# Measurement of genetical differentiation among subpopulations

**H.-R. Gregorius** [*, 1] **and J. H. Roberds** [2]

[1] Department of Genetics, [2] Southeastern Forest Experiment Station, USDA Forest Service; North Carolina State University, Raleigh, NC 27695-7614, USA

**Summary.** The basis for measuring differentiation among subpopulations is discussed and a number of conditions formulated that are desirable for an appropriate measure. These conditions imply that each subpopulation is characterized by the difference in genetic composition between it and its complement. A direct method of determining this difference is described and shown to result from a known measure of genetic distance between populations. The weighted average of the genetic distances between subpopulations and their complements constitutes a measure of differentiation among subpopulations which fulfills all of the desirable conditions and has the additional advantage that its values are directly interpretable. This measure ($\delta$) is equally applicable to gene ($\delta_{ge}$), gametic ($\delta_{ga}$) or genotypic ($\delta_{go}$) frequencies, which guarantees an unequivocal multilocus treatment, provided that the sets of genetic entities to which the frequencies refer are properly defined. The general relationship $\delta_{ge} \leq \delta_{ga} \leq \delta_{go}$ is consistent with the principle that increasing complexity of organization of genetic material results in increased opportunity for differentiation. It is demonstrated that Wright's $F_{ST}$ ($G_{ST}$ in Nei's notation), which is often used to measure subpopulation differentiation, meets some but not all of the conditions formulated for a desirable measure.

* Current address: Abteilung für Forstgenetik, Universität Göttingen, Büsgenweg 2; D-3400 Göttingen, Federal Republic of Germany

## Introduction

Experimental investigations of migration and reproductive behavior in natural populations have repeatedly used the concept of population as the basic unit for analysis. For this purpose a population is defined as a collection of organisms that forms a reproductive entity. Gene flow, however, might be restricted between portions of a collection of organisms so that sections do not function as a single reproductive unit but rather as a system of subpopulations which are more or less reproductively isolated from each other (Ehrlich and Raven 1969). Among the well-documented causes of reproductive isolation are geographic separation, local adaptations combined with restricted migration, and positive assortative mating due to mating preference in animals and asynchronous flowering in plants (Hamrick 1983). The degree to which populations are subdivided into reproductive isolates provides valuable information to breeders and conservationists concerned with selection of appropriate base populations and to evolutionary biologists studying sympatric or allopatric speciation.

From the standpoints of evolution, breeding, and gene conservation, an effective measure of differentiation among subpopulations should be based only on genetic differences. A measure based on phenotypic differences would indicate subpopulation differentiation for situations in which the differences are due entirely or partly to environment. Likewise a measure based on the degree of reproductive isolation would indicate erroneous differentiation in populations in which portions are reproductively isolated from each other but in which all subdivisions are identical genetically. While it is true that in many situations reproductive isolation ultimately leads to genetic differentiation

among subpopulations for some loci, even for these cases, the degree of reproductive isolation cannot be counted on to accurately reflect the amount of genetic differentiation. For example, in such situations some loci may show marked differentiation while others may be almost genetically uniform across subpopulations.

A measure often used to analyze population differentiation with gene frequency data involves one of Wright's $F$ statistics $(F_{ST})$. $F$ statistics were originally developed to measure degrees of inbreeding, assortative mating, and differentiation due to selection (Wright 1969) and were later reinterpreted in terms of differentiation among subpopulations (Wright 1978). Here we adopt a different approach for deriving an appropriate measure. We first determine a number of intuitively reasonable conditions that a measure of subpopulation differentiation should satisfy and then construct a measure that meets these conditions. Genetic units for which this measure can be applied are specified so that the measure can be interpreted without confusion. For multiple-locus data, specification of the genetic unit being considered is especially important to prevent ambiguity. Wright's $F_{ST}$ is also discussed in light of the specified conditions and compared with the derived measure using published data from a natural population.

## Conditions for a measure of subpopulation differentiation

When a population is divided into subpopulations, each method of partition results in a degree of differentiation among the component subpopulations. The amount of differentiation among subpopulations may, however, vary considerably among such partitions. A meaningful determination of the differentiation in the population can then be made only with respect to given subpopulation structure. Here genetic effects are not considered in defining a subpopulation structure since primary interest is focused on determination of the amount of genetic differentiation among subpopulations. Criteria for establishing such a structure may be based on spatial, ecological, behavioral, sociological or demographic factors.

It is desirable to summarize genetic differences among subpopulations in a single measure that properly reflects the differentiation among them. It is appropriate that such a measure have a value of zero in the absence of differentiation, i.e. when all subpopulations have identical genetic compositions for the loci under investigation. Furthermore, only in the case of complete differentiation, i.e. when all subpopulations are completely different genetically with respect to the set of loci being investigated, should the measure attain its maximum value. The maximum value must be finite

so that high and low differentiation can be meaningfully evaluated. For convenience it is desirable to have the measure fall within the interval [0, 1], so we set the maximum equal to 1. A major difficulty lies with the specification of the measure for intermediate values where differentiation among subpopulations occurs but is not complete. The resolution of this problem can be illustrated by the following example. Consider a population that is divided into three subpopulations: one containing 60% of the organisms in the population; another, 30%; and a third, 10%. Now assume that the two larger subpopulations have identical frequencies at a particular genetic level (allele, gamete, genotype) but that the smallest subpopulation differs completely from the other two, i.e. at the respective level it has no genetic components in common with either of the other two. Clearly the subpopulations with 30% of the organisms differs from its complement subpopulation, i.e. the remainder of the population. Furthermore this difference is smaller than the amounts by which the remaining subpopulations differ from their respective complements. The complement of the 30% subpopulation contains the 60% subpopulation which is genetically identical to the 30% subpopulation itself. Therefore $60/(60 + 10)$ or $6/7$ of the complement of this subpopulation consists of a subpopulation that is genetically identical to the reference subpopulation. Three-fourths of the complement of the 60% subpopulation is genetically identical to its reference subpopulation. Hence in terms of differentiation between subpopulations and their complements, the 10% subpopulation is the most differentiated, the 30% subpopulation the least, and the 60% subpopulation ranks between them. The amount of subpopulation differentiation measured on a population basis should be somewhere between the maximum and minimum levels occurring for the individual subpopulations, i.e. it should be some type of mean for the individual subpopulation measures of differentiation. In the example, even thought the 10% subpopulation has the greatest degree of differentiation, it is undesirable to have it influence the measure of differentiation beyond its contribution to the number of individuals in the population. Therefore it is appropriate to express the total amount of subpopulation differentiation as the weighted average of the individual subpopulation differentiation measures, using the relative sizes of the subpopulations as weights.

Differentiation within populations is related to the degree in which each component subpopulation differs from the remainder of the population, or as is emphasized in the example, the degree of differentiation between the subpopulations and their complement subpopulations (subpopulation differentiation). It follows then that the manner in which the complements

are subdivided into component subpopulations is irrelevant to subpopulation differentiation. This is particularly evident for subpopulations which have the same genetic properties as the entire population. For this case, the complement subpopulations have the same genetic properties as the reference subpopulations. These subpopulations thus are not differentiated from their complements and do not contribute to the differentiation within the population, regardless of the subdivision of the complements.

An additional condition that is desirabe for a measure of differentiation to satisfy is also illustrated by the above example. If the two larger, genetically identical subpopulations are considered to be one subpopulation, the subpopulation structure that results has a greater degree of differentiation than the initial three-subpopulation structure. In general, by combining genetically identical subpopulations to reduce their number, the amount of differentiation among subpopulations is inflated and a measure of differentiation should reflect such an increase.

These conditions, which a measure of subpopulation differentiation ($D$) should satisfy, can be succinctly summarized as follows:

a) $D$ should explicitly indicate the subpopulation structure of populations and specifically reflect the homogenizing of subpopulations with identical genetic compositions.

b) $D$ should be in the interval $0 \leq D \leq 1$ with $D = 0$ if and only if all subpopulations have identical genetic compositions, and $D = 1$ if and only if all subpopulations are genetically unique for the set of loci being investigated.

c) Each subpopulation should be characterized by the genetic differences between it and its complement subpopulation. A measure of these differences for a subpopulation indicates the amount of differentiation attributable to that subpopulation. The substructure of the complement of a subpopulation should not have an effect on its measure of differentiation.

d) $D$ should be the weighted mean of the measures of differentiation for the individual subpopulations with weights given by the relative sizes of the subpopulations. Thus if $c_j$ represents the proportion of individuals in the $j$th subpopulation and $D_j$ represents the level of differentiation for the $j$th subpopulation then $D = \sum_j c_j D_j$ where $\sum_j c_j = 1$.

e) Reducing the number of subpopulations by combining those with identical genetic compositions should increase $D$.

It is apparent that a measure $D$ will meet these conditions only if the $D_j$'s are defined properly. Since $D_j$ is a measure of the genetic difference between the $j$th subpopulation and its complement subpopulation,

it is appropriate that it should be a measure of genetic distance. Minimum requirements for such a distance measure are that it fall within the interval [0, 1], is zero if and only if the subpopulation and its complement are genetically identical and is one if and only if the subpopulation and its complement have completely different genetic compositions. Thus $0 \leq \sum_j c_j D_j \leq 1$ holds.

In the following development it will be shown that, among the many genetic distance measures, there is only one with desirable properties which also has an explicit interpretation for its values. This latter characteristic is extremely important when the magnitude of the differences between subpopulations rather than the ranking of differences is of primary concern.

## An appropriate genetic distance

The concept of genetic distance deals with quantifying differences between the frequency of genetic components in two populations. These components may be genes, gametes, or genotypes and the frequencies are specified to be relative frequencies. The information desired is the degree to which the various gentic units (alleles, multilocus gametes, genotypes) do not occur equally in the two populations. It can be measured by comparing absolute frequencies between the two populations. If the populations are not of equal size, the absolute frequencies must be transformed in such a way that differences are based upon absolute frequencies in hypothetical populations of the same size. These reference populations must be of equal size since results are to be in interpreted on a common basis.

Consider two populations

$$P = (n_1, n_2, \ldots, n_k) \quad \text{and} \quad Q = (m_1, m_2, \ldots, m_k),$$

where $n_i$ and $m_i$ indicate the absolute numbers of genetic elements of the $i$th type in the respective populations, both of which have $k$ different types. Let $\sum_{i=1}^{k} n_i = N$ and $\sum_{i=1}^{k} m_i = M$, where $N \neq M$, if the populations are not equal in size.

Now construct the hypothetical populations

$$P' = (n'_1, n'_2, \ldots, n'_k) \quad \text{and} \quad Q' = (m'_1, m'_2, \ldots, m'_k),$$

where $n'_i = \max(N, M) p_i$ and $m'_i = \max(N, M) q_i$. The $p_i$ and $q_i$ are relative frequencies for the $i$th type in $P$ and $Q$, respectively, with $p_i = n_i/N$ and $q_i = m_i/M$. If $N = \max(N, M)$ then $P' = P$ and likewise, if $M = \max(N, M)$ then $Q' = Q$. $P'$ and $Q'$ are genetically similar to $P$ and $Q$ in the sense that $P'$ has relative frequencies that are identical to those of $P$ and, likewise, the relative frequencies in $Q'$ are identical to those of $Q$. Moreover $P'$ and $Q'$ have an equal number

of genetic elements, i.e.,

$$\sum_{i=1}^{k} n_i' = \sum_{i=1}^{k} m_i' = \max(N, M).$$

The number of genetic elements of the $i$th type by which $P'$ and $Q'$ differ can be expressed as $|n_i' - m_i'|$. And the total number of such differences is $S = \sum_{i=1}^{k} |n_i' - m_i'|$. Thus the ratio $0.5\ S \max(N, M)$ is the proportion of genetic elements in $P'$ and $Q'$ by which the two populations differ. Equivalently

$$\frac{1}{2} S \max(N, M) = \frac{1}{2} \sum_{i=1}^{k} |p_i - q_i|.$$

A similar argument can be made for hypothetical populations with number of genetic elements set equal to $\min(N, M)$. For this case,

$$\frac{1}{2} S' \min(N, M) = \frac{1}{2} \sum_{i=1}^{k} |p_i - q_i|,$$

where $S'$ represents the total number of differences in genetic elements between two hypothetical populations having relative frequencies of $p_i$ and $q_i$. The expression on the right in the last two equations is the measure of genetic distance $d_0(p, q)$ introduced by Gregorius (1974). It has properties of a metric and is maximal only for genetically completely different populations. Furthermore $d_0$ is unique for several geometrical properties which are intuitively desirable and useful (Gregorius 1984). The above development reveals an additional important aspect of $d_0$, namely that its values are interpretable without resorting to geometrical concepts. For example, $d_0(p, q) = 0.3$, means that the two populations have a difference in number of genetic elements that is equivalent to a difference involving 30% of the elements in two reference populations of equal size.

## A measure of subpopulation differentiation

In the following treatment, the functional relationship between a measure of subpopulation differentiation and $d_0(p, q)$ is developed. Let $c_j$ denote the relative size of the $j$th subpopulation, where $\sum_j c_j = 1$ and let $p_i(j)$ denote the relative frequency of the $i$th genetic type in the $j$th subpopulation, so that $\sum_i p_i(j) = 1$ in each subpopulation. Then, the relative frequency of the $i$th genetic type in the whole population is given by $p_i = \sum_j p_i(j)\, c_j$, and the relative frequency of this type in the complement of the $j$th subpopulation is $\bar{p}_i(j) = \sum_{k \neq j}^{k} p_i(k)\, c_k/(1 - c_j) = (p_i - c_j\, p_i(j))/(1 - c_j)$. These weights and relative frequencies can be considered as components of the following vectors:

$$p(j) = (p_1(j), p_2(j), \dots);$$
$$\bar{p}(j) = (\bar{p}_1(j), \bar{p}_2(j), \dots);$$

$$p = (p_1, p_2, \dots);$$
$$c = (c_1, c_2, \dots).$$

The amount of differentiation $D_j$ for the $j$th subpopulation is defined to be the genetic distance between the subpopulation and its complement. The measure $D_j$ can be expressed as

$$D_j = d_0(p(j), \bar{p}(j)) = \frac{1}{2} \sum_i |p_i(j) - \bar{p}_i(j)|.$$

Since $p_i(j) - \bar{p}_i(j) = (p_i(j) - p_i)/(1 - c_j)$, $D_j$ may also be written in the form

$$D_j = d_0(p(j), p)/(1 - c_j).$$

This form reveals that subpopulation differentiation is related to the deviation of subpopulations from the entire population. It demonstrates that the $j$th subpopulation is undifferentiated, i.e. $D_j = 0$, if and only if its genetic composition is equivalent to that of the whole population, i.e. $p(j) = p$. It also shows that, since $d_0(p(j), \bar{p}(j)) \leq 1$, the distance of a subpopulation from the whole population is always $\leq (1 - c_j)$. Consequently, a completely differentiated subpopulation has a distance from the total population that is equal to the relative size of its complement subpopulation.

According to the previously given definition, the amount of differentiation among subpopulations is

$$\delta = \sum_j c_j\, d_0(p(j), \bar{p}(j)) = \sum_j d_0(p(j), p)\, c_j/(1 - c_j).$$

Conditions (a) to (d) are satisfied by $\delta$. To prove that (e) is also satisfied by $\delta$, assume, without loss of generality, that $j > 2$ and that subpopulations 1 and 2 are genetically identical, i.e. $p(1) = p(2)$. Considering these two subpopulations as one, the amount of differentiation is

$$\delta' = d_0(p(1), p)(c_1 + c_2)/(1 - c_1 - c_2)$$
$$+ \sum_{j \geq 3} d_0(p(j), p)\, c_j/(1 - c_j).$$

Hence $\delta' - \delta = d_0(p(1), p)\,[(c_1 + c_2)/(1 - c_1 - c_2) - c_1/(1 - c_1) - c_2/(1 - c_2)] = d_0(p(1), p)\, c_1\, c_2(2 - c_1 - c_2)/[(1 - c_1 - c_2)(1 - c_2)(1 - c_2)(1 - c_2)] > 0$ if $p(1)$ is not identical to $p$.

The explicit interpretation of $d_0$ can now be extended to $\delta$. Thus $D_j$ represents the proportion of genetic elements by which the $j$th subpopulation differs from a hypothetical subpopulation of the same size but which has the same relative frequency as the complement subpopulation. It is permissible to express the proportion of differences associated with each subpopulation in terms of subpopulation size since $d_0$ can be interpreted on the basis of either the maximum or minimum sizes of the subpopulations involved. The number of genetic elements by which a subpopulation differs from its hypothetical population constructed as

described above is the effective number of genetic elements by which the subpopulation differs from its complement subpopulation. The maximum effective number of genetic elements by which a subpopulation may differ from its complement is equal to twice the number of genetic elements in the subpopulation. Therefore $D_j$ can also be interpreted as the proportion of the effective number of genetic elements by which the $j$th subpopulation differs from its complement. It follows then that $\delta$ represents the mean proportion of effective numbers of genetic elements by which subpopulations differ from their complements. This average also is a proportion. Consider that $d_0(p(j), \bar{p}(j))$ multiplied by the subpopulation size becomes the effective number of genetic elements by which the $j$th subpopulation differs from its complement. Then $\delta$, which is the sum of these effective numbers divided by the population size, can be interpreted as the proportion of the total effective number of genetic elements in the population by which the subpopulations differ from their respective complements. In this sense, $\delta$ measures the proportion of genetic disparity among the subpopulations.

The development so far has referred to genetic elements and genetic types without specifically addressing genes, gametes, or genotypes. Particularly when multiple loci are considered, the definition of frequencies of genetic types can be ambiguous. This problem occurs when allele frequencies are obtained for each locus separately, or gametic frequencies are computed from frequencies of unordered genotypes (e.g. when male and female contributions cannot be distinguished). For the first case it is common to take the mean of single locus measurements (distance, differentiation) over loci. These means, however, are difficult to interpret unless they are derived from a multilocus concept of gene frequencies; allele frequencies refer to a single locus. Similarly, measures based on gametic frequencies are not useful unless a method of counting gametes is defined.

## Differentiation measures involving multiple loci

Consider the situation in which genes are studied for each locus separately so that the association of genes on chromosomes or in DNA sequences is not taken into consideration. A count of genes in a population studied in this manner results in a description of the population's gene pool. A gene pool is defined to be the set of all genes in a population found at a specified group of loci. If $L$ loci are investigated and each is represented to the same degree, then the set of all individual genes belonging to a particular locus is $1/L$th of the gene

pool. Within this set, the $i$th allele has a relative frequency $p_{in}$, where $n$ indicates the locus ($1 \leq n \leq L$), and $\sum_i p_{in} = 1$.

Consequently, a gene identified as the $i$th allele at the $n$th locus has frequency $p_{in}/L$ in the gene pool. If loci vary in level of representation and $y_n$ is the level for the $n$th locus, the frequency of the $i$th allele at the $n$th locus in the gene pool is $y_n p_{in}/\sum_{n=1}^{L} y_n$. Varying levels of representation would arise, for example, in situations in which some loci are replicated in the genome but others are not. In general, if $t_n$ indicates the proportion of genes in the gene pool that belong to the $n$th locus, $t_n p_{in}$ is the frequency of the $i$th allele at this locus in the gene pool. The measure of differentiation, $\delta$, when based on these frequencies, is a measure of gene-pool differentiation among subpopulations and will be denoted by $\delta_{ge}$. If the loci are represented in the same proportions in all subpopulations, then gene frequencies within subpopulations can be written $t_n p_{in}(j)$ and we obtain

$$\delta_{ge} = \sum_j [c_j/2(1-c_j)] \sum_{i,n} t_n |p_{in}(j) - p_{in}|$$

$$= \sum_{n=1}^{L} t_n \{\sum_j [c_j/2(1-c_j)] \sum_i |p_{in}(j) - p_{in}|\}.$$

The term in braces is the amount of differentiation at the $n$th locus. Hence, a measure of the amount of gene pool differentiation for a given set of loci is the weighted average of the single-locus measures with weights given by the locus frequencies in the gene pool, i.e.

$$\delta_{ge} = \sum_{n=1}^{L} t_n \delta_{ge}^{(n)},$$

where $\delta_{ge}^{(n)}$ is the measure of differentiation at the $n$th locus. It is clear then that averaging over loci is meaningful, provided the gene pool is the reference for defining gene frequencies and that the distance measure on which the differentiation measure is based has the required linearity properties. If neither of these holds, averages over loci are not appropriate.

Gene-pool differentiation measures subpopulation differences at a lower level of organization than the gametic and genotypic levels. Furthermore, it depends on neither genotypic nor gametic associations of genes. Differentiation at higher levels of organization can be studied by redefining the genetic entities to be counted and applying the measure $\delta$ to the frequencies of these entities. Let $\delta_{ga}$ indicate differentiation based on gametic frequencies, and $\delta_{go}$ differentiation based on genotypic frequencies. Use of these measures requires that the pertinent genetic entities can be identified. Often gametic frequencies cannot be determined because only frequencies of unordered, multilocus genotypes are available. However, in some cases it is

possible to determine ordered genotypes at the zygotic-stage, for example in gymnosperm seeds. The genetic constitution of the haploid endosperm can be compared with that of the diploid embryo (Müller 1976). This technique permits the separation of the female contribution from that of the male in the zygote and thus allows the computation of gametic frequencies. However, in later ontogenetic stages, the information obtained from gametic frequencies may be of limited value since genotypic constitution during the later stages of development reflects primarily the outcome of viability selection acting on unordered genotypes rather than on gametic types.

Since gametic frequencies are obtained by taking sums of certain genotypic frequencies, and gene frequencies are obtained in a similar way from gametic frequencies, the representation of $d_0$ implies

$$\delta_{ge} \leq \delta_{ga} \leq \delta_{go}$$

if the same genotypic frequency distribution is used for all three measures of differentiation. This result follows immediately from the triangle inequality in the form of

$$\left| \sum_i a_i - \sum_i b_i \right| \leq \sum_i |a_i - b_i|.$$

Therefore, the monotone relationship among these measures indicates that increasing amounts of differentiation among subpopulations occur with increasing degrees of organization of the genetic material. Such an ordering agrees with biological expectation, since higher degrees of organization imply greater complexity and therefore additional opportunities for differentiation.

Pairwise differences between the measures $\delta_{ge}$, $\delta_{ga}$, and $\delta_{go}$, expressed as $\delta_{ga} - \delta_{ge}$, $\delta_{go} - \delta_{ga}$ and $\delta_{go} - \delta_{ge}$, measure the effect of increases in organization on levels of subpopulation differentiation. If one of these differences equals zero, organization at the higher level does not result in greater subpopulation differentiation. If instead, a difference is greater than zero, it is a measure of the increase in proportion of differences in genetic units associated with greater organization.

For some situations, the use of multilocus frequencies may not be desirable for detecting isolation between portions of a population. This is illustrated by considering the extreme case in which each subpopulation has a unique allele at a particular locus but at the same time is equal in frequency for all alleles of the remaining loci. Clearly, the locus with variable allele frequencies indicates the absence of gene flow between subpopulations and shows that the subpopulations are completely isolated from each other. Yet, the gene-pool measure of differentiation at $L$ loci with identical degrees of representation is $\delta_{ge} = 1/L$, and on this basis the subpopulations would be mistakenly interpreted as

having a low degree of isolation if, for example, $L \geq 10$. Hence, if $\delta$ is to be interpreted in terms of gene flow among subpopulations it is probably more appropriate to compute $\delta_{ge}$ values for each locus separately and then to make comparisons between these values.

## Comparison of the $\delta$ and $F_{ST}$ measures

Wright's $F_{ST}$ measure is based on relative gene frequencies and (applying the present notation for allelic frequencies) can be written in the form

$$F_{ST} = \left[ \sum_{n=1}^{L} \sum_i \left( \sum_j c_j p_{in}^2(j) - p_{in}^2 \right) \right] \left[ \sum_{n=1}^{L} \sum_i p_{in}(1-p_{in}) \right].$$

This is the $L$-locus version (Wright 1978, p. 102) which can also be expressed as $\sum_{n=1}^{L} t_n F_{ST}^{(n)}$, where $F_{ST}^{(n)}$ is the measure for the $n$th locus and $t_n = \sum_i p_{in}(1 - p_{in})/ \sum_{n=1}^{L} \sum_i p_{in}(1 - p_{in})$. Clearly, $0 \leq F_{ST} \leq 1$, and $F_{ST} = 0$ if and only if there is no subpopulation differentiation. However, $F_{ST} = 1$ if and only if all subpopulations are fixed at all loci, so that $F_{ST} < 1$ if the number of alleles exceeds the number of subpopulations for at least one locus. Therefore, $F_{ST}$ does not fulfill the second part of our condition (b). Moreover, it also does not fulfill conditions (a) and (e) since combining subpopulations with identical gene pools does not result in a change in the value of $F_{ST}$. The latter is easily verified, while the first follows from

$$p_{in}(1 - p_{in})$$
$$- (\sum_j c_j p_{in}^2(j) - p_{in}^2) = \sum_j c_j p_{in}(j)(1 - p_{in}(j)) \geq 0.$$

For $F_{ST} = 1$ the last sum equals zero which happens only if $p_{in}(j) = 0$ or $p_{in}(j) = 1$ for each subpopulation. This proves the above assertion for $F_{ST} = 1$.

On the other hand, $F_{ST}$ meets condition (d) if we choose

$$D_j = \sum_{n=1}^{L} \sum_i (p_{in}^2(j) - p_{in}^2) \bigg/ \sum_{n=1}^{L} \sum_i p_{in}(1 - p_{in}),$$

since then $F_{ST} = \sum_j c_j D_j$. The $D_j$'s also satisfy condition (c) in so far as they don't depend on the structure of the respective subpopulation complements, however, they are not based on differences between the subpopulations and their complements and are difficult to interpret in terms of genetic differences.

In this paper we shall not explicitly address statistical problems of estimation of the $\delta$-measures. For the time being it may suffice to point out that applying relative frequencies from samples taken from the individual subpopulations to the $D_j$- and $\delta$-measures leads to statistically consistent esimators. Moreover, the construction of the $\delta$ implies that a significance test for population differentiation is obtained by simply com-

paring the sample frequencies for a subpopulation with the sample frequencies for the remainder of the population. If this is done for each subpopulation, and if at least one subpopulation differs significantly from its respective remainder, then the population as a whole can be considered to be significantly differentiated. At the same time this procedure reveals the subpopulations which contribute to the population differentiation.

## Subpopulation differentiation in an experimental data set

Use of $\delta$ for measuring differentiation among subpopulations is illustrated with allelic and genotypic data reported by Linhart et al. (1981) for the following seven protein loci, peroxidase (PER), flourescent esterase (FE), colorimetric esterase (CE) phosphoglucomutase-1 (PGM-1), phosphoglucomutase-2 (PGM-2), phosphohexose isomerase (PHI) and glutamate dehydrogenase (GDH). These data were collected for investigation of genetic structure within a sexually mature population of ponderosa pine (Pinus ponderosa Laws.). The population was composed of six spatially discrete subpopulations of varying sizes, all located within an area of approximately two hectares. Data were obtained from all individuals of reproductive age in the area occupied by each subpopulation, resulting in a complete genetic description of these groups in terms of allelic and genotypic frequencies for the seven loci studied.

On the basis of tests for heterogeneity of allele frequencies as well as values computed for Nei's

distance and Wright's $F_{ST}$, the original investigators concluded that the subpopulations were genetically differentiated at the genic level for two loci, FE and CE. However, $F_{ST}$ values recalculated using relative subpopulation sizes as weights reveal little differentiation among subpopulations (Table 1). Here we present values for the $\delta$ measure of differentiation for this population. Values for differentiation among subpopulations ($D_j$) and population differentiation at the genic level ($\delta_{ge}^{(n)}$) are given in Table 1 and are illustrated in Fig. 1. Values for $D_j$ and for differentiation at the genotypic level ($\delta_{go}^{(n)}$) are presented in Table 2. As previously discussed these values can be interpreted in terms of the proportion of the effective numbers of genetic elements (genes or genotypes) that differ between the subpopulations.

Values for $\delta$, both at the genic and genotypic levels, indicate substantial differentiation at the FE, CE, and PER loci. At the FE locus, subpopulations have differences at 11% of the total effective number of genes ($\delta_{ge}^{(n)} = 0.113$) and at 17% of the total effective number of genotypic units ($\delta_{go}^{(n)} = 0.170$) in the population. The remainder of the loci (PGM-1, PGM-2, PHI, and GDH) showed only small amounts of differentiation.

Comparison of the $D_j$'s for the individual subpopulations reveals that subpopulation C is substantially more differentiated from the remainder of the population than the remaining subpopulations. On the gene pool basis, the mean $D_j$ computed over loci indicates that subpopulation C differs from the remainder of the population at over 11% of the effective number of genes at these loci. The same pattern holds for the genotypic level of organization. In view of the small

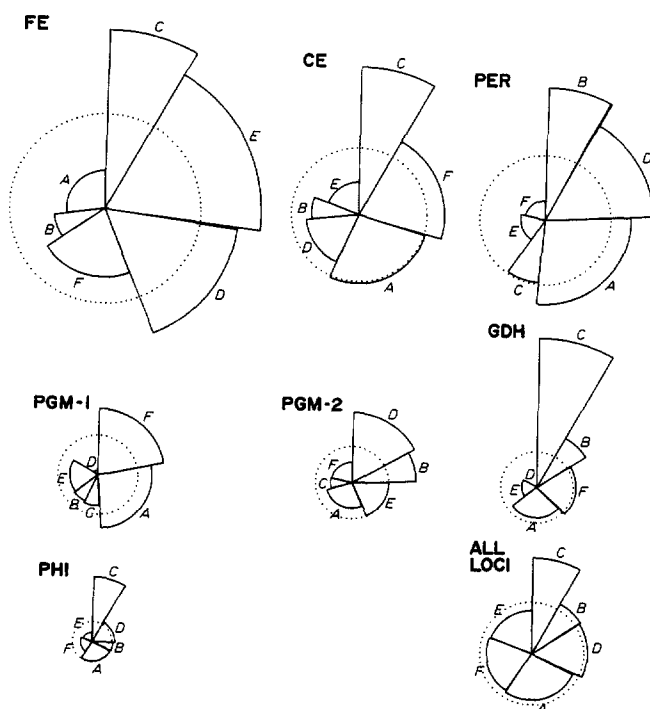**Table 1.** Subpopulation differentiation at the genic level

| Loci | | | | | | | | | Subpopulation | |
|---|---|---|---|---|---|---|---|---|---|---|
| Sub-population | PER[a] | FE | CE | PGM-1 | PGM-2 | PHI | GDH | Mean | Size | Relative weight |
| | $D_j$ | | | | | | | | | |
| A | 0.101 | 0.046 | 0.081 | 0.064 | 0.031 | 0.023 | 0.036 | 0.055 | 54 | 0.271 |
| B | 0.159 | 0.061 | 0.056 | 0.035 | 0.076 | 0.024 | 0.068 | 0.068 | 15 | 0.075 |
| C | 0.075 | 0.215 | 0.178 | 0.037 | 0.025 | 0.077 | 0.179 | 0.112 | 16 | 0.080 |
| D | 0.129 | 0.158 | 0.062 | 0.006 | 0.084 | 0.027 | 0.001 | 0.067 | 33 | 0.166 |
| E | 0.029 | 0.185 | 0.039 | 0.033 | 0.043 | 0.012 | 0.017 | 0.051 | 38 | 0.191 |
| F | 0.024 | 0.082 | 0.101 | 0.080 | 0.025 | 0.014 | 0.047 | 0.053 | 43 | 0.216 |
| $\delta_{ge}^{(n)}$ | 0.077 | 0.113 | 0.080 | 0.048 | 0.044 | 0.024 | 0.043 | 0.061 | | |
| $F_{ST}$ | 0.024 | 0.034 | 0.020 | 0.029 | 0.020 | 0.024 | 0.014 | 0.024 | | |

[a] Abbreviations indicate the following allozyme loci: PER = peroxidase; FE = fluorescent esterase; CE = colorimetric esterase; PGM-1 = phosphoglucomutase-1; PGM-2 = phosphoglucomutase-2; PHI = phosphohexose isomerase; GDH = glutamate dehydrogenase

**Table 2.** Differentiation among subpopulations at the genotypic level

| Sub-population | PER[a] | FE | CE | PGM-1 | PGM-2 | PHI | GDH | Mean | Size | Relative weight |
|---|---|---|---|---|---|---|---|---|---|---|
| | $D_j$ | | | | | | | | | |
| A | 0.181 | 0.095 | 0.171 | 0.128 | 0.056 | 0.046 | 0.119 | 0.114 | 54 | 0.271 |
| B | 0.350 | 0.188 | 0.170 | 0.069 | 0.135 | 0.049 | 0.072 | 0.148 | 15 | 0.075 |
| C | 0.133 | 0.326 | 0.264 | 0.074 | 0.085 | 0.155 | 0.190 | 0.175 | 16 | 0.080 |
| D | 0.240 | 0.204 | 0.076 | 0.011 | 0.151 | 0.054 | 0.033 | 0.110 | 33 | 0.166 |
| E | 0.039 | 0.230 | 0.062 | 0.066 | 0.072 | 0.023 | 0.085 | 0.082 | 38 | 0.191 |
| F | 0.099 | 0.121 | 0.180 | 0.160 | 0.088 | 0.028 | 0.066 | 0.106 | 43 | 0.216 |
| $\delta_{go}^{(n)}$ | 0.155 | 0.170 | 0.144 | 0.095 | 0.090 | 0.048 | 0.089 | 0.113 | | |
| $\delta_{ge}^{(n)}/\delta_{go}^{(n)}$ | 0.500 | 0.667 | 0.557 | 0.500 | 0.484 | 0.500 | 0.482 | 0.542 | | |

[a] Abbreviations indicate the following allozyme loci: $PER$ = peroxidase; $FE$ = fluorescent esterase; $CE$ = colorimetric esterase; $PGM$-1 = phosphoglucomutase-1; $PGM$-2 = phosphoglucomutase-2; $PHI$ = phosphohexose isomerase; $GDH$ = glutamate dehydrogenase



**Fig. 1.** Subpopulation differentiation at the genic level. The dotted circles have radii equal to the genic level of differentiation ($\delta_{ge}^{(n)}$). The solid sectors represent the contributions of the subpopulations to the total subpopulation differentiation. Radii of the solid sectors are equal to the amounts of differentiation for the individual subpopulations ($D_j$'s) and the angles of the sectors represent the subpopulation weights ($c_j$'s). Abbreviations represent the following allozyme loci: $FE$ = fluorescent esterase, $CE$ = colorimetric esterase, $PER$ = peroxidase, $PGM$-1 = phosphoglucomutase-1, $PGM$-2 = phosphoglucomutase-2, $PHI$ = phosphohexose isomerase, $GDH$ = glutamate dehydrogenase

size of this population, only 16 individuals, the possibility that genetic drift contributed to the observed differentiation cannot be overlooked.

At the gene-pool level of organization as well as at the genotypic level, values for $\delta$ indicate that substantial variation in differentiation existed among loci. It thus seems unlikely that the level of differentiation among subpopulations is entirely controlled by gene flow since migration alone would be expected to influence differentiation at all loci to about the same degree. It should be pointed out in this regard that $F_{ST}$ values are uniformly small for all loci and that this measure does not reflect the variation in differentiation across loci as does the $\delta_{ge}^{(n)}$ measure. In this respect, results for $\delta_{ge}^{(n)}$ are in accord with inferences drawn from the tests for heterogeneity of allele frequencies, whereas the $F_{ST}$ results are not. Clearly for this experiment, use of $F_{ST}$ as a measure of differentiation results in an interpretation of the data that is quite different from that obtained when $\delta_{ge}^{(n)}$ is used.

Other differences in $\delta_{ge}^{(n)}$ and $F_{ST}$ for this data set are $F_{ST} \leq \delta_{ge}^{(n)}$ for all loci and the rank order of the loci, with respect to the amount of differentiation, is different for $\delta_{ge}^{(n)}$ and $F_{ST}$. These differences make it evident that $F_{ST}$ and $\delta_{ge}^{(n)}$ are not equivalent measures of subpopulation differentiation.

Differences between the $\delta$'s for the genotypic and gene pool level of organization reveal that there was a substantial increase in differentiation with the increase in complexity of organization. The amount of differentiation at the genotypic level was approximately twice that at the genic level for all loci except $FE$ and $CE$.

For these two loci, although the increase in differentiation was substantial, the relative increase in differentiation was not as great as for the other loci.

# References

Ehrlich PR, Raven PH (1969) Differentiation of populations. Science 165:1227–1232

Gregorius H-R (1974) Genetischer Abstand zwischen Populationen. 1. Zur Konzeption der genetischen Abstandsmessung. Silvae Genet 23:22–27

Gregorius H-R (1984) A unique genetic distance. Biom J 26:13–18

Hamrick JL (1983) The distribution of genetic variation within and among natural plant populations. In: Schonewald-Cox CM, Chambers SM, MacBride B, Thomas L (eds) Genetics and conservation. Benjamin/Cummings, Menlo Park, pp 335–349

Linhart YB, Mitton JB, Sturgeon KB, Davis ML (1981) Genetic variation in space and time in a population of ponderosa pine. Heredity 46:407–426

Müller G (1976) A simple method of estimating rates of self-fertilization by analyzing isozymes in tree seeds. Silvae Genet 25:15–17

Wright S (1969) Evolution and the genetics of populations, vol 2. University of Chicago Press, Chicago

Wright S (1978) Evolution and the genetics of populations, vol 4. University of Chicago Press, Chicago