

# CHƯƠNG 2

## MÔ TẢ CÁC TẬP DỮ LIỆU

*Về chương này:*

Đôi lúc dữ liệu chúng ta đã thu thập thể hiện một mẫu được chọn từ một tổng thể. Những lúc khác (chẳng hạn như một cuộc tổng điều tra dân số toàn quốc), dữ liệu có thể thể hiện toàn bộ tổng thể. Trong cả hai trường hợp, chúng ta đều cần phải có khả năng mô tả tập dữ liệu đó. Mục tiêu của chương này là trình bày hai loại phương pháp mô tả các tập dữ liệu: (1) các phương pháp mô tả bằng đồ thị và (2) các phương pháp mô tả bằng số. Phương pháp mô tả bằng đồ thị mô tả dữ liệu bằng cách sử dụng biểu đồ và đồ thị. Phương pháp mô tả bằng số sử dụng các con số để giúp chúng ta xây dựng một hình ảnh trong trí óc về dữ liệu.



## NGHIÊN CỨU TÌNH HUỐNG

### VẬY LÀ ANH/CHỊ MUỐN TRỞ THÀNH NHÀ TRIỆU PHÚ?

Vào thập niên 1980, các luật thuế mới đã dẫn đến việc tạo ra Tài khoản Hưu trí Cá nhân (Individual Retirement Accounts-IRA), đó là những tài khoản tiết kiệm miễn thuế đối với thu nhập hưu trí. Dựa theo nhiều mẫu quảng cáo trên báo chí vào lúc đó, nếu mà chúng ta đầu tư 2.000USD mỗi năm vào một Tài khoản Hưu trí Cá nhân (IRA), thì sau 40 năm tham gia, tiền dự trữ của chúng ta sẽ tăng lên đến trên một triệu đô la. Tất nhiên là kể từ đó, các luật thuế đã được thay đổi vài lần, và với việc xem xét lại thuế gần đây nhất, IRA miễn thuế sẽ không còn có sẵn cho hầu hết chúng ta. Dù vậy, cái nguyên tắc được thể hiện trong những mẫu quảng cáo đó vẫn còn có giá trị. Cách thức tốt nhất để tích lũy một số tiền lớn là tham gia vào một chương trình tiết kiệm và đầu tư có hệ thống và tính lãi kép những số tiền đầu tư qua nhiều năm.

Nếu anh/chị đang để dành tiền cho thời kỳ hưu trí hay nhằm mục đích nào khác, thì số tiền anh/chị tích lũy được sẽ phụ thuộc vào số tiền anh/chị đầu tư mỗi năm, nơi anh/chị đầu tư (tài khoản tiết kiệm tại ngân hàng, quỹ đầu tư thị trường vốn ngắn hạn, hay một trong những quỹ cổ phiếu thường khác nhau), và ai quản lý tài khoản của anh/chị. Về cơ bản, mức tăng trưởng của tài khoản của anh/chị và giá trị cuối cùng của nó sẽ phụ thuộc vào suất sinh lợi hàng năm mà nhà quản lý tài khoản của anh/chị có thể thu nhận được cho anh/chị.

Mặc dù suất sinh lợi từ tiền đầu tư của anh/chị sẽ thay đổi từ ngày này sang ngày khác, nhưng Bảng 2.1 cho anh/chị biết số tiền mình có thể kỳ vọng sẽ tích lũy được sau 40 năm. Những số tiền được trình bày trong bảng này dựa vào giả định rằng anh/chị đầu tư 2.000USD vào đầu mỗi năm trong thời kỳ 40 năm và tiền này được tính lãi kép hàng tháng với lãi suất hàng năm cố định là  $I$ .

Số tiền rút ra của một tài khoản sau khi thực hiện đầu tư hàng năm 2.000USD với suất sinh lợi hàng năm cố định $I$ (%) trong 40 năm	Số tiền trong Tài khoản	
	Lãi suất $I$ (%)	Sau 40 Năm (\$)
	4	197.652
	6	328.095
	8	559.562
	10	973.704
	12	1.718.285

Suất sinh lợi hàng tháng thay đổi không chỉ giữa các công cụ đầu tư (ngân hàng, quỹ cổ phiếu thường, quỹ đầu tư thị trường vốn ngắn hạn, v.v.) mà còn trong phạm vi bản thân một công cụ đầu tư. Thí dụ, hãy xét một quỹ đầu tư thị trường vốn ngắn hạn, mà nó đầu tư vào kỳ phiếu thương mại ngắn hạn. Bởi vì những kỳ phiếu này sẽ được thương lượng vào những thời điểm khác nhau và sẽ có khoảng thời gian đến khi đáo hạn khác nhau, nên suất sinh lợi đối với toàn bộ quỹ này tại bất kỳ thời điểm nào đều sẽ phụ thuộc vào kỹ năng của nhà quản lý quỹ trong việc cho vay. Nếu lãi suất tăng trong tương lai, thì nắm giữ các kỳ phiếu có thời gian đến khi đáo hạn trung bình ngắn sẽ có lợi. Nếu lãi suất giảm, thì nắm giữ các kỳ phiếu có thời gian đến khi đáo hạn trung bình dài sẽ có lợi.





Những đặc điểm của quỹ đầu tư thị trường vốn ngắn hạn như là một công cụ đầu tư được cho thấy trong dữ liệu của Bảng 2.2. Bảng 2.2 trình bày qui mô tài sản (tính bằng triệu đô la), thời gian đáo hạn trung bình (tính bằng ngày) của kỳ phiếu, và lợi suất 7–ngày trung bình (%) trong thời kỳ kết thúc vào ngày 13/7/1994, đối với 604 quỹ đầu tư thị trường vốn ngắn hạn lớn và có sẵn cho các nhà đầu tư. Xem xét Bảng 2.2 thì chúng ta sẽ thấy rõ vấn đề khó khăn về thống kê. Mặc dù có thể có được cảm nhận tổng quát về qui mô tài sản, thời gian đáo hạn trung bình, và suất sinh lợi trung bình qua việc xem xét dữ liệu trong bảng này, nhưng khó mà có được một hình ảnh rõ ràng về những đặc điểm của các tập dữ liệu này bằng cách chỉ xem xét kỹ bảng này. Vấn đề này thúc đẩy chúng ta nghiên cứu đề tài của Chương 2. Trong chương này, chúng ta xem xét những phương pháp mô tả các tập dữ liệu. Sau đó, trong Mục 2.14 (trong nguyên bản tiếng Anh), chúng ta áp dụng những kỹ thuật này vào dữ liệu về quỹ đầu tư thị trường vốn ngắn hạn nói trên và xem thông tin có tính mô tả này phù hợp như thế nào với triển vọng trở thành nhà triệu phú của chúng ta.

## 2.1 Biến (Variables) và Dữ liệu (Data)

Mục tiêu chủ yếu của chúng ta trong Chương 2 sẽ là trình bày một số kỹ thuật căn bản trong **thống kê mô tả (descriptive statistics)**—ngành thống kê liên quan đến việc mô tả những tập hợp các giá trị đo lường, cả **mẫu (sample)** và **tổng thể (population)**. Sau khi chúng ta đã thu thập một tập hợp các giá trị đo lường (measurements), làm sao chúng ta có thể trình bày tập hợp này dưới một hình thức rõ ràng, có thể hiểu được và dễ đọc? Trước tiên, chúng ta phải có thể định nghĩa giá trị đo lường hay dữ liệu là gì và phân loại các loại dữ liệu chúng ta có khả năng gặp phải trong đời sống thực. Chúng ta bắt đầu bằng việc giới thiệu một số định nghĩa, một số thuật ngữ mới trong ngôn ngữ thống kê mà anh/chị cần biết.

**ĐỊNH NGHĨA** • **Biến** là một đặc trưng thay đổi hay biến đổi theo thời gian, hay một đặc trưng mà biến đổi giữa các cá nhân hay các đối tượng khác nhau được xem xét tại một thời điểm nhất định •

Thí dụ, giá cổ phiếu là một biến thay đổi theo thời gian trong phạm vi một cổ phiếu đơn lẻ; nó cũng thay đổi từ cổ phiếu này sang cổ phiếu khác tại một thời điểm cho trước. Sự liên kết chính trị, nguồn gốc dân tộc, thu nhập, tuổi, và số con cái đều là biến – đó là những đặc trưng mà khác nhau tùy thuộc vào cá nhân được chọn.

Trong phần giới thiệu, chúng ta đã định nghĩa một **đơn vị thí nghiệm (experimental unit)** là đối tượng mà người ta lấy giá trị đo lường. Một cách tương đương, chúng ta có thể định nghĩa một đơn vị thí nghiệm là đối tượng mà trên đó một biến được đo lường. Khi một biến được đo lường thật sự trên một tập hợp các đơn vị thí nghiệm, thì một tập hợp các giá trị đo lường hay **dữ liệu** được tạo ra.

**ĐỊNH NGHĨA** • **Một đơn vị thí nghiệm** là cá nhân hay đối tượng mà trên đó một biến được đo lường. Một **giá trị đo lường** đơn lẻ hay một giá trị dữ liệu được tạo ra khi một biến được đo lường thật sự trên một đơn vị thí nghiệm •

Nếu một giá trị đo lường được tạo ra đối với mọi đơn vị thí nghiệm trong toàn bộ tập hợp, thì tập dữ liệu được tạo ra là **tổng thể** được quan tâm. Bất kỳ một tập hợp con nhỏ hơn nào của những giá trị đo lường cũng là một **mẫu**.

**THÍ DỤ 2.1** Một tập hợp gồm năm người làm công được chọn từ những người làm công tại một công ty lớn, và những giá trị đo lường sau đây được ghi chép. Hãy thảo luận về các biến được đo đối với năm người làm công này.

Người làm công	Điểm số về thành quả	Giới tính	Số năm phục vụ	Phân loại việc làm	Tiền lương (nghìn đô la)
1	18	Nữ	12	Bán hàng	35
2	15	Nữ	9	Quản lý	55
3	10	Nam	2	Hành chính	23
4	19	Nam	15	Quản lý	58
5	15	Nữ	13	Bán hàng	36

**Lời giải** Có một số biến trong thí dụ này. **Đơn vị thí nghiệm** mà trên đó mỗi biến được đo lường là một người làm công nhất định trong công ty. Đối với mỗi người làm công, có năm biến được đo lường: điểm số về thành quả, giới tính, số năm phục vụ, phân loại việc làm, và tiền lương. Mỗi trong những đặc trưng này thay đổi từ người làm công này sang người làm công khác. Nếu chúng ta xem những điểm số về thành quả của tất cả người làm công tại công ty này là tổng thể được quan tâm, thì năm điểm số về thành quả đó thể hiện một **mẫu** từ tổng thể này. Nếu như điểm số về thành quả của mỗi người làm công của công ty này đều được đo lường, thì chúng ta lẽ ra đã tạo ra toàn bộ **tổng thể** các giá trị đo lường cho biến này.

Biến thứ hai được đo lường trên những người làm công này là *giới tính*, mà có thể được xếp vào một trong hai loại – nam hay nữ. Nó không phải là một biến được đánh giá bằng số, và như thế nó có phần khác với *điểm số về thành quả*. Nếu có thể được nêu từng người, thì tổng thể sẽ gồm có một tập hợp những chữ Nam và Nữ, mỗi chữ đại diện cho mỗi người làm công tại công ty này. Tương tự, biến thứ tư, *phân loại việc làm*, tạo ra dữ liệu không phải bằng số, với một loại cho mỗi phân loại việc làm tại công ty này. Các biến thứ ba và thứ năm, *số năm đã làm việc* và *tiền lương*, đều được đánh giá bằng số, chúng ta tạo ra một tập hợp số chứ không phải một tập hợp các loại.

Mặc dù chúng ta đã thảo luận về từng biến một, hãy nhớ rằng chúng ta đã đo lường từng biến trong năm biến này trên năm đơn vị thí nghiệm – đó là năm người làm công. Vì thế, trong thí dụ này, một quan sát trên một cá nhân gồm có năm giá trị đo lường. Thí dụ, quan sát được thực hiện trên người làm công 2 mang lại kết quả đo lường sau đây:

(15, Nữ, 9, quản lý, 55) •

Anh/Chị có thể thấy rằng có sự khác biệt giữa một biến *đơn lẻ* được đo lường trên một đơn vị thí nghiệm đơn lẻ và *nhiều* biến được đo lường trên một đơn vị thí nghiệm đơn lẻ. Nếu một biến đơn lẻ được đo lường, thì dữ liệu tạo ra được gọi là **dữ liệu đơn biến**. Nếu hai biến được đo lường trên một đơn vị thí nghiệm đơn lẻ (chẳng hạn như giới tính và tiền lương), thì dữ liệu tạo ra được gọi là **dữ liệu nhị biến**. Nếu nhiều hơn hai biến được đo lường, như trong Thí dụ 2.1, thì dữ liệu được gọi là **dữ liệu đa biến**.

## 2.2 Các Loại Biến

Thí dụ 2.1 chứng tỏ rằng việc đo lường các biến tạo ra dữ liệu có thể bằng số hoặc không phải bằng số. Các biến mà dẫn đến dữ liệu không phải bằng số, trong đó các quan sát được phân loại dựa theo những điểm tương tự hay những điểm khác biệt về loại, thì được gọi là **biến định tính (qualitative variables)**. Sự liên kết chính trị, nghề nghiệp, tình trạng gia đình, và số năm học trung học phổ thông đều là những thí dụ về biến định tính, cũng như các biến “giới tính” và “phân loại việc làm” trong Thí dụ 2.1. Các biến được sử dụng để đo lường một đặc điểm mà tạo ra những quan sát bằng số thì được gọi là **biến định lượng (quantitative variables)**. Chỉ số Công nghiệp Dow–Jones, lãi suất cơ bản, số xe taxi không đăng ký ở một thành phố, mức sử dụng điện hàng ngày cho một nhà máy công nghiệp đều là những thí dụ về các biến định lượng, vốn dẫn đến dữ liệu định lượng.

**ĐỊNH NGHĨA** • **Các biến định lượng** dẫn đến các quan sát bằng số thể hiện một số lượng. **Các biến định tính** dẫn đến các quan sát không phải bằng số mà có thể được phân loại •

Những biến định lượng, mà thường được biểu hiện bằng chữ cái  $x$ , có thể được phân loại thêm nữa dựa vào miền giá trị bằng số mà một giá trị đo lường có thể có. Các biến, chẳng hạn như số thành viên trong các gia đình ở Arizona, doanh số xe hơi mới tại Trung tâm Mua sắm Xe hơi Riverfront, và số lốp xe có khiếm khuyết được trả lại cho nhà sản xuất để thay thế, có các giá trị tương ứng với một tập hợp con của số đếm 0, 1, 2, .... Cụ thể là các biến này có thể nhận một số có thể đếm được các giá trị và được gọi là **biến rời rạc (discrete variables)**. Cái tên *rời rạc* phản ánh thực tế là có những khoảng trống rời rạc giữa các giá trị khả dĩ mà dữ liệu có thể có. Mặt khác, những giá trị đo lường trên các biến chẳng hạn như chiều cao, trọng lượng, thời gian, khoảng cách, hay thể tích có thể có những giá trị tương ứng với tất cả các điểm trên một khoảng vạch (line interval). Loại biến này được gọi là **biến liên tục (continuous variables)**. Giữa bất kỳ hai giá trị nào của một biến liên tục, luôn luôn có thể tìm thấy một giá trị thứ ba.

**ĐỊNH NGHĨA** • **Biến liên tục** là một biến có thể nhận tất cả giá trị nhiều vô hạn tương ứng với một khoảng vạch. **Biến rời rạc** chỉ có thể nhận một số có thể đếm được các giá trị •

**THÍ DỤ 2.2** Hãy xác định mỗi biến trong các biến sau đây là *định tính* hay *định lượng*.

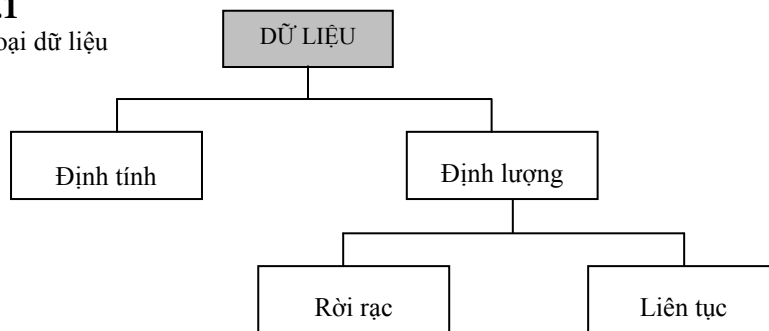
- a. Mục đích sử dụng thường xuyên nhất của lò vi ba của Anh/Chị (hâm lại, làm hết đông lạnh, đun nóng, mục đích khác) suốt tháng 12/2000.
- b. Số người tiêu dùng từ chối trả lời cuộc điều tra bằng điện thoại.
- c. Loại dịch vụ cấp được cung cấp cho nơi cư trú (cấp tiêu chuẩn, cấp cao cấp hay chỉ có anten) ở Atlanta.
- d. Thời gian hoàn tất đối với một nhiệm vụ nào đó được thực hiện bởi một chương trình phần mềm máy tính.
- e. Số cổ phiếu ở Sở Giao dịch Chứng khoán New York cho thấy có tăng giá từ 1/3/2000 đến 1/7/2000.

**Lời giải** Các biến (a) và (c) đều là biến định tính, bởi vì chỉ có một đặc điểm được đo lường trên mỗi đơn vị thí nghiệm. Các loại đối với hai biến này được trình bày trong các ngoặc đơn. Ba biến còn lại là biến định lượng. Số người tiêu dùng là biến rời rạc; nó có thể nhận bất kỳ giá trị nào trong các giá trị 0, 1, 2, ..., với giá trị tối đa phụ thuộc vào số người tiêu dùng được gọi điện thoại phỏng vấn. Tương tự, số cổ phiếu cho thấy có tăng giá có thể nhận bất kỳ giá trị nào trong các giá trị 0, 1, 2, ..., với giá trị tối đa phụ thuộc vào số cổ phiếu ở Sở Giao dịch Chứng khoán New York. Biến (d), thời gian hoàn tất đối với một nhiệm vụ nào đó, là biến liên tục duy nhất trong danh sách ở trên. Thời gian hoàn tất có thể là 121 giây, 121,25 giây, hay một giá trị nằm giữa hai giá trị bất kỳ được liệt kê. •

Tại sao chúng ta phải quan tâm đến các loại khác nhau của biến và dữ liệu chúng tạo ra? Các kỹ thuật được sử dụng để tổng hợp (summarizing) và mô tả các tập dữ liệu phụ thuộc vào loại dữ liệu được thu thập. Dữ liệu định tính thường được tổng hợp bằng cách xác định số lượng hay tỷ lệ những quan sát trong mỗi một trong một số loại. Sau đó các kết quả được biểu hiện bằng cách sử dụng bảng và đồ thị. Những biểu hiện bằng đồ thị có phần khác nhau đối với các biến định lượng rời rạc và liên tục, nhưng nhìn chung chúng tập trung vào những đồ thị trong đó số quan sát trong một lớp hay loại được vẽ theo các lớp hay các loại. Đối với mỗi tập dữ liệu Anh/Chị gặp phải, thì kỹ xảo sẽ là xác định loại dữ liệu nào liên quan và làm sao anh/chị có thể biểu hiện nó theo một cách thức rõ ràng và có thể hiểu được đối với cử tọa của mình (xem Hình 2.1)

HÌNH 2.1

Các loại dữ liệu



### 2.3 Các Phương pháp Bảng số để Mô tả Một Tập Dữ liệu

Các phương pháp bằng đồ thị hết sức hữu ích trong việc biểu hiện dữ liệu và trong việc truyền tải sự mô tả tổng quát và nhanh chóng về dữ liệu được thu thập. Điều này chứng minh, trong nhiều khía cạnh, cho câu tục ngữ một bức họa đáng giá cả ngàn từ. Tuy nhiên, có những hạn chế đối với việc sử dụng kỹ thuật bằng đồ thị để mô tả và phân tích dữ liệu. Ví dụ như, giả sử chúng ta muốn thảo luận về dữ liệu của mình trước một nhóm người và không có sẵn máy chiếu phóng đại! Chúng ta sẽ buộc phải sử dụng những thước đo mô tả khác mà sẽ truyền tải cho người nghe một hình ảnh trong trí óc về biểu đồ tần suất. Một hạn chế thứ hai và không thật là hiển nhiên của biểu đồ tần suất và các kỹ thuật bằng đồ thị khác, đó là chúng khó sử dụng nhằm những mục đích về suy luận thống kê



(statistical inference). Giả sử chúng ta sử dụng biểu đồ tần suất của mẫu để đưa ra những suy luận về hình dạng và vị trí của biểu đồ tần suất của tổng thể, dùng để mô tả tổng thể này và chúng ta chưa biết. Sự suy luận của chúng ta dựa vào giả định đúng, đó là một mức độ tương tự nào đó sẽ tồn tại giữa hai biểu đồ tần suất này, nhưng rồi chúng ta phải đối mặt với vấn đề đo lường mức độ tương tự này. Chúng ta biết rõ khi hai hình vẽ giống hệt nhau, nhưng tình hình này sẽ không có khả năng xảy ra trong thực tiễn. Nếu chúng giống hệt nhau, chúng ta có thể nói “Chúng giống nhau.” Nhưng, nếu chúng khác nhau, thì khó mà mô tả được “mức độ khác biệt.”

Những hạn chế của phương pháp mô tả dữ liệu bằng đồ thị có thể được khắc phục bằng việc sử dụng những **thước đo mô tả bằng số**. Thước đo mô tả bằng số dành cho một tổng thể được gọi là **tham số**. Thước đo mô tả bằng số tương ứng được tính toán từ một mẫu thì được gọi là **trị thống kê**. Như thế, chúng ta muốn sử dụng dữ liệu của mẫu để tính toán một tập hợp các con số, các trị thống kê, mà sẽ truyền tải một hình ảnh trong trí óc thật tốt về phân phối tần suất tương đối của mẫu và sẽ hữu ích trong việc đưa ra những suy luận về phân phối tần suất tương đối của tổng thể.

**ĐỊNH NGHĨA** • Các thước đo mô tả bằng số được tính từ những giá trị đo lường của tổng thể được gọi là **tham số** •

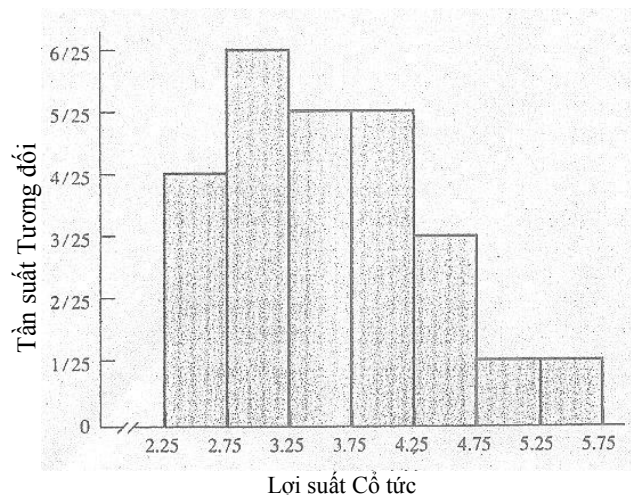
**ĐỊNH NGHĨA** • Các thước đo mô tả bằng số được tính từ những giá trị đo lường của mẫu được gọi là **trị thống kê** •

## 2.4 Các Thước đo Hướng Tâm

Trong việc xây dựng một hình ảnh trong trí óc về phân phối tần suất cho một tập hợp giá trị đo lường trên một biến định lượng,  $x$ , chúng ta rất có thể hình dung ra một biểu đồ tần suất tương tự với biểu đồ được trình bày trong Hình 2.2, đối với dữ liệu về lợi suất cổ tức của cổ phiếu ngân hàng. Một trong những thước đo mô tả đầu tiên được quan tâm là **thước đo hướng tâm (measure of central tendency)**, đó là một thước đo, chẳng hạn như một số trung bình, xác định vị trí trung tâm của phân phối. Chúng ta lưu ý rằng lợi suất cổ tức thay đổi trong khoảng từ mức thấp là 2,3 lên mức cao là 5,3, với trung tâm của biểu đồ tần suất nằm gần 3,6. Bây giờ chúng ta hãy xem xét một số quy tắc rõ ràng để xác định vị trí trung tâm của một phân phối dữ liệu.

HÌNH 2.2

Biểu đồ Tần suất Tương đối



Một trong những thước đo hướng tâm hữu ích và thông dụng nhất, đó là trị số trung bình số học của một tập hợp các giá trị đo lường. Trị số này thường cũng được gọi là **trung bình số học (arithmetic mean)**, hay chỉ đơn giản là **trung bình (mean)**, của một tập hợp các giá trị đo lường. Bởi vì chúng ta sẽ muốn phân biệt giữa trung bình của một mẫu và trung bình của một tổng thể, nên chúng ta sẽ sử dụng ký hiệu  $\bar{x}$  ( $x$  gạch ngang trên đầu) để biểu hiện trung bình của mẫu và  $\mu$  (chữ mu thường của Hy Lạp) để biểu hiện trung bình của tổng thể.

**ĐỊNH NGHĨA** • **Trung bình số học** của một tập hợp các giá trị đo lường bằng tổng số các giá trị đo lường này chia cho số lượng giá trị đo lường •

Những quy trình tính toán trung bình mẫu và nhiều trị thống kê khác được thể hiện một cách thuận lợi thành các công thức. Do vậy, chúng ta sẽ cần một ký hiệu để biểu hiện quy trình tính tổng số. Nếu chúng ta biểu thị  $n$  số lượng phải được tính tổng số là  $x_1, x_2, \dots, x_n$ , thì tổng số của chúng được biểu thị bằng ký hiệu

$$\sum_{i=1}^n x_i$$

Chữ sigma viết hoa của Hy Lạp ( $\Sigma$ ) là chỉ dẫn *cộng lại*. Số lượng  $x_i$  ở bên phải của  $\Sigma$  là phân tử tiêu biểu sẽ được cộng lại. Những ký hiệu  $i = 1$  ở dưới và  $n$  ở bên trên chữ  $\Sigma$  chỉ ra rằng  $i$  là biến của phép tính tổng số và bắt đầu bằng trị số 1, tăng dần thêm 1, và kết thúc bằng trị số  $n$ . Thí dụ,

$$\sum_{i=1}^3 x_i = x_1 + x_2 + x_3$$

Sử dụng ký hiệu này, chúng ta có thể biểu hiện các công thức cho trung bình mẫu và trung bình tổng thể như sau:

### Các Công thức Tính Trị số Trung bình

$$\text{Trung bình mẫu: } \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\text{Trung bình tổng thể: } \mu = \frac{\sum_{i=1}^N x_i}{N}$$

**THÍ DỤ 2.3** Tìm trung bình của tập hợp các giá trị đo lường 2, 9, 11, 5, 6.

#### Lời giải

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{2+9+11+5+6}{5} = 6,6$$

Thậm chí quan trọng hơn việc xác định vị trí trung tâm của một tập hợp các giá trị đo lường của mẫu,  $\bar{x}$  sẽ được sử dụng làm một hàm ước lượng (hàm tiên đoán) về giá trị của trung bình  $\mu$  chưa biết của tổng thể. Thí dụ, trung bình của dữ liệu trong Bảng 2.3 bằng

BẢNG 2.3

Lợi suất cổ tức (%) đối với 25 cổ phiếu thường của ngân hàng	3,1	4,2	2,3	3,3	2,8
	5,3	3,5	3,1	2,6	3,3
	4,7	3,7	3,0	2,6	4,0
	3,8	4,4	3,2	3,2	3,8
	5,1	3,7	2,3	4,3	3,9

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{89,2}{25} = 3,568$$

Hãy lưu ý giá trị này xấp xỉ rơi vào trung tâm của tập hợp các giá trị đo lường. Trung bình của toàn bộ tổng thể lợi suất cổ tức,  $\mu$ , thì chúng ta chưa biết; nhưng nếu như chúng ta ước lượng giá trị của nó, thì giá trị ước lượng của chúng ta về  $\mu$  sẽ là 3,586.

Một thước đo hướng tâm thứ hai là **trung vị**.

**ĐỊNH NGHĨA** • **Trung vị**  $m$  của một tập hợp  $n$  giá trị đo lường  $x_1, x_2, x_3, \dots, x_n$  là giá trị của  $x$  mà nằm ở giữa khi các giá trị đo lường này được xếp theo thứ tự từ nhỏ nhất đến lớn nhất •

Nếu các giá trị đo lường trong một tập dữ liệu được xếp từ nhỏ nhất đến lớn nhất, thì trung vị sẽ là giá trị của  $x$  nằm ở giữa. Nếu số  $n$  giá trị đo lường là lẻ, thì số trung vị sẽ là giá trị đo lường có thứ hạng bằng  $(n + 1)/2$ . Nếu số  $n$  giá trị đo lường là chẵn, thì số trung vị được chọn là giá trị của  $x$  nằm ở điểm giữa hai giá trị đo lường ở khoảng giữa – đó là ở điểm giữa giá trị đo lường có thứ hạng  $n/2$  và giá trị đo lường có thứ hạng  $(n / 2) + 1$ . Quy tắc tính toán số trung vị được trình bày trong hộp sau đây:

### Quy tắc Tính toán Số Trung vị

Xếp hạng  $n$  giá trị đo lường từ nhỏ nhất đến lớn nhất

1. Nếu  $n$  lẻ, số trung vị  $m$  là giá trị đo lường có thứ hạng  $(n + 1)/2$
2. Nếu  $n$  chẵn, số trung vị  $m$  là giá trị của  $x$  nằm ở điểm giữa giá trị đo lường có thứ hạng  $n/2$  và giá trị đo lường có thứ hạng  $(n/2) + 1$ .

**THÍ DỤ 2.4** Hãy tìm số trung vị của tập hợp năm giá trị đo lường sau đây.

9, 2, 7, 11, 14

**Lời giải** Trước tiên, chúng ta xếp hạng  $n = 5$  giá trị đo lường từ nhỏ nhất đến lớn nhất, 2, 7, 9, 11, 14. Như thế, vì  $n = 5$  là số lẻ, nên chúng ta chọn 9 là số trung vị. Giá trị này là giá trị đo lường có thứ hạng là  $(n + 1)/2 = (5 + 1)/2 = 3$  •

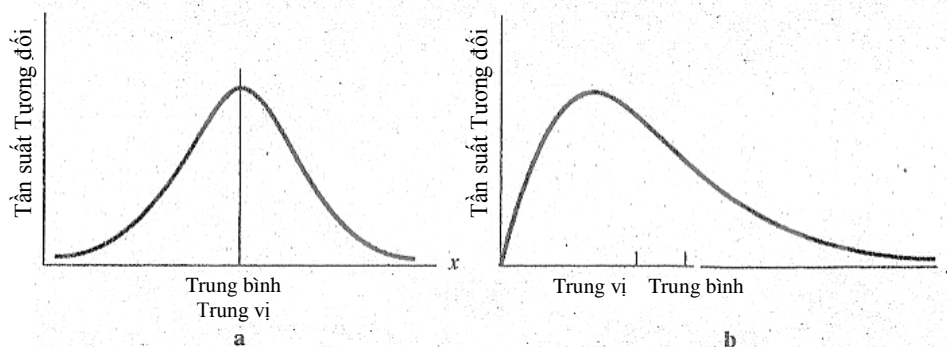
**THÍ DỤ 2.5** Hãy tìm số trung vị của tập hợp các giá trị đo lường sau đây.

9, 2, 7, 11, 14, 6

**Lời giải** Vì  $n = 6$  là số chẵn, nên chúng ta xếp hạng các giá trị đo lường thành 2, 6, 7, 9, 11, 14 và chọn số trung vị là điểm giữa của hai giá trị đo lường ở khoảng giữa, 7 và 9. Vì thế, số trung vị bằng 8 •

Mặc dù cả trung bình lẫn trung vị đều là hai thước đo tốt đối với trung tâm của một phân phối các giá trị đo lường, nhưng trung vị kém nhạy cảm với những giá trị thái cực (cực trị). Thí dụ, nếu phân phối này đối xứng qua trung bình của nó – nghĩa là hai nửa bên trái và bên phải của phân phối này là những hình ảnh phản chiếu – thì số trung bình và số trung vị bằng nhau [xem Hình 2.3 (a)]. Nếu một phân phối không đối xứng và có những quan sát thái cực nằm ở đuôi bên phải của phân phối này, thì phân phối này được gọi là bị *lệch xiên về bên phải* [xem Hình 2.3(b)]. Bởi vì những giá trị thái cực lớn ở đuôi trên của phân phối này làm tăng tổng số của các giá trị đo lường, nên số trung bình chuyển dịch sang phải. Số trung vị không bị ảnh hưởng bởi những giá trị thái cực này, bởi vì giá trị bằng số của các giá trị đo lường không được sử dụng trong việc tính toán số trung vị. Cuối cùng, nếu một phân phối bị *lệch xiên về bên trái*, thì số trung bình chuyển dịch sang trái.

**HÌNH 2.3**  
 Các phân phối tần suất tương đối cho thấy tác động của các giá trị thái cực đối với trung bình và trung vị



Một thước đo hướng tâm khác là **yếu vị (cao tần)**, được định nghĩa là giá trị quan sát xảy ra thường nhất trong một tập dữ liệu.

**ĐỊNH NGHĨA** • **Số Yếu vị** của một tập hợp  $n$  giá trị đo lường  $x_1, x_2, x_3, \dots, x_n$  là giá trị của  $x$  xảy ra với tần suất lớn nhất •

Khi các giá trị đo lường được phân nhóm trong một biểu đồ tần suất tương đối, thì lớp có tần suất tương đối lớn nhất được gọi là **lớp yếu vị**, và điểm giữa của lớp yếu vị được lấy làm giá trị của yếu vị

**THÍ DỤ 2.6** Cho trước những giá trị đo lường của mẫu

5, 5, 7, 7, 7, 10, 15

giá trị 7 xảy ra ba lần, giá trị 5 xảy ra hai lần, và các giá trị 10 và 15 thì mỗi số xảy ra một lần. Vì thế, số yếu vị của những giá trị đo lường của mẫu này là 7.

Đối với những phân phối đối xứng, thì các số trung bình, trung vị, và yếu vị đều bằng nhau. Trong những phân phối bị lệch xiên về bên phải, số yếu vị nằm bên trái số trung vị và số trung bình. Xem các Hình 2.3(a) và 2.3(b). Nếu phân phối bị lệch về bên trái, thì vị trí của ba thước đo này được đảo ngược, số yếu vị nằm bên phải số trung bình và số trung vị.

Một phân phối các giá trị đo lường có thể có nhiều hơn 1 số yếu vị. Thí dụ, việc phân phối tiền lương đối với một nhóm nhiều người làm công có thể tạo ra một *phân phối có hai yếu vị*, có thể phản ánh một hỗn hợp các giá trị đo lường được lấy trên những người làm công cỏ xanh và cỏ trắng.

## Bài tập

### Các Kỹ thuật Căn bản

2.1 Hãy xét  $n = 5$  giá trị đo lường, 0, 5, 1, 1, 3.

- Hãy vẽ một đồ thị phân tán cho dữ liệu này. [*Gợi ý*: Nếu hai giá trị đo lường giống nhau, hãy đặt chấm này ở trên chấm kia]. Hãy phỏng đoán “trung tâm” xấp xỉ.
- Hãy tìm số trung bình, số trung vị, và số yếu vị.
- Hãy xác định vị trí của ba thước đo vừa tìm ra trong phần (b) trên đồ thị phân tán trong phần (a). Dựa trên các vị trí tương đối của số trung bình và số trung vị, Anh/Chị cho là những giá trị đo lường này đối xứng hay bị lệch xiên?

2.2 Hãy xét  $n = 8$  giá trị đo lường, 3, 1, 5, 4, 4, 3, 5.

- Tìm  $\bar{x}$
- Tìm  $m$
- Dựa trên kết quả của các phần (a) và (b), những giá trị đo lường này bị lệch xiên hay đối xứng? Hãy vẽ đồ thị phân tán để xác nhận câu trả lời của anh/chị

2.3 Cho trước  $n = 10$  giá trị đo lường, 3, 5, 4, 6, 10, 5, 6, 9, 2, 8, hãy tìm:

- $\bar{x}$
- $m$
- số yếu vị

### Ứng dụng

2.4 Nhiều người mua máy tính đã phát hiện ra rằng họ có thể tiết kiệm được một số tiền đáng kể bằng việc mua máy tính cá nhân từ một công ty nhận đặt và giao hàng qua đường bưu điện – trung bình là 900USD theo giá trị ước lượng của họ. (“Who’s Tops,” 1992). Điểm xếp hạng về sự thỏa mãn của khách hàng (trên thang đo từ 1 đến 9) đối với bảy công ty như thế, dựa trên cuộc điều tra 4.000 người mua, được trình bày dưới đây.

Công ty	Xếp hạng	Công ty	Xếp hạng
CompuAdd	7,5	Insight	7,8
Dell	7,9	Northgate	7,7
FastMicro	7,4	Zeos	8,0
Gateway	8,2		

- Điểm xếp hạng trung bình về sự thỏa mãn của khách hàng đối với bảy công ty này là bao nhiêu?
- Hãy cho biết số trung vị của những điểm xếp hạng về sự thỏa mãn của khách hàng.
- Nếu anh/chị là một người mua máy tính, anh/chị có quan tâm đến điểm xếp hạng trung bình về sự thỏa mãn của khách hàng hay không? Nếu không, thước đo nào anh/chị quan tâm? Hãy giải thích.

**2.5** Thu nhập bình quân mỗi cổ phiếu trong quý hai, năm 1994, đối với một mẫu gồm 20 công ty được trình bày dưới đây:

\$ 0,72	0,56	0,21	0,54	0,32
1,28	0,10	1,64	0,29	0,33
0,29	0,73	0,29	0,33	0,43
0,56	0,89	0,84	0,62	0,44

Nguồn: Dữ liệu trích từ *Press-Enterprise*, Riverside, Calif, 20 tháng 7, 1994

- Anh/Chị cho rằng phân phối về thu nhập bình quân mỗi cổ phiếu là đối xứng hay bị lệch xiên?
- Hãy tính số trung bình, số trung vị và số yếu vị cho những giá trị ước lượng này.
- Hãy vẽ một biểu đồ tần suất tương đối cho tập dữ liệu này. Hãy xác định vị trí của số trung bình, số trung vị và số yếu vị dọc theo trục hoành. Câu trả lời của anh/chị đối với phần (a) có đúng hay không?

**2.6** Tạp chí *PC World* cung cấp một nguồn thông tin tuyệt vời cho những người sử dụng máy tính muốn nâng cấp hệ điều hành hiện tại của họ hay mua những hệ điều hành mới. Số gần đây của tạp chí *PC World* (“Top 10,” 1994) đã liệt kê mười bộ tăng tốc dựa trên Windows hàng đầu, cùng với điểm xếp hạng giá trị toàn bộ và giá ngoài đường ước lượng, như được trình bày trong bảng sau đây:

Bộ tăng tốc	Điểm Xếp hạng Giá trị Toàn bộ	Giá Ngoài đường Ước lượng
Diamond Stealth	87	\$249
Number Nine	86	275
Genoa Phantom	85	245
Hercules Dynamite Pro	82	210
miroCrystal8S	82	195
Orchid Kelvin	75	275
Hercules Graphite	73	335
Matrox MGA	73	475
Hercules Dynamite Power	72	237
Paradis Ports o’Call	72	235

- Điểm xếp hạng giá trị toàn bộ trung bình cho mười sản phẩm này là bao nhiêu?
- Giá ngoài đường ước lượng trung bình là bao nhiêu?
- Nếu anh/chị sắp mua một bộ tăng tốc, thì những số trung bình này có quan trọng đối với anh/chị hay không? Hãy giải thích.

**2.7** Trong một bài báo với nhan đề “Bạn không phải là người bị hoang tưởng nếu bạn nghĩ ai đó quan sát mọi hành động của bạn”, *Tạp chí Phố Wall* (19 tháng 3, 1985) lưu ý rằng doanh nghiệp lớn thu thập số liệu thống kê chi tiết về hành vi của bạn. Công ty Jockey International biết bạn có bao nhiêu quần lót; công ty Frito-Lay, Inc., biết bạn ăn cái nào trước – những cái

bánh quy bẽ trong gói bánh hay những cái còn nguyên; và bắt tay vào những điều cụ thể, công ty Coca-Cola biết rằng bạn bỏ vào ly 3,2 cục nước đá. Bạn đã bao giờ bỏ 3,2 cục nước đá vào cái ly chưa? Bài báo của *Tạp chí Phố Wall* muốn nói gì qua lời phát biểu đó?

**2.8** Bảng sau đây trình bày nợ bình quân đầu người đối với từng bang trong 50 bang trong năm tài chính 1992.

Bang	Nợ bình quân đầu người	Thuế bình quân đầu người	Bang	Nợ bình quân đầu người	Thuế bình quân đầu người
AL	998	1019	MT	2266	1153
AK	8418	2730	NE	1092	1176
AZ	743	1259	NV	1457	1369
AR	809	1145	NH	3882	770
CA	1225	1495	NJ	2540	1643
CO	857	1018	NM	1015	1415
CT	3644	1846	NY	3083	1661
DE	5140	1944	NC	558	1316
FL	911	1068	ND	1615	1186
GA	662	1076	OH	1106	1099
HI	4040	2335	OK	1138	1206
ID	1210	1303	OR	2114	1113
IL	1611	1157	PA	1079	1354
IN	913	1143	RI	5125	1270
IA	669	1280	SC	1300	1092
KS	192	1110	SD	2657	794
KY	1762	1353	TN	558	900
LA	2331	991	TX	453	964
ME	2135	1347	UT	1187	1096
MD	1698	1324	VT	2706	1339
MA	4002	1651	VA	1160	1101
MI	1097	1195	WA	1400	1648
MN	924	1662	WV	1431	1297
MS	621	954	WI	1457	1380
MO	1213	988	WY	1920	1386

Nguồn: Dữ liệu từ Bộ Thương mại Hoa Kỳ, Cục Điều tra Dân số, *The World Almanac and Book of Facts*, ấn bản 1994, trang 105

- a Hãy tìm số nợ bình quân đầu người trung bình cho 50 bang.
- b Hãy tìm số nợ bình quân đầu người trung vị cho 50 bang này và so sánh nó với số trung vị đã tính trong phần (a)
- c Dựa trên sự so sánh của anh/chị trong phần (b), anh/chị có kết luận rằng phân phối của nợ bình quân đầu người bị lệch xiên? Hãy giải thích.

**2.9** Việc định giá đơn vị đã trở thành một tiêu chuẩn toàn ngành trong hoạt động kinh doanh tạp hóa. Công việc của người tiêu dùng là cân nhắc chất lượng của sản phẩm so với giá đơn vị để cố gắng xác định “món hời nhất”. Những giá trị đo lường (measurements) sau đây là giá mỗi túi nhựa lót thùng rác, được ghi nhận đối với 10 nhãn hiệu khác nhau của túi nhựa lót thùng rác 13-gallon và cao (*Báo cáo Người tiêu dùng*, Tháng 2, 1994).

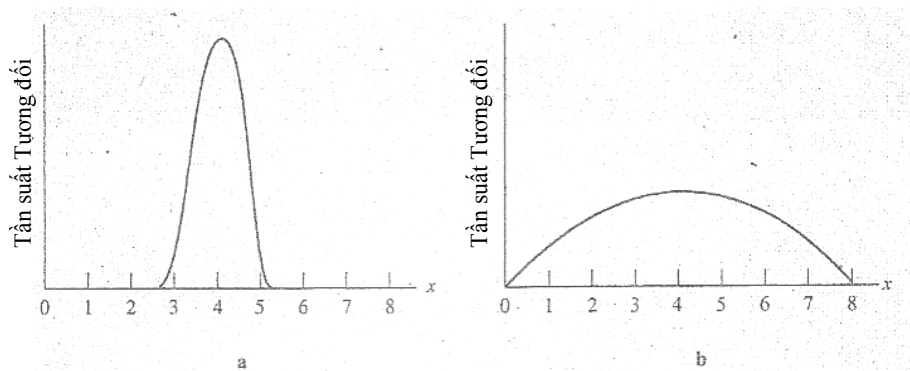
10	9	13	8	9
10	10	6	5	11

- a Hãy tìm giá trung bình mỗi túi nhựa lót
- b Hãy tìm giá trung vị mỗi túi nhựa lót
- c Nếu anh/chị đang viết báo cáo để mô tả những dữ liệu này, thước đo hướng tâm nào anh/chị sẽ sử dụng? Hãy giải thích.

## 2.5 Những Thước đo về Độ Biến thiên

Một khi chúng ta đã xác định trung tâm của một phân phối dữ liệu, bước tiếp theo là cung cấp một thước đo về **độ biến thiên (variability)**, hay **độ phân tán (dispersion)**, của dữ liệu này. Hãy xét hai phân phối được trình bày trong Hình 2.4. Cả hai phân phối đều được đặt ở vị trí có trung tâm tại  $x = 4$ , nhưng có sự khác biệt lớn về độ biến thiên của những giá trị đo lường xung quanh số trung bình đối với hai phân phối này. Các giá trị đo lường trong Hình 2.4(a) thay đổi xấp xỉ từ 3 đến 5; trong Hình 2.4(b), các giá trị đo lường thay đổi từ 0 đến 8.

**HÌNH 2.4**  
Độ biến thiên hay độ phân tán của dữ liệu



Sự biến thiên là một đặc trưng quan trọng của dữ liệu. Thí dụ, nếu chúng ta đang chế tạo bu lông, thì sự biến thiên quá mức trong đường kính của bu lông sẽ kéo theo một tỷ lệ phần trăm cao của sản phẩm có khiếm khuyết. Mặt khác, khi chúng ta sử dụng một bài kiểm tra để phân biệt giữa những kẻ toán viên giỏi và kém, thì chúng ta sẽ không vui nhất nếu bài kiểm tra này lúc nào cũng mang lại những điểm kiểm tra với ít biến thiên, bởi vì điều này sẽ làm cho việc phân biệt trở nên rất khó khăn.

Ngoài tầm quan trọng trên thực tế của sự biến thiên trong dữ liệu, một thước đo về đặc trưng này còn cần thiết cho việc xây dựng một hình ảnh trong trí óc về phân phối tần suất. Chúng ta sẽ chỉ thảo luận về vài trong số nhiều thước đo về sự biến thiên.

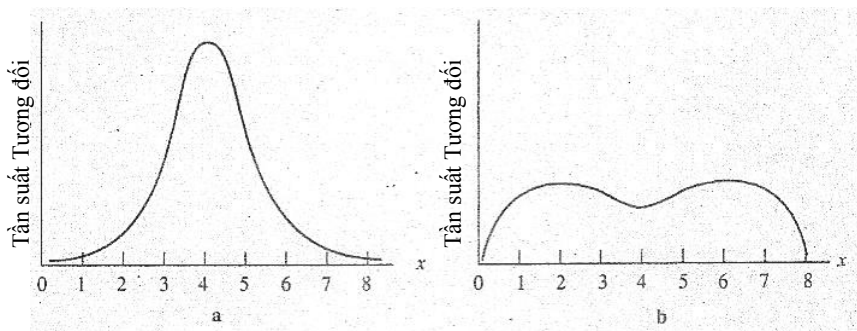
Thước đo đơn giản nhất về sự biến thiên là **khoảng biến thiên (miền)**.

**ĐỊNH NGHĨA** • **Khoảng biến thiên (range)** của một tập hợp  $n$  giá trị đo lường  $x_1, x_2, x_3, \dots, x_n$  được định nghĩa là chênh lệch giữa giá trị đo lường lớn nhất và giá trị đo lường nhỏ nhất •

Dữ liệu về lợi suất cổ tức thay đổi từ 2,3 đến 5,3. Như thế, khoảng biến thiên là  $(5,3 - 2,3) = 3,0$ . Khoảng biến thiên dễ tính toán, dễ diễn giải, và hoàn toàn thỏa đáng trong vai trò một thước đo về sự biến thiên cho những tập dữ liệu nhỏ. Nhưng đối với những tập dữ liệu lớn thì khoảng biến thiên không phải là một thước đo thỏa đáng về độ biến thiên. Thí dụ, hai phân phối tần suất tương đối trong Hình 2.5 có cùng khoảng biến thiên nhưng lại có hình dạng và độ biến thiên rất khác nhau.



**HÌNH 2.5**  
Những phân phối có  
khoảng biến thiên  
bằng nhau và độ biến  
thiên khác nhau



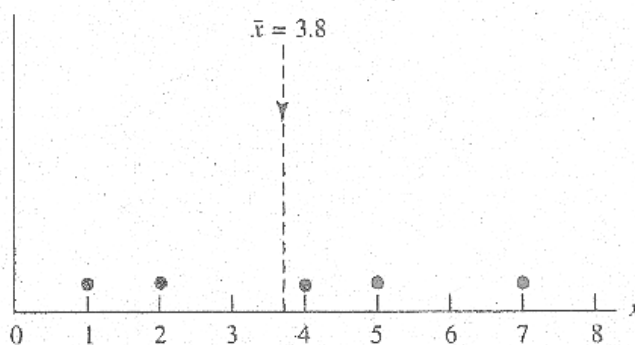
Chúng ta có thể tìm ra một thước đo về độ biến thiên nhạy cảm hơn khoảng biến thiên hay không? Lấy thí dụ, hãy xét các giá trị đo lường của mẫu 5, 7, 1, 2, 4, được biểu hiện thành đồ thị phân tán trong Hình 2.6. Số trung bình của năm giá trị đo lường này là

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{19}{5} = 3,8$$

như được chỉ ra trong đồ thị phân tán này

Bây giờ chúng ta có thể xem xét độ biến thiên theo khoảng cách giữa mỗi điểm (giá trị đo lường) và trung bình  $\bar{x}$ . Nếu những khoảng cách này lớn thì chúng ta có thể nói rằng dữ liệu biến thiên nhiều hơn so với khi những khoảng cách này nhỏ. Nói rõ hơn, chúng ta định nghĩa **độ lệch** của một giá trị đo lường khỏi số trung bình của nó là lượng  $(x_i - \bar{x})$ . Những giá trị đo lường nằm bên phải của số trung bình tạo ra độ lệch dương, và những giá trị đo lường nằm bên trái tạo ra độ lệch âm. Đối với thí dụ của chúng ta, các giá trị của  $x$  và các độ lệch được trình bày trong cột thứ nhất và cột thứ hai của Bảng 2.4.

**HÌNH 2.6**  
Đồ thị phân tán



**BẢNG 2.4**  
Những phép tính  
liên quan đến độ  
lệch của mẫu

$x_i$	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$	$ x_i - \bar{x} $
5	1,2	1,44	1,2
7	3,2	10,24	3,2
1	-2,8	7,84	2,8
2	-1,8	3,24	1,8
4	0,2	0,04	0,2
19	0,0	22,80	9,2

Nếu chúng ta đồng ý rằng các độ lệch chứa đựng thông tin về sự biến thiên, thì bước tiếp theo của chúng ta là xây dựng một thước đo về sự biến thiên dựa trên các độ lệch xung quanh số trung bình. Khả năng đầu tiên là chúng ta có thể chọn trung bình của các độ lệch. Đáng tiếc là trung bình này sẽ không có tác dụng, bởi vì một số độ lệch thì dương, một số thì âm, và tổng số luôn luôn bằng không (trừ khi những sai số làm tròn số đã được đưa vào các phép tính). Hãy lưu ý rằng các độ lệch trên cột thứ hai của Bảng 2.4 có tổng bằng không.

Có hai cách để tránh được vấn đề này. Tại sao không tính số trung bình của giá trị tuyệt đối của các độ lệch? Thước đo này được gọi là **độ lệch tuyệt đối trung bình (mean absolute deviation, MAD)**.

**ĐỊNH NGHĨA** • **Độ lệch tuyệt đối trung bình** của một tập hợp  $n$  giá trị đo lường  $x_1, x_2, \dots, x_n$  là số trung bình của giá trị tuyệt đối của các độ lệch xung quanh trung bình mẫu và được cho bởi công thức

$$\text{MAD} = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n} \quad \bullet$$

Các độ lệch tuyệt đối của tập hợp  $n = 5$  giá trị quan sát (observations) của chúng ta cùng với tổng số của chúng được trình bày trong Bảng 2.4. Vì thế cho nên,

$$\text{MAD} = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n} = \frac{9,2}{5} = 1,84$$

Mặc dù MAD đôi khi được sử dụng làm thước đo về độ biến thiên cho một tập dữ liệu, nhưng nó chủ yếu được dùng trong việc đánh giá độ chính xác của tiên đoán.

Cách thứ hai để sử dụng độ lệch là làm việc với tổng các bình phương của độ lệch. Bằng việc sử dụng tổng của các độ lệch bình phương, chúng ta tính toán một thước đo đơn lẻ gọi là **phương sai (variance)** của một tập hợp các giá trị đo lường. Để phân biệt giữa phương sai của *mẫu* và phương sai của *tổng thể*, chúng ta sử dụng ký hiệu  $s^2$  để biểu hiện phương sai mẫu và  $\sigma^2$  (chữ sigma thường của Hy Lạp) để biểu hiện phương sai tổng thể. *Thước đo này sẽ tương đối lớn đối với dữ liệu biến thiên nhiều và tương đối nhỏ đối với dữ liệu biến thiên ít.*

- ĐỊNH NGHĨA** • **Phương sai của tổng thể** gồm  $N$  giá trị đo lường  $x_1, x_2, \dots, x_N$  được định nghĩa là trị trung bình của các bình phương của độ lệch của các giá trị đo lường xung quanh số trung bình  $\mu$  của chúng. Phương sai của tổng thể (phương sai tổng thể) được cho bởi công thức

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \quad \bullet$$

Rất thường là anh/chị sẽ không có sẵn tất cả các giá trị đo lường của tổng thể, mà sẽ cần tính toán *phương sai của một mẫu* gồm  $n$  giá trị đo lường.

- ĐỊNH NGHĨA** • **Phương sai của mẫu** gồm  $n$  giá trị đo lường  $x_1, x_2, \dots, x_n$  được định nghĩa là tổng các độ lệch bình phương của các giá trị đo lường này xung quanh số trung bình  $\bar{x}$  của chúng, chia cho  $(n - 1)$ . Phương sai mẫu được ký hiệu bằng chữ  $s^2$  và được cho bởi công thức

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

Thí dụ, chúng ta có thể tính phương sai cho tập hợp gồm  $n = 5$  giá trị đo lường của mẫu, được trình bày trong Bảng 2.4. Bình phương của độ lệch của mỗi giá trị đo lường được ghi trên cột thứ ba của Bảng 2.4. Cộng lại, chúng ta thu được

$$\sum_{i=1}^5 (x_i - \bar{x})^2 = 22,80$$

Phương sai mẫu là

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = \frac{22,80}{4} = 5,70$$

Phương sai được đo theo bình phương của các đơn vị đo lường ban đầu. Nếu các giá trị đo lường ban đầu được tính bằng inch, thì phương sai được biểu hiện bằng inch bình phương. Lấy căn bậc hai của phương sai này, thì chúng ta có **độ lệch chuẩn (standard deviation)**, mà chuyển thước đo độ biến thiên trở lại các đơn vị đo lường ban đầu.

- ĐỊNH NGHĨA** • **Độ lệch chuẩn** của một tập hợp các giá trị đo lường bằng căn bậc hai dương của phương sai. •

### Hệ thống ký hiệu

$n$ : số giá trị đo lường trong mẫu

$s^2$ : phương sai mẫu

$s = \sqrt{s^2}$ : độ lệch chuẩn của mẫu

$N$ : số giá trị đo lường trong tổng thể

$\sigma^2$ : phương sai tổng thể

$\sigma = \sqrt{\sigma^2}$ : độ lệch chuẩn của tổng thể

Đối với tập hợp  $n = 5$  giá trị đo lường của mẫu trong Bảng 2.4, phương sai mẫu là  $s^2 = 5,70$ , do đó độ lệch chuẩn của mẫu là  $s = \sqrt{s^2} = \sqrt{5,70} = 2,39$ . Tập dữ liệu càng biến thiên, thì giá trị của  $s$  sẽ càng lớn

Đối với tập hợp nhỏ của các giá trị đo lường chúng ta đã sử dụng, thì việc tính toán phương sai không quá khó. Tuy nhiên, đối với một tập hợp lớn hơn, những tính toán có thể trở nên rất nhàm chán. Hầu hết máy tính cầm tay có khả năng thống kê đều có các chương trình cài sẵn mà sẽ tính  $\bar{x}$  và  $s$  hay  $\mu$  và  $\sigma$ , do đó công việc tính toán của anh/chị sẽ được giảm đến mức thấp nhất. Phím trung bình của mẫu hay tổng thể thường được đánh dấu bằng chữ  $\bar{x}$ . Phím độ lệch chuẩn của mẫu thường được đánh dấu bằng chữ  $s$  hay  $\sigma_{n-1}$ , và phím độ lệch chuẩn của tổng thể thường được đánh dấu bằng chữ  $\sigma$  hay  $\sigma_N$ . Khi sử dụng bất kỳ máy tính cầm tay nào có những phím chức năng cài sẵn này, hãy nắm chắc rằng anh/chị biết phép tính toán nào đang được thực hiện bởi mỗi phím!

Nếu anh/chị cần tính  $s^2$  và  $s$  bằng tay, thì sẽ dễ dàng hơn nhiều nếu sử dụng công thức tính thay thế được cho dưới đây. Hình thức tính toán này đôi khi được gọi là phương pháp đi tắt để tính toán  $s^2$ .

**Công thức tính toán đối với  $s^2$**

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

trong đó

$$\sum_{i=1}^n x_i^2 = \text{tổng các bình phương của những giá trị quan sát đơn lẻ}$$

$$\left(\sum_{i=1}^n x_i\right)^2 = \text{bình phương của tổng các giá trị quan sát đơn lẻ}$$

*Độ lệch chuẩn của mẫu,  $s$ , là căn bậc hai dương của  $s^2$ .*

**THÍ DỤ 2.7** Hãy tính phương sai và độ lệch chuẩn cho năm giá trị đo lường trong Bảng 2.4 mà được cho trước là 5, 7, 1, 2, và 4. Hãy sử dụng công thức tính toán đối với  $s^2$  và so sánh các kết quả của anh/chị với các kết quả thu được bằng cách sử dụng định nghĩa nguyên thủy của  $s^2$ .

**BẢNG 2.5**  
Bảng dành để tính toán  $s^2$  và  $s$  theo cách đã đơn giản hóa

$x_i$	$x_i^2$
5	25
7	49
1	1
2	4
4	16
19	95

**Lời giải** Những số ghi trong Bảng 2.5 là các giá trị đo lường đơn lẻ,  $x_i$ , và bình phương của chúng,  $x_i^2$ , cùng với tổng của chúng. Bằng việc sử dụng công thức tính toán đối với  $s^2$ , chúng ta có

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

$$= \frac{95 - \frac{(19)^2}{5}}{4} = \frac{22,80}{4} = 5,70$$

và  $s = \sqrt{s^2} = \sqrt{5,70} = 2,39$ , như đã tính trước đây •

**THÍ DỤ 2.8** Hãy tính phương sai mẫu và độ lệch chuẩn cho  $n = 25$  lợi suất trong Bảng 2.3

**BẢNG 2.3**

Lợi suất cổ tức	3,1	4,2	2,3	3,3	2,8
(%) đối với	5,3	3,5	3,1	2,6	3,3
25 cổ phiếu	4,7	3,7	3,0	2,6	4,0
thường của	3,8	4,4	3,2	3,2	3,8
ngân hàng	5,1	3,7	2,3	4,3	3,9

**Lời giải** Bằng việc sử dụng một máy tính cầm tay có các chức năng thống kê cài sẵn, anh/chị có thể kiểm tra những kết quả sau đây:

$$\sum_{i=1}^n x_i = 89,2$$

$$\sum_{i=1}^n x_i^2 = 333,82$$

Sử dụng công thức tính toán

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

$$= \frac{333,82 - \frac{(89,2)^2}{25}}{24} = \frac{15,5544}{24} = 0,6481$$

và  $s = \sqrt{s^2} = \sqrt{0,6481} = 0,81$  •

Anh/Chị có thể tự hỏi tại sao chúng ta chia cho  $n - 1$  chứ không phải  $n$ , khi chúng ta tính toán phương sai mẫu. Trung bình mẫu  $\bar{x}$  được sử dụng như là một hàm ước lượng của trung bình tổng thể, bởi vì nó cung cấp một giá trị ước lượng tốt về  $\mu$ . Nếu chúng ta muốn sử dụng phương sai mẫu như là hàm ước lượng của phương sai tổng thể  $\sigma^2$ , thì phương sai mẫu  $s^2$  với  $n = 1$  ở mẫu số sẽ cho ra những giá trị ước lượng về  $\sigma^2$  tốt hơn so với một hàm ước lượng được tính với  $n$  ở mẫu số. **Vì lý do này, chúng ta sẽ luôn luôn chia cho  $n - 1$  khi tính toán phương sai mẫu  $s^2$  và độ lệch chuẩn của mẫu  $s$ .**

Vào lúc này, anh/chị đã biết cách thức tính toán phương sai và độ lệch chuẩn của một tập hợp các giá trị đo lường. Hãy nhớ những điểm sau đây:

- Giá trị của  $s^2$  hay  $s$  càng lớn, thì độ biến thiên của tập dữ liệu càng lớn
- Nếu  $s^2$  hay  $s$  bằng số không, thì tất cả các giá trị đo lường phải có cùng giá trị
- Độ lệch chuẩn  $s$  được tính toán để có một thước đo về độ biến thiên mà được đo lường bằng cùng đơn vị như các giá trị quan sát.

Thông tin này cho phép chúng ta so sánh vài tập dữ liệu xét theo vị trí và độ biến thiên của chúng. Chúng ta có thể sử dụng những thước đo này như thế nào để nói điều gì đó cụ thể hơn về một tập dữ liệu duy nhất? Định lý và quy tắc được trình bày trong phần sau sẽ giúp chúng ta trả lời câu hỏi này.

## 2.6 Các Thước đo về Vị trí Tương đối (Measures of Relative Standing)

Đôi khi chúng ta muốn biết vị trí của một giá trị quan sát so với những giá trị quan sát khác trong một tập dữ liệu. Thí dụ, nếu anh/chị dự một kỳ thi tìm việc làm và đạt số điểm là 640, anh/chị có thể muốn biết tỷ lệ phần trăm những người tham dự đạt số điểm thấp hơn 640. Một **thước đo về vị trí tương đối** như thế của một giá trị quan sát trong một tập dữ liệu được gọi là **phân vị**.

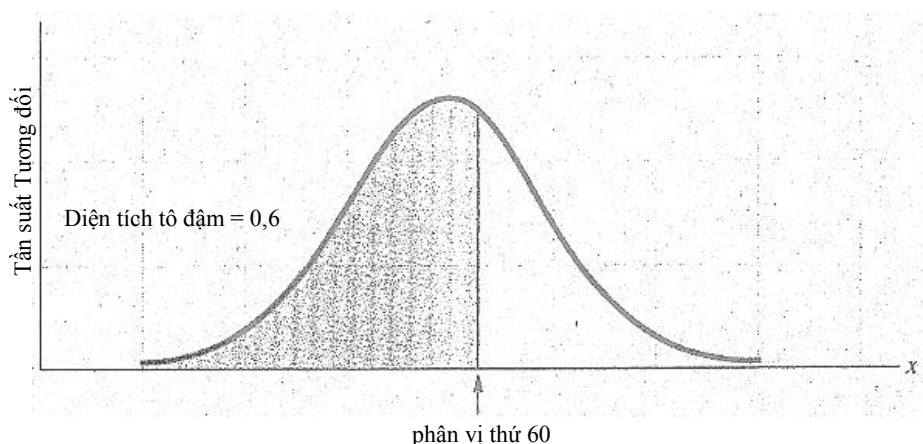
- ĐỊNH NGHĨA** • Cho  $x_1, x_2, \dots, x_n$  là một tập hợp  $n$  giá trị đo lường được sắp xếp theo thứ tự tăng dần. **Phân vị** thứ  $p$  là giá trị của  $x$  sao cho nhiều nhất là  $p$  phần trăm các giá trị đo lường là thấp hơn giá trị đó của  $x$  và nhiều nhất là  $(100-p)$  phần trăm là lớn hơn. •

**THÍ DỤ 2.9** Trước khi được nhận vào học một chương trình thạc sĩ quản trị kinh doanh (MBA) tại một trường đại học, anh/chị đã được thông báo rằng số điểm của anh/chị là 610 trong Kỳ Kiểm tra Miệng về Thành tích của Người Tốt nghiệp Đại học đã đặt anh/chị tại phân vị thứ 60 trong phân phối của những số điểm. Số điểm 610 của anh/chị đứng ở đâu so với những số điểm của những người khác cùng dự kỳ thi kiểm tra với anh/chị?

**Lời giải** Đạt số điểm tại phân vị thứ 60 có nghĩa là 60% những số điểm kiểm tra khác là thấp hơn số điểm của anh/chị và 40% là cao hơn.

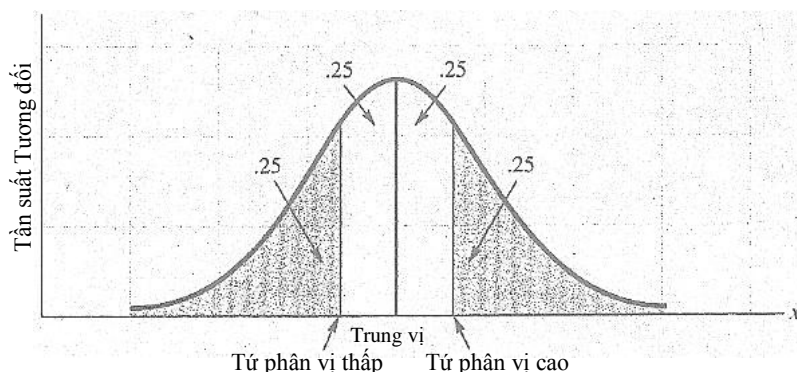
Xem xét theo đồ thị thì một phân vị nào đó, ví dụ phân vị thứ 60, là một điểm trên trục hoành  $x$  và nằm ở vị trí sao cho 60% diện tích bên dưới biểu đồ tần suất tương đối của dữ liệu nằm bên trái phân vị thứ 60 (xem Hình 2.7) và 40% diện tích này nằm bên phải. Như thế, theo định nghĩa, trung vị của một tập dữ liệu là phân vị thứ 50, bởi vì một nửa của các giá trị đo lường trong tập dữ liệu là nhỏ hơn số trung vị và một nửa là lớn hơn.

**HÌNH 2.7**  
Phân vị thứ 60 được trình bày trên biểu đồ tần suất tương đối của một tập dữ liệu



Phân vị thứ 25 và phân vị thứ 75, được gọi là **tứ phân vị thấp** và **tứ phân vị cao (lower and upper quartiles)**, cùng với trung vị (phân vị thứ 50), xác định vị trí những điểm mà chia dữ liệu thành bốn tập hợp có cỡ bằng nhau. Hai mươi lăm phần trăm các giá trị đo lường sẽ thấp hơn tứ phân vị thấp (đầu tiên), 50% sẽ thấp hơn trung vị (tứ phân vị thứ hai), và 75% các giá trị đo lường sẽ thấp hơn tứ phân vị cao (thứ ba). Như thế, trung vị và các tứ phân vị cao và thấp nằm tại những điểm trên trục  $x$  sao cho diện tích bên dưới biểu đồ tần suất tương đối của dữ liệu được phân chia thành bốn diện tích bằng nhau, như được cho thấy trong Hình 2.8. Anh/Chị có thể nhận thấy (trong Hình 2.8) rằng  $\frac{1}{4}$  diện tích bên dưới biểu đồ nằm bên trái của tứ phân vị thấp và  $\frac{3}{4}$  nằm bên phải. Tứ phân vị cao là giá trị của  $x$  sao cho  $\frac{3}{4}$  của diện tích này nằm bên trái và  $\frac{1}{4}$  nằm bên phải.

**HÌNH 2.8**  
Vị trí của các tứ phân vị



Còn có **Giá trị  $z$  (z-score)** là một thước đo khác về vị trí tương đối; nó sử dụng cả trung bình và độ lệch chuẩn của tập dữ liệu.

**ĐỊNH NGHĨA** • **Giá trị  $z$  của mẫu** tương ứng với một giá trị quan sát  $x$  là một thước đo về vị trí tương đối và được định nghĩa bằng công thức

$$\text{giá trị } z = \frac{x - \bar{x}}{s} \quad \bullet$$

Một giá trị  $z$  đo lường số lượng độ lệch chuẩn giữa một giá trị quan sát và trung bình của tập dữ liệu. Giả sử chúng ta biết rằng trung bình và độ lệch chuẩn của một tập hợp các số điểm kiểm tra, dựa trên một tổng số là 100 điểm, là  $\bar{x} = 74$  và  $s = 8$ . Giá trị  $z$  đối với điểm kiểm tra 92 của anh/chị được tính là

$$\text{giá trị } z = \frac{x - \bar{x}}{s} = \frac{92 - 74}{8} = 2,25$$

Vì thế số điểm của anh/chị nằm cao hơn 2,25 độ lệch chuẩn so với trung bình; đó là,  $92 = 74 + 2,25(8)$ .

Bản thân các giá trị  $z$  chỉ đơn thuần cho thấy số điểm kiểm tra cao hơn hay thấp hơn trung bình bao nhiêu độ lệch chuẩn. Tuy nhiên, khi giá trị  $z$  được sử dụng cùng với Định lý Tchebysheff, thì có thể đưa ra một số lời phát biểu thận trọng về vị trí tương đối của một giá trị quan sát. Hơn nữa, nếu dữ liệu có hình dạng cái gò, thì Quy tắc Thực nghiệm có thể được dùng để đưa ra những lời phát biểu mạnh hơn về vị trí tương đối của một giá trị quan sát xét theo giá trị  $z$  của nó. Bởi vì ít nhất là 75%, và rất có thể là 95%, các giá trị quan sát trong một tập dữ liệu nằm trong phạm vi hai độ lệch chuẩn so với trung bình, nên các giá trị  $z$  trong khoảng từ  $-2$  đến  $+2$  là rất có khả năng xảy ra, và như thế không phải là không bình thường. Tuy nhiên, ít nhất là 8/9, hay rất có thể là tất cả, các giá trị quan sát nằm trong phạm vi ba độ lệch chuẩn so với trung bình. Vì thế, các giá trị  $z$  trong khoảng từ 2 đến 3, tính theo giá trị tuyệt đối, ít có khả năng xảy ra hơn nhiều, trong khi đó các giá trị  $z$  cao hơn 3, tính theo giá trị tuyệt đối, rất không có khả năng xảy ra và phải được xem xét cẩn thận. Một điểm kiểm tra có giá trị  $z$  cao hơn 3 là xuất sắc, trong khi đó một cổ phiếu mà tỷ số giá trên thu nhập của nó (giá của cổ phiếu chia cho thu nhập bình quân mỗi cổ phiếu hàng năm) có giá trị  $z$  là  $-3$  sẽ được xem là một cuộc đầu tư có tiềm năng thu nhập tốt.

Giá trị  $z$  cực kỳ lớn và giá trị  $z$  cực kỳ nhỏ nêu lên câu hỏi về hiệu lực (validity) của một giá trị quan sát. Có thể giá trị quan sát này chỉ là hết sức lớn hoặc hết sức nhỏ so với những giá trị quan sát khác. Tuy nhiên, giá trị quan sát này có thể đã được ghi nhận không đúng, hoặc vì lý do nào đó, nó có thể không thuộc về tổng thể mà chúng ta đã mong muốn lấy mẫu. Những giá trị quan sát với các giá trị  $z$  hết sức lớn hoặc nhỏ thường được gọi là **giá trị dị biệt** bởi vì chúng nằm cách xa trung tâm của tập dữ liệu. Những giá trị quan sát nằm cao hơn hay thấp hơn trung bình trong khoảng từ hai đến ba độ lệch chuẩn là những giá trị dị biệt có thể có, trong khi đó những giá trị quan sát nằm cao hơn hay thấp hơn trung bình nhiều hơn ba độ lệch chuẩn thì được xem là những giá trị dị biệt rõ ràng.

**THÍ DỤ 2.10** Hãy xét một mẫu gồm  $n = 10$  giá trị đo lường:

3, 2, 0, 15, 2, 3, 4, 0, 1, 3



Thoạt nhìn anh/chị có thể thấy giá trị đo lường  $x = 15$  dường như là một giá trị dị biệt. Hãy tính giá trị  $z$  cho giá trị quan sát này, và hãy trình bày các kết luận của anh/chị.

**Lời giải** Đối với mẫu này, chúng ta có những phép tính toán sau đây:

$$\sum_{i=1}^{10} x_i = 33 \quad \text{và} \quad \sum_{i=1}^{10} x_i^2 = 277$$

Như thế

$$\begin{aligned} \bar{x} &= \frac{\sum_{i=1}^{10} x_i}{10} = \frac{33}{10} = 3,3 \\ s^2 &= \frac{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}}{n-1} \\ &= \frac{277 - \frac{(33)^2}{10}}{9} = \frac{168,1}{9} = 18,6778 \\ s &= 4,32 \end{aligned}$$

Bằng việc sử dụng những số lượng này để tính giá trị  $z$  cho giá trị dị biệt nghi ngờ  $x = 15$ , chúng ta tìm ra

$$\text{giá trị } z = \frac{x - \bar{x}}{s} = \frac{15 - 3,3}{4,32} = 2,71$$

Như thế giá trị đo lường  $x = 15$  nằm cách trung bình mẫu  $\bar{x} = 3,3$  một khoảng cách là 2,71 độ lệch chuẩn mẫu. Bởi vì giá trị  $z$  này cao hơn 2, nên chúng ta xác định  $x = 15$  là một giá trị dị biệt có thể có. Chúng ta phải xem xét thủ tục lấy mẫu của mình để xem liệu có bằng chứng cho thấy  $x = 15$  là một giá trị quan sát bị sai hay không •

Anh/Chị có thể sử dụng Minitab hay Excel để tạo ra nhiều trong số những thước đo mô tả bằng số mà chúng ta đã thảo luận. Trong Minitab, hãy dùng **Stat** → **Basic Statistics** → **Display Descriptive Statistics**, và chọn các biến thích hợp để mô tả. Trong Excel, hãy dùng **Tool** → **Data Analysis** → **Descriptive Statistics**, chọn dãy những ô có chứa dữ liệu, nhấn vào hộp có đánh dấu “Số liệu Thống kê Tổng hợp” (“Summary Statistics”), và chọn một ô trong đó thể hiện số liệu ra. Anh/Chị sẽ nhận ra trung bình, độ lệch chuẩn, trung vị, và những giá trị quan sát tối thiểu và tối đa trong cả hai bản in kết quả ra. Excel cũng tính số yếu vị, phương sai, và khoảng biến thiên (miền), trong khi đó Minitab bao gồm cả các tứ phân vị thấp và cao. Cả hai bản in kết quả ra đều có một số trị thống kê khác mà chúng ta chưa thảo luận.

Bản in kết quả ra của Minitab, được trình bày trong Hình 2.9, tổng hợp (summarize) các lợi suất cổ tức của Thí dụ 2.8 (dữ liệu được cho trong Bảng 2.3) và những giá trị quan sát trong Thí dụ 2.10. Anh/Chị có thể so sánh các giá trị của những trị thống kê được tính trong các thí dụ đó với các giá trị được trình bày trong bản in kết quả ra.

**HÌNH 2.9**

Bản in kết quả ra của Minitab sử dụng lệnh DESCRIBE (MÔ TẢ) cho dữ liệu trong Thí dụ 2.8 (C1) và dữ liệu của Thí dụ 2.10 (C2).		N	MEAN	MEDIAN	TRMEAN	STDEV	SEMEAN
C1		25	3,568	3,500	3,548	0,805	0,161
C2		10	3,30	2,50	2,25	4,32	1,37
		MIN	MAX	Q1	Q3		
C1		2,300	5,300	3,050	4,100		
C2		0,00	15,00	0,75	3,25		

## BÀI TẬP

### Các Kỹ thuật Căn bản

**2.10** Hãy sử dụng tập dữ liệu sau đây:

3, 9, 6, 5, 5, 4, 7, 6, 8, 2, 6, 7, 3

- a Hãy tính  $\bar{x}$  và  $s$
- b Hãy tính giá trị  $z$  cho những giá trị quan sát nhỏ nhất và lớn nhất. Có giá trị nào trong hai giá trị quan sát này lớn hay nhỏ bất thường không?

**2.11** Hãy tìm giá trị  $z$  cho giá trị quan sát lớn nhất trong tập dữ liệu sau đây:

19, 12, 16, 0, 14, 9, 6, 1, 12, 13, 10, 19, 7, 5, 8

**2.12** Nếu Anh/Chị đạt số điểm trong phân vị thứ 90 trong kỳ kiểm tra tuyển sinh vào lớp cao học, số điểm của anh/chị đứng ở vị trí ra sao so với những người khác cùng dự kỳ kiểm tra?

### Ứng dụng

**2.13** Tham khảo dữ liệu về số nợ bình quân đầu người trong Bài tập 2.8.

- a Hãy tìm trung bình và độ lệch chuẩn của các số nợ bình quân đầu người này.
- b Tham khảo bài tập 2.8 để tìm số nợ bình quân đầu người ở bang của anh/chị trong năm 1992. Hãy sử dụng một giá trị  $z$  để mô tả số nợ bình quân đầu người ở bang của anh/chị so sánh như thế nào với những số nợ tương ứng ở những bang khác.

**2.14** Một bài báo trong Tạp chí *American Demographics* (*Nhân Khẩu học Hoa Kỳ*) (Kirchner, R., và Thomas, R., “Những Thị trường Mới đối với Bảo hiểm Y tế,” tháng 12, 1990, trang 40) cung cấp một số số liệu thống kê thú vị về số dân Mỹ không có loại bảo hiểm y tế nào trong năm 1988. Trong nhiều trường hợp, những người Mỹ không được bảo hiểm này ít nhất cũng có đủ nguồn lực để trả cho bảo hiểm y tế. Theo bài báo này, “gần 40 phần trăm

số người không được bảo hiểm có thu nhập là 20.000US hay nhiều hơn; 22 phần trăm có thu nhập là 30.000USD hay nhiều hơn; và 13 phần trăm, hay trên 4 triệu, sống trong những hộ gia đình có thu nhập là 40.000USD hay nhiều hơn.” Hãy nhận dạng những phân vị nào có thể được xác định từ thông tin này.

- 2.15** Theo *Consumer Reports (Báo cáo Người Tiêu dùng)* (Tháng 3/1994), giá trung bình của một Sony SLV-700HF stereo VCR là 410USD, với độ lệch chuẩn là 14USD. Nếu anh/chị mua loại VCR này với giá 430USD, hãy tính giá trị  $z$  đối với giá mua của anh/chị. Giá này có cao bất thường không?

## 2.7 Tóm tắt

Những phương pháp mô tả tập hợp các giá trị đo lường có thể chia thành hai loại, đó là phương pháp bằng đồ thị và phương pháp bằng số. Biểu đồ tần suất tương đối là một phương pháp bằng đồ thị cực kỳ hữu ích để biểu thị đặc trưng một tập hợp các giá trị đo lường. Các thước đo mô tả bằng số là các con số mà cố gắng tạo ra một hình ảnh trong trí óc về biểu đồ tần suất (hay phân phối tần suất). Chúng ta đã hạn chế nội dung thảo luận trong các thước đo hướng tâm và sự biến thiên, mà hữu ích nhất trong các thước đo này là trung bình và độ lệch chuẩn. Mặc dù trung bình có ý nghĩa mô tả theo trục giá, nhưng độ lệch chuẩn chỉ có ý nghĩa khi được sử dụng cùng với Định lý Tchebysheff và Quy tắc Thực nghiệm. Mục tiêu của việc lấy mẫu là mô tả (đưa ra những suy luận về) về tổng thể từ đó mẫu này đã được lấy ra. Mục tiêu này được hoàn thành bằng việc sử dụng trung bình mẫu  $\bar{x}$  và số lượng  $s^2$  như là các hàm ước lượng về trung bình tổng thể  $\mu$  và phương sai  $\sigma^2$ . Khi dữ liệu gồm có những cặp giá trị quan sát, thì đồ thị nhị biến được dùng để đánh giá bằng hình ảnh cách thức  $x$  thay đổi theo  $y$ , trong khi đó hệ số tương quan được dùng để xác định sức mạnh của mối quan hệ tuyến tính giữa  $x$  và  $y$ . Các thước đo khác, chẳng hạn như phân vị hay giá trị  $z$ , được dùng để xác định vị trí tương đối của một quan sát trong tổng thể hay trong một mẫu. Các đồ thị hộp là những tóm lược dữ liệu bằng hình ảnh và chúng hữu ích trong việc phát hiện các giá trị dị biệt.

Các phương pháp mô tả và các thước đo bằng số đã được trình bày chỉ là một số nhỏ trong những phương pháp và thước đo lẽ ra có thể được thảo luận. Ngoài ra, ở đây cũng đã bỏ qua nhiều kỹ thuật tính toán đặc biệt thường được tìm thấy trong những cuốn sách giáo khoa sơ cấp. Sự bỏ qua này bắt buộc phải có do thời gian hạn chế trong bất kỳ khóa học sơ cấp nào. Ngoài ra, việc sử dụng phổ biến máy tính cầm tay và máy tính đã tối thiểu hóa tầm quan trọng của những công thức tính toán đặc biệt. Nhưng quan trọng hơn, việc bao gồm vào những kỹ thuật như thế sẽ thường hay làm lu mờ và gây khó hiểu cho mục tiêu chính của thống kê hiện đại và cuốn sách này: đó là suy luận thống kê.