

平均

- 標本平均 (\bar{x} で表す) は次のように計算できる：

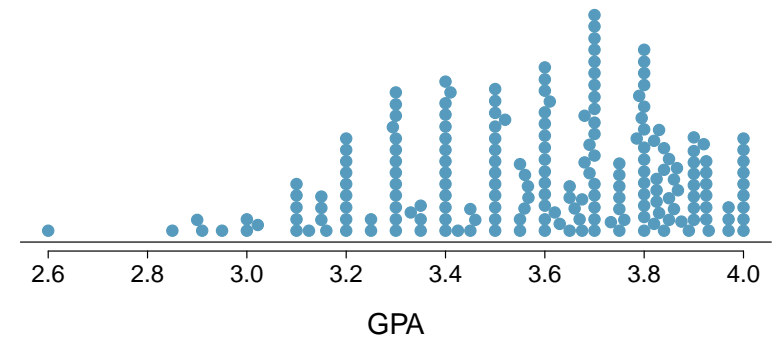
$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n},$$

ここで x_1, x_2, \dots, x_n は n 個の観測値を表す。

- 母平均も同様に計算されるが μ で表される。母集団データが得られることはまれなため、 μ を計算することはしばしば不可能である。
- 標本平均は**標本統計量**であり、母平均の**点推定値**として機能する。この推定は完璧ではないかもしれないが、標本が良好(母集団を代表している)であれば、通常は十分な推定値となる。

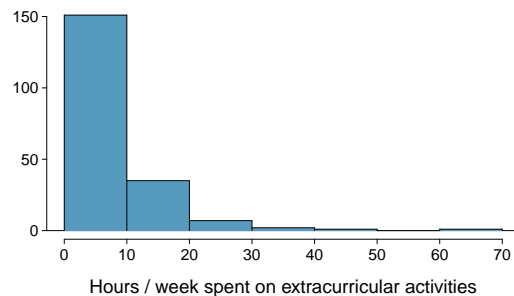
積み上げドットプロット

棒が高い部分はより多くの観測値が集中しており、分布の中心と形状をやや判断しやすくなる。



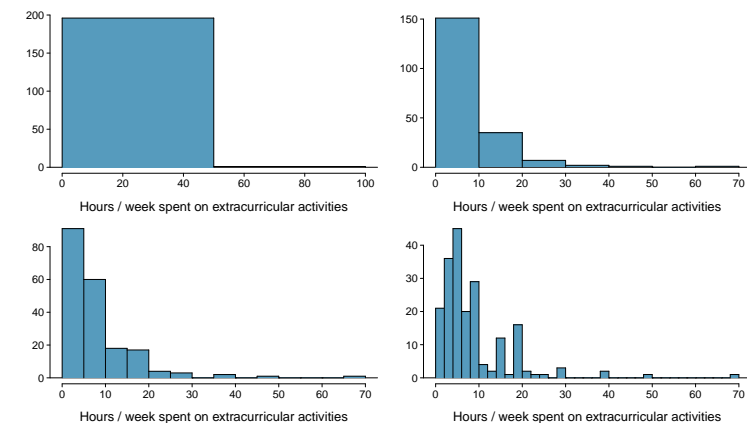
ヒストグラム—課外活動時間

- ヒストグラムは**データの密度**を表す。棒が高い部分はデータが相対的に多い。
- ヒストグラムはデータ分布の**形状**を説明するのに特に便利である。
- 選択する**階級幅**によってヒストグラムが示す内容が変わる。



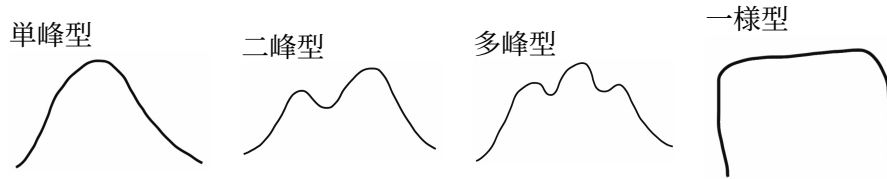
階級幅

これらのヒストグラムのうち有用なものはどれか？ データを明かしすぎているのはどれか？ 隠しすぎているのはどれか？

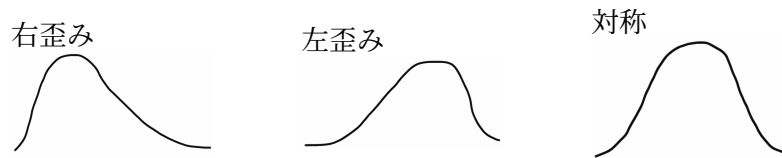


よく見られる分布の形状

- 最頻値の個数



- 歪み



練習問題

次の変数のうち一様分布すると予想されるものはどれか？

- (a) 成人女性の体重
- (b) ノースカロライナ州の無作為標本の給与
- (c) 住宅価格
- (d) クラスメートの誕生日（日付）

応用活動：分布の形状

次の変数の予想される分布をスケッチせよ：

- ピアスの数
- 試験の得点
- IQ スコア

任意の変数の予想分布を判断する方法を簡潔に（1～2文で）説明せよ。

あなたは「典型的な」人か？



<http://www.youtube.com/watch?v=4B2xOvKFFz4>

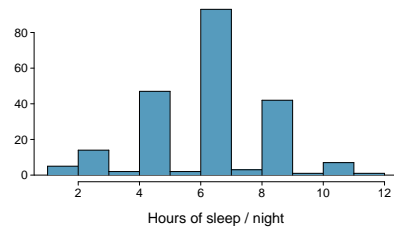
分布の真の特徴を伝えるうえで、中心だけがどれほど有用か？

分散

分散は平均からの偏差の2乗のおおよその平均である。

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

- 標本平均は $\bar{x} = 6.71$ 、標本サイズは $n = 217$ である。
- 学生が1晩に睡眠をとる時間の分散は次のように計算できる：



$$s^2 = \frac{(5 - 6.71)^2 + (9 - 6.71)^2 + \dots + (7 - 6.71)^2}{217 - 1} = 4.11 \text{ 時間}^2$$

分散（続き）

なぜ分散の計算に偏差の2乗を使うのか？

標準偏差

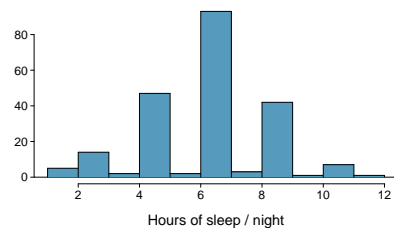
標準偏差は分散の平方根であり、データと同じ単位を持つ。

$$s = \sqrt{s^2}$$

- 学生が1晩に睡眠をとる時間の標準偏差は次のように計算できる：

$$s = \sqrt{4.11} = 2.03 \text{ 時間}$$

- すべてのデータが平均の3標準偏差以内にあることが分かる。



中央値

- 中央値は、データを昇順に並べたときに真ん中を2分割する値である。

0, 1, 2, 3, 4

- 観測値の数が偶数の場合、中央値は真ん中2つの値の平均である。

$$0, 1, \underline{2}, \underline{3}, 4, 5 \rightarrow \frac{2 + 3}{2} = 2.5$$

- 中央値はデータの間接点であるため、値の50%がその下にある。したがって、第50パーセンタイルとも呼ばれる。

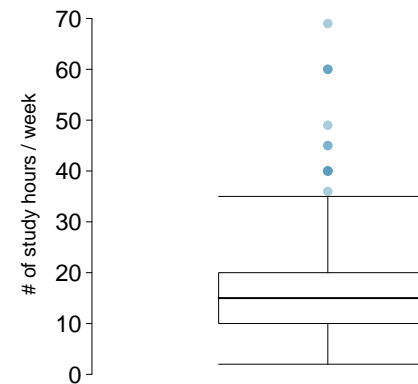
Q1、Q3、IQR

- 第25パーセンタイルは第1四分位数 (Q1) と呼ばれる。
- 第50パーセンタイルは中央値とも呼ばれる。
- 第75パーセンタイルは第3四分位数 (Q3) と呼ばれる。
- Q1 と Q3 の間にはデータの中央50%がある。このデータが示す範囲を四分位範囲 (IQR) と呼ぶ。

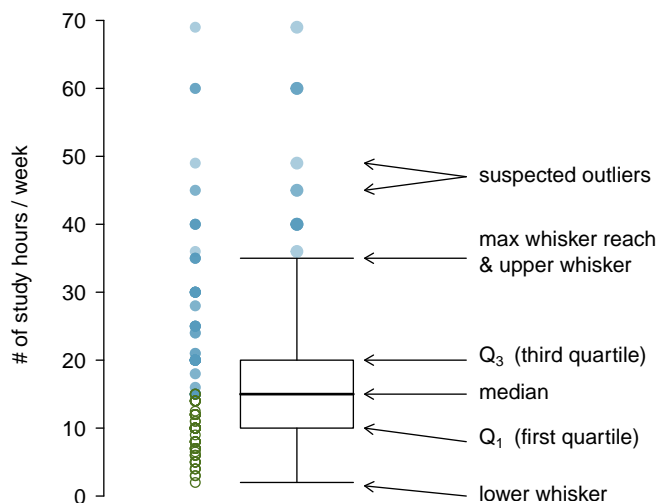
$$IQR = Q3 - Q1$$

箱ひげ図

箱ひげ図の箱はデータの中央50%を表し、箱内の太線は中央値である。



箱ひげ図の構造



ひげと外れ値

- 箱ひげ図のひげは四分位数から最大 $1.5 \times IQR$ まで伸びる。

$$\text{上ひげの最大リーチ} = Q3 + 1.5 \times IQR$$

$$\text{下ひげの最大リーチ} = Q1 - 1.5 \times IQR$$

$$IQR : 20 - 10 = 10$$

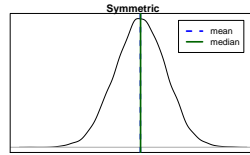
$$\text{上ひげの最大リーチ} = 20 + 1.5 \times 10 = 35$$

$$\text{下ひげの最大リーチ} = 10 - 1.5 \times 10 = -5$$

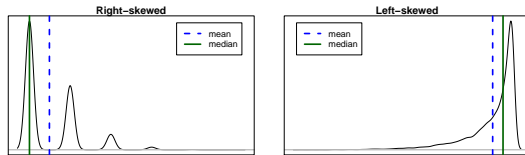
- 潜在的な外れ値とは、ひげの最大リーチを超えた観測値である。残りのデータに対して極端に見える観測値である。

平均と中央値

- 分布が対称な場合、中心はしばしば平均として定義される：平均 \approx 中央値

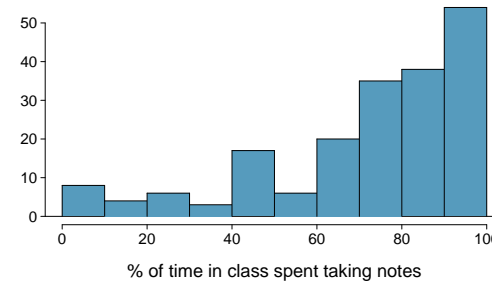


- 分布が歪んでいるか極端な外れ値がある場合、中心はしばしば中央値として定義される
 - 右歪み：平均 $>$ 中央値
 - 左歪み：平均 $<$ 中央値



練習問題

授業中にノートを取ることに費やした時間の割合（Facebook や Twitter などと比較した）の分布について最も正しいと思われるのはどれか？



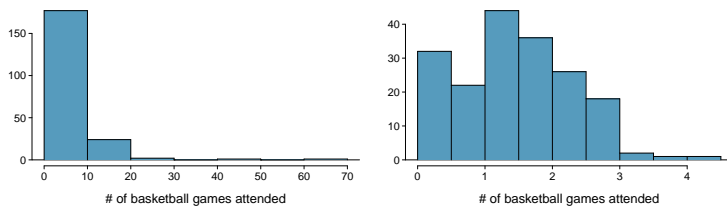
中央値：80%
平均：76%

- (a) 平均 $>$ 中央値
- (b) 平均 $<$ 中央値
- (c) 平均 \approx 中央値
- (d) 判断不可能

極端に歪んだデータ

データが極端に歪んでいる場合、変換すると分析しやすくなることもある。一般的な変換は対数変換である。

左のヒストグラムは学生が観戦したバスケットボールの試合数の分布を示す。右のヒストグラムは試合数の対数の分布を示す。



変換の利点と欠点

- 歪んだデータは適切な変換によって外れ値が目立たなくなるため、変換後の方が分析しやすい。

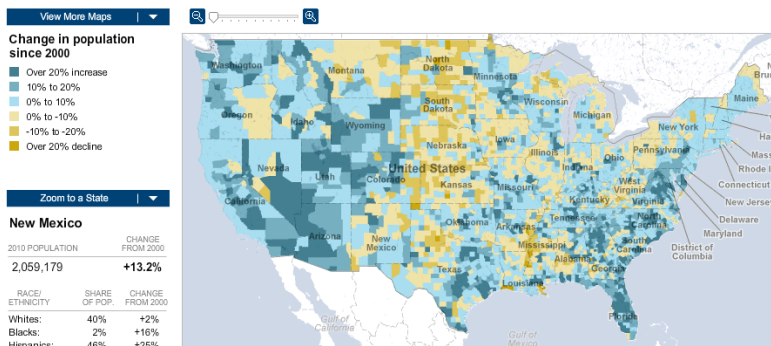
試合数	70	50	25	...
log (試合数)	4.25	3.91	3.22	...

- ただし、測定変数の対数単位での分析結果は解釈が難しいことがある。

他にどのような変数が極端に歪むと予想されるか？

強度マップ

2000年から2010年にかけての人口変化において、どのようなパターンが見られるか？



<http://projects.nytimes.com/census/2010/map>

分割表

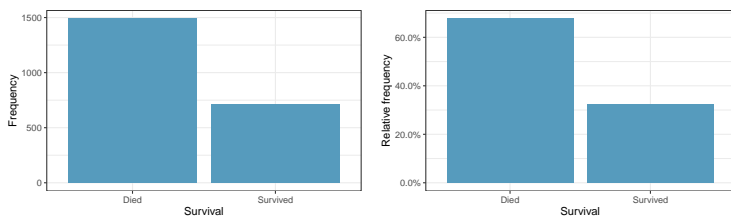
2つのカテゴリ変数のデータをまとめた表を分割表と呼ぶ。

下の分割表はタイタニック号の乗客の生存と年齢の分布を示す。

		生存		合計
		死亡	生存	
年齢	大人	1438	654	2092
	子供	52	57	109
	合計	1490	711	2201

棒グラフ

棒グラフは1つのカテゴリ変数を表示する一般的な方法である。頻度の代わりに比率を示す棒グラフを相対頻度棒グラフと呼ぶ。



棒グラフとヒストグラムの違いは何か？

適切な比率の選択

タイタニック号の乗客において、年齢と生存の間に関係があるように見えるか？

		生存		合計
		死亡	生存	
年齢	大人	1438	654	2092
	子供	52	57	109
	合計	1490	711	2201

この問いに答えるために行比率を調べる：

- 生存した大人の割合：654 / 2092 ≈ 0.31
- 生存した子供の割合：57 / 109 ≈ 0.52

