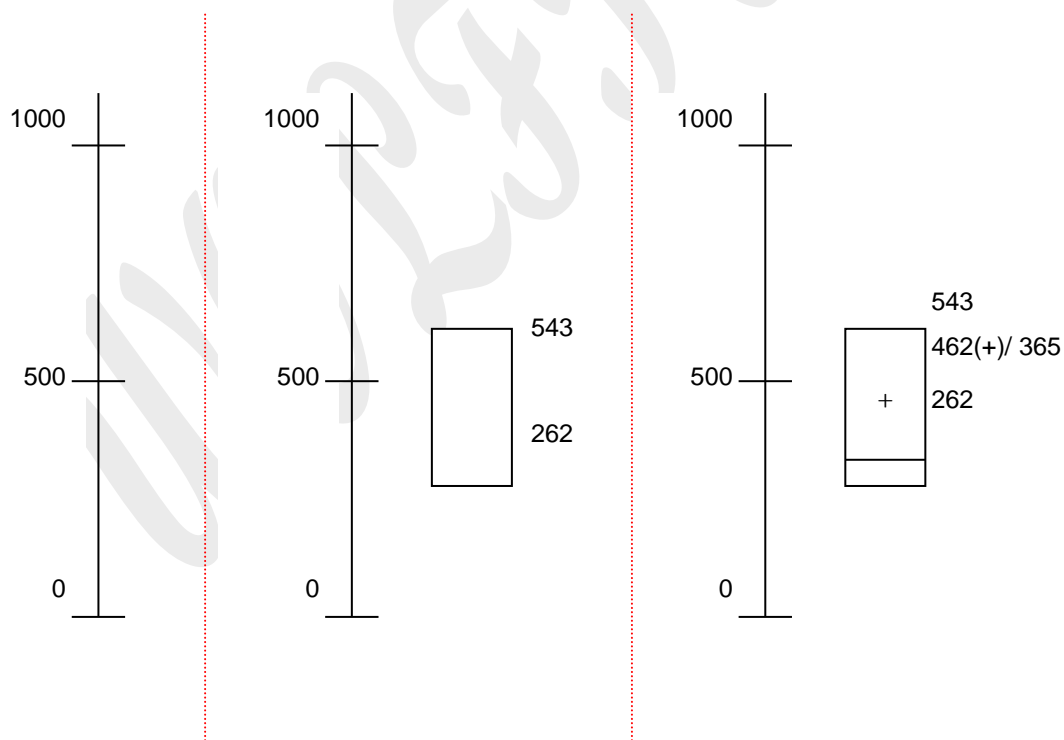
+ Box and whisker plot 그리기

상자로부터 나온 선이 수염처럼 생겨 Box and whisker plot 이라고 하는 Box plot 을 그리는 순서는 다음과 같다. [CEO 자료 이용]

[순서 1] 자료의 최소값, 최대값을 이용하여 y 축 선을 그린다.

[순서 2] Q1, Q3 를 이용하여 상자를 그린다. 상자의 넓이는 아무 의미가 없다.

[순서 3] 상자 가운데 중앙값을 그리고 평균은 기호로 (+) 표시한다.



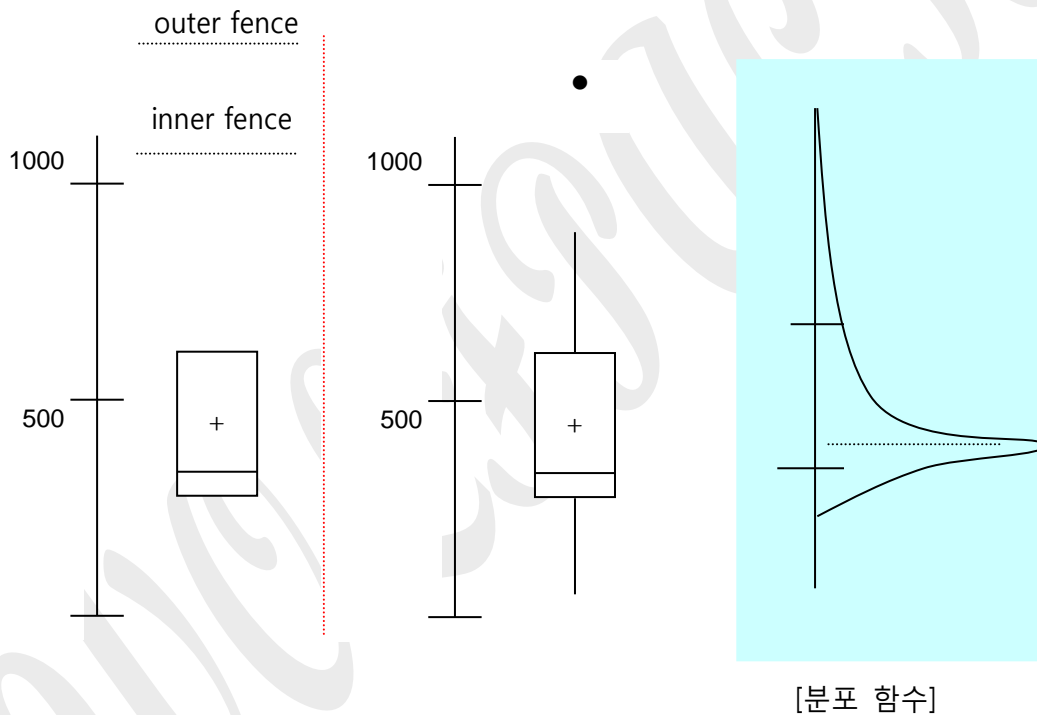
[순서 4] IQR 을 이용하여 가상 선 (imaginary line) Inner fence, Outer fence 를 그린다. 가상 선은 실제 상자 그림에 표시되지 않는다. 이상치 존재 여부를 표현하기 위한 임시 선이다.

$$IQR = (543 - 262) = 281$$

$$\text{Inner fence } (Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR) = (-159.5, 964.5)$$

$$\text{Outer fence } (Q_1 - 3 \times IQR, Q_3 + 3 \times IQR) = (-581, 1386)$$

[순서 5] 수염과 이상치를 표시한다. 관측치 중 Inner fence 를 넘지 않는 최대, 최소값까지 수염을 그린다. Fence 를 넘는 관측치를 이상치라 (outlier) 한다. outer fence 까지 넘는 관측치는 severe (극심한) 이상치, inner fence 만 넘으면 mild 이상치라 한다. [CEO 에서 1103 은 mild 이상치]



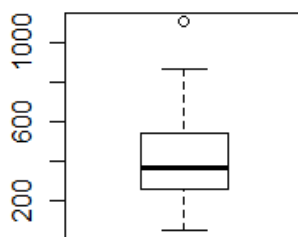


줄기-잎 그림

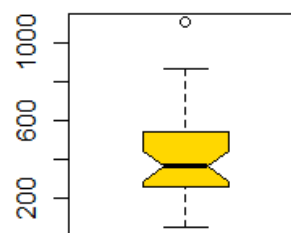
```
ds=read.csv('ceo.csv')
boxplot(ds$salary,main="연봉 상자-수염그림")

boxplot(ds$salary,notch=TRUE,
col=(c("gold","darkgreen")),main="연봉 상자-수염그림")
```

연봉 상자-수염그림



연봉 상자-수염그림



+ Box Plot 해석하기

분포의 형태

박스, 박스 안의 선(중앙값), 수염의 길이를 이용하여 분포의 형태를 짐작할 수 있다. 박스에 50% 자료가 있고 박스 위 부분에 25%, 박스 아래 부분에 25%가 있다. 박스 내에서도 중앙 선 위 부분에 25%, 아래 부분에 25%가 있으므로 분포의 형태를 알 수 있다. CEO 연봉은 우측으로 치우친 형태를 띠고 있음을 알 수 있다. (위쪽, 즉 확률밀도함수의 오른쪽 상자와 수염 부분이 왼쪽에 비해 살짝 길다) 그러므로 평균이 중앙값보다 크므로 역시 치우친 형태임을 알 수 있다. 확률분포함수는 빨간 선이다.

단점은 봉우리의 개수를 알지 못하는 단점이 있으므로 보완적으로 줄기-잎 그림 그린다.

중앙값, 산포 정도, 군집

자료 관측치의 중앙 위치, 그리고 관측치들의 어디에 모여 있는지 (군집), 자료 값들의 흩어진 정도를 파악할 수 있다. 중앙 값이 350 부근(실제로는 365)임을 알 수 있다. 값의 범위(range), 사분위 값을 대략적으로 알 수 있다.



이상치 존재 판단

자료 관측치 중 다른 값들에 비해 지나치게 크거나 작은 관측치, 이상치의 존재 여부를 파악할 수 있다. 연봉이 1100 이상을 받는 CEO 는 이상치임을 알 수 있다.

줄기-잎 그림을 통해서는 이상치 존재 여부를 짐작할 수 있지만, 상자 수염 그림을 통하여 판단 가능하다.

stem-leaf plot 과는 관측치 값들에 대한 정보를 얻을 수 없다는 단점이 있으나 표본 분포의 형태도 파악할 수 있고 중앙값과 이상치를 표시하여 표본 자료의 정보 파악이 용이하다. 특히 box-whisker plot 은 집단간 자료의 분포 차이 비교, 모평균 차이 검정 시 매우 유용한 plot 이다.



상자 수염 그리기

<http://lib.stat.cmu.edu/DASL/Stories/SingerHeights.html>

합창단원의 키에 대한 (단위: inch) 데이터이다.

- 1) 키 데이터 전체에 대한 줄기-잎을 그리고 해석하시오.
- 2) 각 파트의 상자 수염 그림을 한 그래프에 그리고 해석하시오.



상자 수염 그리기 FASTFOOD.xls

미국 fast-food 레스토랑 5 개의 Drive-through 서비스 소요시간을 측정한 것이다. 레스토랑 상자-수염 그림을 한 화면에 그리고 결과를 해석하자.

