# STATISTICAL PROCEDURES

## AND THEIR

# MATHEMATICAL BASES

BY

CHARLES C. PETERS, Ph.D.

*Director of Educational Research,*
*The Pennsylvania State College*

AND

WALTER R. VAN VOORHIS, M.A.

*Schuylkill Administrative Head and*
*Assistant Professor of Mathematics,*
*The Pennsylvania State College*

First Edition
Second Impression

# STATISTICAL PROCEDURES

AND THEIR

# MATHEMATICAL BASES

# PREFACE

This volume is a revision and extension of a book by the same title published privately in lithoprinted form in 1935. The wide demand for the preliminary edition showed that there was a need for the type of presentation of statistics here offered.

The characteristic feature of the book is the effort to explain the mathematical origins of the most widely used statistical formulas in terms that persons with comparatively little mathematical training can easily follow. We believe that, if statistical workers do not take their tools as magic but understand them in the light of their origins and assumptions, they will use these tools more intelligently and more safely. In order to make such understanding available to persons of little mathematical training we give the derivations in much detail. It is a well-known fact that the source of difficulty in mathematical reading by relatively untrained persons is largely the omission of steps which are supposed to be obvious. When these steps are supplied and when the use of specialized mathematical terminology is reduced to a minimum, much that would otherwise be closed to the reader is readily understandable. In order to make calculus available as a tool for those who do not have a command of it, we open this volume with a chapter on calculus. This is, of course, only "a little calculus," but it is enough to prepare the reader who has not hitherto studied calculus to follow the derivations in which we must draw upon this branch of mathematics. Our experience with this presentation, as well as that reported by some others, shows that this chapter on calculus can be mastered in about 10 per cent of the time normally allotted to a one-semester course in advanced statistics.

The title of the book is somewhat too pretentious. It might better be called *Some Statistical Procedures and a Little Insight into the Mathematical Bases of a Few of Them.* It is not, of course, a comprehensive treatment of the mathematical bases of statistics. It is intended to bridge the gap between the elementary courses, in which the formulas are given purely authoritatively,

# CONTENTS

# STATISTICAL PROCEDURES AND THEIR MATHEMATICAL BASES

## CHAPTER I

### A LITTLE CALCULUS

This chapter is intended for persons who have not previously studied calculus. It presents, in a way that a reader who has had only a limited training in mathematics should be able to follow, practically all the calculus upon which we shall have occasion to draw in this volume on statistics, which includes many fundamental elements common also to applications in other fields. We trust that this simple presentation of the elements of differential and integral calculus may not only prove useful to the student of statistics but that it may also give to laymen in mathematics an interesting and culturally enriching insight into the nature and applications of this fascinating mathematical discipline.

In every case where one quantity varies in a manner that is definitely related to the variation in a second, the relation between the two may be represented geometrically by a curve of some shape (including a straight line). Take first the simple relation $y = x$, where $x$ may be represented by any number and $y$ will therefore of necessity be the same number. We may lay off this relation on the adjacent diagram.



Fig. 1.—Straight line relation: Slope = 1.

When $x = 1$, $y = 1$; when $x = 2$, $y = 2$, etc. If we go to the right one unit for $x$ and then up one unit for $y$, we shall have a point $x_1y_1$ that shows the relation between the two series at that value of $x$. If we go two units to the right and two units up, we shall have a second point $x_2y_2$; etc.

1

A straight line may be drawn through all these points. Its slope will be $\frac{1}{1}$, $\frac{2}{2}$, $\frac{3}{3}$, $\frac{4}{4}$, $\cdots = 1$. Its slope will always be the same at all values of $x$.

Suppose, now, that $y = 0.4x$. We shall then have a line representing the relation as follows:



| | |
|---|---|
| $x = 0,$ | $y = 0$ |
| $x = 1,$ | $y = 0.4$ |
| $x = 2,$ | $y = 0.8$ |
| $x = 3,$ | $y = 1.2$ |
| $x = 4,$ | $y = 1.6$ |
| $x = 5,$ | $y = 2.0$ |
| $\cdots\cdots\cdots\cdots\cdots$ | |

Fig. 2.—Straight line relation: Slope = 0.4.

Here likewise we can represent the relation by a straight line, and this line will have the same slope at every value of $x$; at each point a change of $n$ units in $x$ will be accompanied by a change of $0.4n$ units in $y$.

But let us now take a more complicated case, $y = x^2$.

| | |
|---|---|
| $x = 0,$ | $y = 0$ |
| $x = \frac{1}{2},$ | $y = \frac{1}{4}$ |
| $x = 1,$ | $y = 1$ |
| $x = 1\frac{1}{2},$ | $y = 2\frac{1}{4}$ |
| $x = 2,$ | $y = 4$ |
| $x = 3,$ | $y = 9$ |
| $x = -\frac{1}{2},$ | $y = \frac{1}{4}$ |
| $x = -1,$ | $y = 1$ |
| $x = -1\frac{1}{2},$ | $y = 2\frac{1}{4}$ |
| $x = -2,$ | $y = 4$ |
| $x = -3,$ | $y = 9$ |
| $x = -4,$ | $y = 16$ |
| $\cdots\cdots\cdots\cdots\cdots$ | |



Fig. 3.—Curved line relation: Slope changes with $x$.

Here the line is not a straight one; it does not have the same slope at all values of $x$. As we proceed out from the $y$ axis, the slope is at first very small; at $x = 1$ it is moderate; and at $x = 3$ the slope is very steep. We have a similar behavior on the side where $x$ has negative values. We have, in fact, very great difficulty in saying what the slope is, because it is always changing. We could draw a straight line between $A$ and $B$, where $x = 1$ and $x = 2$, respectively, but the slope of this line would

not precisely describe the slope of the curve. If we took a smaller change in $x$, say that represented by the distance $AB'$ on our scale, our secant line would more nearly coincide with the curve. We may consider $A$ as any fixed point on the graph and allow $B$ to move along the curve and approach $A$ as a limiting position The changes in $x$ would become smaller and smaller and approach zero as a limit. The secant line drawn through $A$ and $B$ would turn about $A$, approaching the tangent line at $A$ as its limiting position.

Now the basic task of the differential calculus (except in regions of discontinuity and other similar matters which lie beyond the scope of this chapter) is to ascertain the slope of a curve at various points by determining the slope of a secant which, in a limiting position, becomes the tangent to the curve at the point in question. This same idea may be expressed in other terms by saying that it is the task of the differential calculus to ascertain the amount of change in a variable $y$ that corresponds to a certain change in a related variable $x$ as these increments in the independent variable $x$ become so small as to approach zero in value. At certain times in its history this discipline has been called the *infinitesimal calculus* in recognition of the fact that it deals with the relation of infinitesimal increments of one variable to infinitesimal increments of another. In operating with the calculus we are often operating algebraically with no curve in sight; but usually we can represent these algebraic operations geometrically and show that what we are seeking is something about the slope of that curve at some point.

### DIFFERENTIATION

Let us proceed with that algebraic process with which we said, in our preceding paragraph, we shall often be operating with no curve before us visually. We have the equation

$$y = x^2$$

We wish to find what change in $y$ goes with a change in $x$ at any value of $x$ in which we are interested. Let $\Delta x$ be an increment to be added to $x$ (algebraically), and $\Delta y$ be the corresponding increment that would need to be added to $y$ in order to maintain the equation.

(1) $$y + \Delta y = (x + \Delta x)^2$$

Performing the indicated square,

$$y + \Delta y = x^2 + 2x\,\Delta x + \overline{\Delta x}^2$$

But our original equation gave us $y = x^2$. We may subtract the terms of this equation from the corresponding ones of our last equation above on the basis of the axiom, "If equals be subtracted from equals the remainders are equal." We shall then have

(2) $$\Delta y = 2x\,\Delta x + \overline{\Delta x}^2$$

Dividing through by $\Delta x$, we shall have

(3) $$\frac{\Delta y}{\Delta x} = 2x + \Delta x$$

We said $\Delta x$ should be an increment added to $x$ and $\Delta y$ an increment added to $y$, but we did not commit ourselves as to the particular size of the increment. Let us now conceive of $\Delta x$ as decreasing until it becomes infinitesimal in size. It will necessarily drag $\Delta y$ down with it, since the equation must continue to hold for all values of $\Delta x$. When $\Delta x$ has become so small as to have approached zero as its limit let us replace $\Delta y/\Delta x$ by $dy/dx$. At this limit the $\Delta x$ in the last term of our equation will approach zero in value and thus disappear from consideration. The reason why $dy/dx$ can not be similarly dropped as of zero value is that both its numerator and its denominator become small together so that the fraction has a value that may be of considerable dimension. And so as the limit zero is approached by $\Delta x$ we have

(4) $$\frac{dy}{dx} = 2x$$

This $2x$ is called the *derivative* of the expression $y = x^2$. The process of getting it is called *differentiation*. If the $dx$ appears in the denominator of the fraction expressing the derivative, we say that we are differentiating "with respect to $x$"; if the $dy$ appears in the denominator, as it will sometimes do in our later developments, we say we are differentiating "with respect to $y$." This process alone, in more or less complicated forms, constitutes essentially all there is to the differential calculus. In terms of the slope of a curve a derivative equal to $2x$ means that, at the

point where $x = 1$, the slope of the curve is 2 times 1 or 2 (which means that at that point the $y$ values are changing twice as rapidly as the $x$ values through the infinitesimal distance to which we have shrunken our $\Delta x$ at its limit). At the point where $x = 3$, the slope of the curve relating the two is 2 times 3 or 6, which means that $y$ changes 6 units for each unit of change in $x$.

After a few more concrete examples we shall seek a general rule for differentiating an expression directly without going each time through a long process of algebraic manipulation. But in the meantime the reader may be interested in observing the relation of the form of the $2x$ to the $x^2$ of which it is the derivative. He will notice that the exponent of the $x$ has dropped from 2 to 1, a decrease of one unit. He will also observe that the coefficient of the derivative has become 2, possibly the same 2 that was lost from the original exponent; about that we shall see later.

Let us now try differentiating the expression $y = x^3$. We shall go through the same four fundamental steps, through which we passed in our previous example, as follows: (1) Add $\Delta y$ to the $y$ and $\Delta x$ to the $x$ and perform the indicated involution. (2) Subtract from the resultant equation our original equation. (3) Divide through by $\Delta x$. (4) Let $\Delta x$ approach zero as a limit, and, as the limit is approached, substitute $dy/dx$, the symbol for the derivative at the limit, for $\Delta y/\Delta x$; drop from the equation any of these $\Delta x$ values that stand without a $\Delta$ denominator, on the ground that even in the first power their values are approximately zero and that in any of the higher powers the values are lower than in the first power.

$$y = x^3$$

Adding $\Delta y$ to $y$ and $\Delta x$ to $x$,

(1) $$y + \Delta y = (x + \Delta x)^3$$

Expanding the second term,

$$y + \Delta y = x^3 + 3x^2\Delta x + 3x\,\overline{\Delta x}^2 + \overline{\Delta x}^3$$

Subtracting,

(2) $$\Delta y = 3x^2\Delta x + 3x\,\overline{\Delta x}^2 + \overline{\Delta x}^3$$

Dividing by $\Delta x$,

(3) $$\frac{\Delta y}{\Delta x} = 3x^2 + 3x\,\Delta x + \overline{\Delta x}^2$$

Letting $\Delta x$ approach zero,

$$(4) \qquad \frac{dy}{dx} = 3x^2$$

In this last expression the $3x\,\Delta x$ dropped out in the limit because as $\Delta x$ approaches zero as a limit, any product formed by multiplying it by any factor approaches zero and, in the limit, disappears from the equation. For a similar reason the $\overline{\Delta x}^2$ becomes zero in the limit. In fact since $\Delta x$ becomes, as it decreases, a very small quantity (*i.e.*, a decimal quantity less than 1), when raised to any power (including 1) and multiplied by 1 or by any other factor it will approach zero and vanish from the equation as $\Delta x$ approaches zero as its limit.

The derivative of $y = x^3$ is, therefore, $3x^2$. Notice that here, again, the original exponent has become the coefficient of the derivative and that the exponent of the $x$ in the derivative is one less than that of the original quantity. Let us now take a more generalized example,

$$y = x^n$$

where $n$ may represent any positive integer.[1] Performing in succession our four fundamental steps,

$$(1) \qquad y + \Delta y = (x + \Delta x)^n$$

Expanding,

$$y + \Delta y = x^n + nx^{n-1}\Delta x + \frac{n(n-1)}{1 \cdot 2}\,x^{n-2}\overline{\Delta x}^2$$
$$+ \frac{n(n-1)(n-2)}{1 \cdot 2 \cdot 3}\,x^{n-3}\overline{\Delta x}^3 + \cdots$$
$$y = x^n$$

was our original equation, to be subtracted,

$$(2) \quad \Delta y = nx^{n-1}\Delta x + \frac{n(n-1)}{1 \cdot 2}\,x^{n-2}\overline{\Delta x}^2$$
$$+ \frac{n(n-1)(n-2)}{1 \cdot 2 \cdot 3}\,x^{n-3}\overline{\Delta x}^3 + \cdots$$

[1] It may be shown that the rule for differentiating functions of this form will hold for any real value of $n$.

$$(3) \quad \frac{\Delta y}{\Delta x} = nx^{n-1} + \frac{n(n-1)}{1 \cdot 2} x^{n-2}\Delta x$$
$$+ \frac{n(n-1)(n-2)}{1 \cdot 2 \cdot 3} x^{n-3}\overline{\Delta x}^2 + \cdots$$

Letting $\Delta x$ approach zero as limit and observing what was said above about the vanishing of all powers of $\Delta x$ that stand without a $\Delta$ denominator, we have remaining

$$\frac{dy}{dx} = nx^{n-1} \qquad \text{(Derivative of the function } y = x^n) \quad (1)$$

From this general case it is now obvious that what we inferred as a possibility in our two previous examples is in fact the rule: *the derivative in respect to x at any power has as its coefficient the original power of x and as its exponent the original exponent decreased by* 1. This must be so because only the second term in the binomial expansion is free from the $\Delta x$ after the subtraction of step (2) and the division of step (3) and because the coefficient of the second term in a binomial expansion is always the power to which the binomial is being raised and its exponent is always the original exponent less 1.

Suppose now we try the effect of a constant as coefficient of our variable $x$,

$$y = ax^n$$

Here *a* may represent any coefficient we please, whether integral or fraction, whether positive or negative. Going through our four steps,

$$(1) \qquad y + \Delta y = a(x + \Delta x)^n$$

$$y + \Delta y = ax^n + anx^{n-1}\Delta x + \frac{an(n-1)}{1 \cdot 2} x^{n-2}\overline{\Delta x}^2$$
$$+ \frac{an(n-1)(n-2)}{1 \cdot 2 \cdot 3} x^{n-3}\overline{\Delta x}^3 + \cdots$$

$$(2) \quad \Delta y = anx^{n-1}\Delta x + \frac{an(n-1)}{1 \cdot 2} x^{n-2}\overline{\Delta x}^2$$
$$+ \frac{an(n-1)(n-2)}{1 \cdot 2 \cdot 3} x^{n-3}\overline{\Delta x}^3 + \cdots$$

$$(3) \quad \frac{\Delta y}{\Delta x} = anx^{n-1} + \frac{an(n-1)}{1 \cdot 2} x^{n-2}\Delta x$$
$$+ \frac{an(n-1)(n-2)}{1 \cdot 2 \cdot 3} x^{n-3}\overline{\Delta x}^2 + \cdots$$

(4)    $\dfrac{dy}{dx} = anx^{n-1}$      (Derivative of a constant times   (2)
                                  a function of the form $x^n$)

Here the constant reappears in the derivative unchanged. Hence, since $a$ may represent any constant, we may say: *The derivative of a constant times a function is the same constant times the derivative of the function.*

Let us now take a more complicated expression, one involving $x$ in each of several terms with different powers in each term, the total function being the sum of these several functions.

$$y = ax^3 + bx^2 + cx$$

(1)   $y + \Delta y = ax^3 + 3ax^2\Delta x + 3ax\,\overline{\Delta x}^2 + a\,\overline{\Delta x}^3 + bx^2$
$$+ 2bx\,\Delta x + b\,\overline{\Delta x}^2 + cx + c\,\Delta x$$

$y \qquad\quad = ax^3 \qquad\qquad\qquad\qquad\qquad\qquad + bx^2 + cx$

(2)   $\Delta y = 3ax^2\Delta x + 3ax\,\overline{\Delta x}^2 + a\,\overline{\Delta x}^3 + 2bx\,\Delta x + b\,\overline{\Delta x}^2 + c\,\Delta x$

(3)    $\dfrac{\Delta y}{\Delta x} = 3ax^2 + 3ax\,\Delta x + a\,\overline{\Delta x}^2 + 2bx + b\,\Delta x + c$

(4)   $\dfrac{dy}{dx} = 3ax^2 + 2bx + c$      (Derivative of the sum of   (3)
                                    functions)

If the reader will compare this derivative with the expression we started out to differentiate, he will observe that the derivative of the complex quantity made up of the sum of three terms is precisely the sum of the derivatives of the several terms if differentiated separately. If he will carry through on paper the generalized case or visualize to himself how it would work out, he will easily convince himself that that same conclusion would hold universally for any values for which the binomial law holds. Therefore, *the derivative of the sum of any number of functions is the sum of their derivatives.*

Let us now try differentiating a constant, $y = a$. Since $x^0 = 1$, the above equation might be written

$$y = x^0 a$$

Going through our four steps,

(1)                  $y + \Delta y = (x + \Delta x)^0 a = a$

Subtracting our original equation, $y = a$, we get

(2) $\Delta y = 0.$    Then (3) $\dfrac{\Delta y}{\Delta x} = 0$; and (4) $\dfrac{dy}{dx} = 0$