# How to use ANOVA in Qlucore Omics Explorer

## Contents

# 1 INTRODUCTION

This how to document outlines one method for selecting a suitable statistical model in Qlucore Omics Explorer, for analyzing data acquired using a given experimental design. It is assumed that the reader has a certain understanding of statistical tests, p-values and q-values, and we refer to statistical textbooks for more comprehensive treatments of these concepts. Normally a data set is described by a matrix where the columns represent samples and the rows represent variables.

# 2 INTRODUCTION TO ANOVA

The ANOVA framework that is used as the basis for the statistical tests in Qlucore Omics Explorer is very flexible and can be used to analyze data from many different experimental designs, with the aim of finding variables (for example, genes) that are highly related to a specific experimental factor. In this document we cover the following ANOVA models:

## 2.1 MODEL A:  INDEPENDENT MEASURES, SINGLE FACTOR DESIGN (A)

This model is used when one experimental factor is varied in a collection of subjects, and each subject is randomly allocated to receive a specific level of the factor. For example, consider a study where a group of patients are randomly allocated either to a group receiving a treatment or to a control group receiving placebo.

## 2.2 MODEL B: INDEPENDENT MEASURES, FACTORIAL DESIGN (B)

This model is used when more than one experimental factor is varied, and each subject is randomly assigned a specific level of each factor. As an example, consider a study where each patient is randomly allocated to either a treatment or a control group, and furthermore is studied either one or two weeks after the treatment or placebo administration.

## 2.3 MODEL C: REPEATED MEASURES, SINGLE FACTOR DESIGN (C)

This model is used when a single experimental factor is varied in a group of subjects, and each subject receives all levels of the factor. An example of this type is given by a study where each patient receives both the treatment and placebo. The tests used for this type of designs are sometimes referred to as paired tests.

## 2.4 MODEL D: REPEATED MEASURES, FULLY FACTORIAL DESIGN (D)

This model is suitable for experimental designs where several experimental factors are varied and each subject receives all combinations of the factor levels. Consider for example a study where each patient is given both the treatment and placebo and furthermore is studied both one and two weeks after each administration.

## 2.5 MODEL E: REPEATED MEASURES, MIXED DESIGN (E)

This model is used when several experimental factors are varied, and each subject receives all levels of some factors but only one level of other factors. For example, this model can be used in a situation where the patients are given either the treatment or placebo, but all patients are studied both one and two weeks after the administration.

We will also cover the situation where there are covariates, that is, factors that are not controlled in the experiment but still may influence the measured variables, and hence need to be considered.

The built statistical methods handle all examples discussed above. It is also possible to enable the Extended statistics functionality using the Open API to R. See the "How to use the R interface" – document for more information.

## 3 TERMINOLOGY

We will assume that we have collected a data set, consisting of several variables and a number of samples, and that we have a data file containing the experimental data, in a format suitable for Qlucore Omics Explorer (see the Reference Manual for more information about file formats). We will use the following terminology:

• An **independent variable** is a variable that is being controlled in the experiment. These are also called factors. Typical examples are treatments and doses. In Qlucore Omics Explorer, they are represented by sample annotations.

• A **dependent variable** is a variable that is affected by changes in the independent variables. The dependent variables are often those that are observed in the experiment, for example expression levels of genes or protein levels. In Qlucore Omics Explorer, a separate statistical test is performed for each dependent variable, to assess the effect of the independent variables.

• A **covariate** is a variable that is not controlled in the experiment, but still may affect the dependent variables. Typical examples of covariates are age, gender and batch effects that can **be linked to, for** example, the time of data acquisition. In Qlucore Omics Explorer, covariates are represented by sample annotations.

• A **main effect** is the effect of a single factor on the dependent variable.

• An **interaction effect** is an effect obtained by combining more than one factor. For example, the expression levels of a gene may increase with time after the administration of one treatment, and decrease with time after the administration of another treatment. In this case, it is the interaction between the chosen treatment and the time after administration that determines the expression level of the gene in a specific sample.

We will furthermore make a distinction between three different types of variables:

• **Nominal** variables are categorical variables, that is, variables whose value in a sample can be one of several categories. The possible categories of a nominal variable do not induce a natural ordering of the samples. Typical examples of nominal variables are gender and disease type.

• **Ordinal** variables are categorical variables that can be used to order the samples. The actual values of an ordinal variable are not important, only their relative orders. For example, coding the possible values of an ordinal variable as 1, 2 and 3 gives the same result as coding them as

10, 200 and 4000. An example of an ordinal variable is dose level, which can be coded as low, medium or high dose.

• **Continuous** variables are quantitative variables, where the actual, observed values are meaningful. Typical examples of continuous variables are age and the given dose of a drug (measured in, for example, mg/kg).

The statistical framework in Qlucore Omics Explorer supports differentiated treatment of variables of different types. For example, depending on the type of an experimental factor F, different types of statistical tests can be used to find those dependent variables that are associated with F. In Qlucore Omics Explorer, the type of statistical test is selected in the Statistics window.

• If F is nominal (for example disease subtype) with two levels, select Two Group Comparison as the test type. For this test type, it is also possible to select a test direction, that is, whether to search for any dependent variables showing a difference between the two levels of F (by choosing the test direction $\neq$), or variables showing a difference in a specific direction (by choosing > or <). This test is equivalent to a standard t-test between the groups defined by the two levels of F (a one-sided or two-sided test, depending on the chosen test direction).

• If F is nominal with multiple levels, select Multi Group Comparison as the test type. A Multi Group Comparison for a factor with two levels is equivalent to a non-directed Two Group Comparison.

• If F is ordinal (for example dose level), Rank Regression can be used to find dependent variables which either increase or decrease following the level of F. Also in this case we can select a test direction. Note that for rank regression the order of the factor levels, as they appear in the Samples window, is important (see the Reference Manual for more information).

• If F is continuous (for example dose), we can use Linear Regression or Quadratic Regression to find variables whose values are linearly or quadratically related to the value of F. For linear regression, a test direction can be selected.
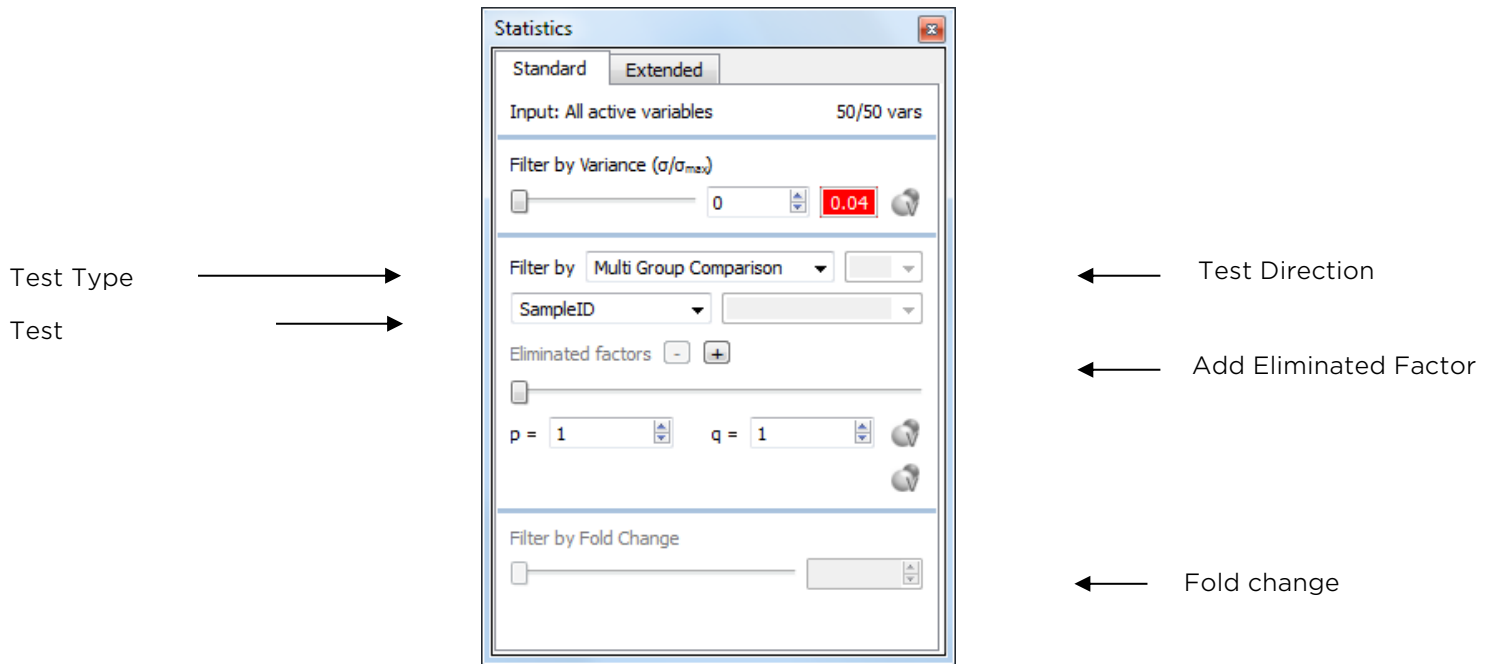
Test Type

Test

Test Direction

Add Eliminated Factor

Fold change

Figure 1

When an eliminated factor is added the 'Data Type' for the factor can be chosen.

## 4 METHOD

We will now outline one method for choosing a suitable statistical model for analyzing data obtained with a given experimental design. In this method, you answer up to five questions regarding your experimental design. The answers to these questions lead you to a suitable statistical model for your data. The green circles, numbered 1-5 in the flow chart, represent the five questions. The blue circles, labeled A-E, represent the statistical models. When given a new data set, begin with Question 1. In figure 1 we present an overview of the decision tree. The questions are marked with a green circle.
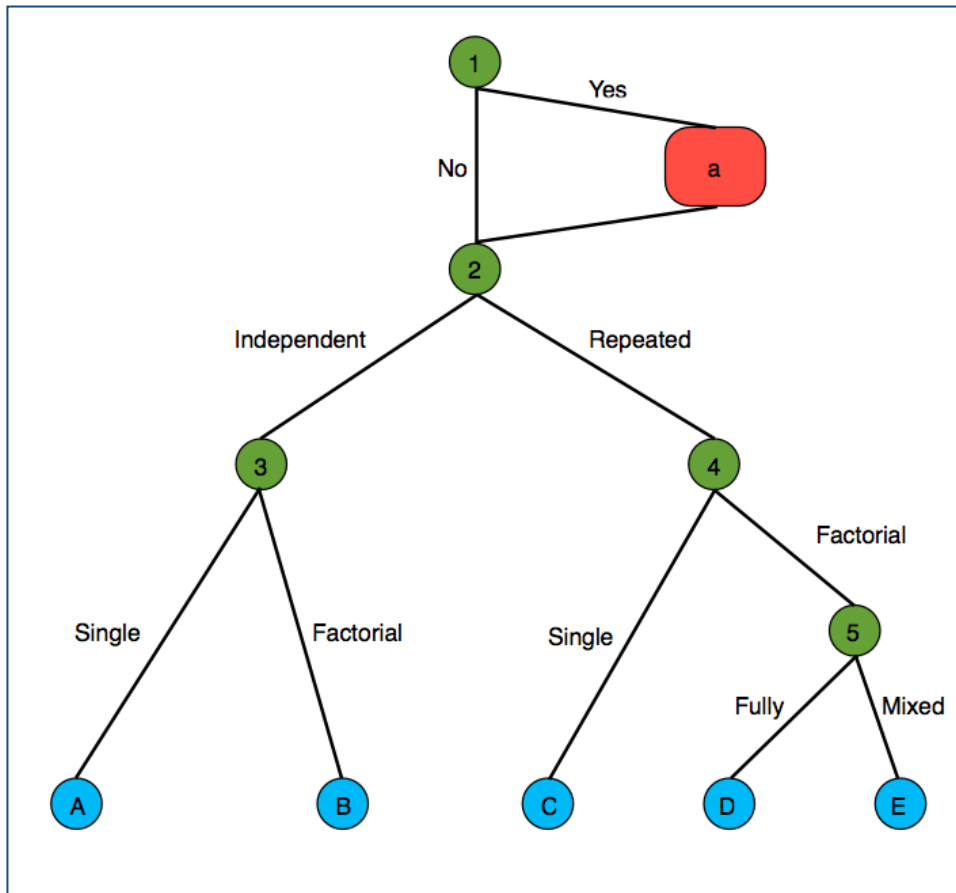
Figure 2

## 4.1 QUESTION 1: ARE THERE ANY COVARIATES IN THE DATA?

A covariate is a variable that is not controlled in the experiment, but may affect the dependent variables. In the data file, each covariate should be represented by an annotation. An example of a data set with two covariates (age and gender) can look as follows:

|  | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
|  | Treatment (F) | C | C | C | C | T | T | T | T |
|  | Age | 34 | 23 | 29 | 56 | 40 | 52 | 19 | 41 |
|  | Gender | M | F | F | F | M | M | F | F |
| Variable ID |  |  |  |  |  |  |  |  |  |
| v1 |  | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 |  | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

**Yes** – If there are covariates that potentially affect the dependent variables but are independent of any of the factors which are controlled in the experiment, the effect of the covariates on the dependent variables can be eliminated, which can make the subsequent statistical analysis more powerful. In this case, proceed to Step a**. (a)**

**No** – If there are no covariates, proceed to Question 2. **(2)**

## 4.2 STEP A: ELIMINATE THE EFFECT OF COVARIATES. (a)

In Qlucore Omics Explorer, the effect of a covariate is eliminated with *the Eliminated factors* tool in the Statistics window. Each covariate should be encoded by an annotation in the data file. Depending on the type of the covariate, different types of eliminations are suitable (see the Terminology section). To eliminate the effect of the covariates in a data set, proceed as follows for each covariate:

• Add an eliminated factor.

• Select the corresponding annotation and the proper variable type.

This procedure fits a multiple regression model to all covariates, and subtracts the expression values predicted by this model from the observed values. This means that the covariate elimination follows the same formalism as the general ANOVA, and a covariate can be seen as a factor that does not have any interactions with any other factor in the model. When all covariates are eliminated, proceed to Question 2. **(2)**

## 4.3 QUESTION 2: DO YOU HAVE AN INDEPENDENT OR A REPEATED MEASURES DESIGN? (2)

In an *independent* measures design, the subjects are randomly allocated to different groups, for example receiving different treatments. In a *repeated measures* design, there are multiple measurements for each subject, under different conditions.

**Independent measures** – Proceed to Question 3. **(3)**

**Repeated measures** – Proceed to Question 4. **(4)**

## 4.4 QUESTION 3: HOW MANY FACTORS ARE VARIED IN THE EXPERIMENT? (3)

In a *single factor* experiment, only one factor is varied. An example of this type is a dose-response study, where some subjects receive a low treatment dose while others receive a high dose, and all subjects are studied under otherwise similar conditions. Another example in this category is a case-control study, where subjects with a specific disease are compared to control subjects, following a treatment given to both groups.

In a *multifactor experiment*, several factors are varied simultaneously. For examples, the subjects may be given either a low or a high dose of a drug and furthermore be studied either

one or two weeks after the drug administration. If we vary more than one factor, we can also study interaction effects, for example if the time course of the drug effect is different between the subjects given high or low doses.

**Single factor** – Proceed to Model A. **(A)**

**Multiple factors** – Proceed to Model B. **(B)**

## 4.5 QUESTION 4: HOW MANY FACTORS ARE VARIED IN THE EXPERIMENT? **(4)**

In a *single factor* experiment, only one factor is varied. An example of this type is a dose-response study, where all subjects receive both a low treatment dose and a high treatment dose.

In a *multifactor* experiment, several factors are varied simultaneously. For examples, the subjects may be given either a low or a high dose of a drug and then be studied both one and two weeks after the drug administration. If we vary more than one factor, we can also study interaction effects, for example if the time courses of the drug effect are different when the subjects are given high or low doses, respectively.

**Single factor** – Proceed to Model C. **(C)**

**Multiple factors** – Proceed to Question 5. **(5)**

## 4.6 QUESTION 5: IS THE DESIGN FULLY FACTORIAL OR MIXED? **(5)**

Here, we have an experimental design where more than one factor is varied, and where each subject contributes with multiple measurements. In a fully factorial design, each subject contributes with one measurement for each combination of the levels of the different factors. For example, we could conduct an experiment where each subject receives both a high and a low dose of a drug, and furthermore is studied both one and two weeks after each drug administration. In a mixed design, each subject contributes with measurements for all levels of some factors, but only one level of other factors. As an example of such a study, consider the case where the subjects receive different treatment doses, but where all subjects are studied both one and two weeks after the treatment.

**Fully factorial** – Proceed to Model D. **(D)**

**Mixed design** – Proceed to Model E. **(E)**

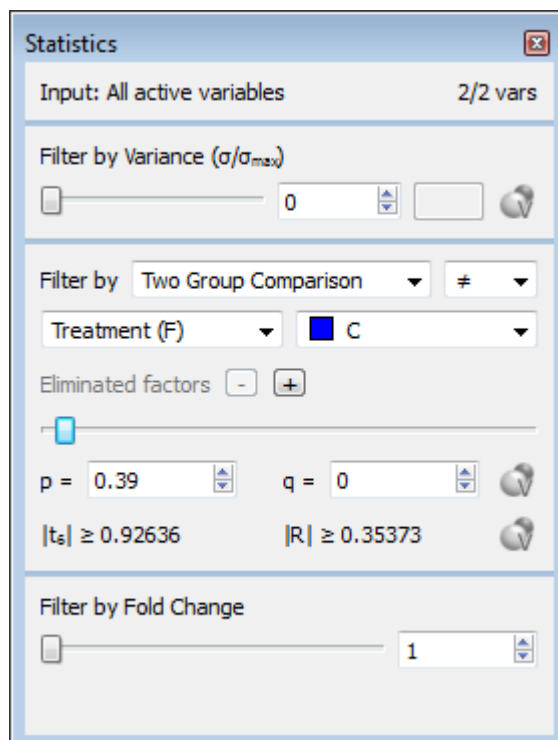## 4.7 MODEL A: INDEPENDENT MEASURES, SINGLE FACTOR DESIGN **(A)**

In this ANOVA model, a single factor (F) is varied in a collection of subjects. It is assumed that the allocation of a subject to a specific level of F is done randomly. The aim is to find those

dependent variables which are related to the level of the factor F. The data file should contain one annotation that defines the value of F for each sample. An example data file could look as follows, where C (control) and T (treated) denote the two levels of the factor F:

| | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| | Treatment (F) | C | C | C | C | T | T | T | T |
| Variable ID | | | | | | | | | |
| v1 | | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 | | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

In Qlucore Omics Explorer, we can find dependent variables related to F in the following way:

• Select a test type and a test direction.

• Identify the relevant annotation, corresponding to the factor F. Select this annotation in the "Selected annotation" combobox.

•Drag the slider to filter the variables by their association to the factor F, or select a desirable significance level (p-value) or false discovery rate (q-value) to select all variables which are significantly associated to F at this level.

**4.8 MODEL B: INDEPENDENT MEASURES, FACTORIAL DESIGN.**

In this ANOVA model, several factors (we will describe the situation where we have two factors, F1 and F2) are varied in a collection of subjects. It is assumed that each subject is randomly allocated to receive one of the levels of each factor. The data file should contain annotations that define the value of each of the factors F1 and F2 in the samples. To study interaction effects, there should also be one annotation corresponding to the combination F1xF2. This annotation should have one level for each combination of the levels of F1 and F2. An example data file could look as follows, where C (control) and T (treated) denote the two levels of F1, and 1 and 2 denote the two levels of F2:
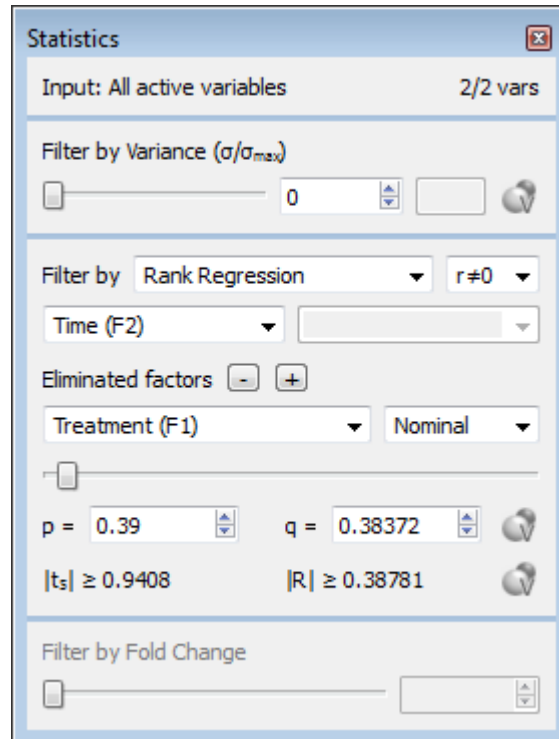
| | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| | Treatment (F1) | C | C | C | C | T | T | T | T |
| | Time (F2) | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 2 |
| | Treatment x Time (F1xF2) | Cx1 | Cx2 | Cx1 | Cx2 | Tx1 | Tx2 | Tx1 | Tx2 |
| Variable ID | | | | | | | | | |
| v1 | | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 | | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

In this situation, we can ask different questions concerning the dependent variables, for example:

• Are there any variables that are related to F2 after eliminating the effect of F1? (This corresponds to performing a two-way ANOVA without interaction effect.)

• Are there any variables that are related to F2 given a specific value of F1? We could for example be interested in knowing if there are any variable that show a time dependence in the treated samples (that is, for samples where F1 has the value T).

• Are there any variables that are related to the interaction between F1 and F2? For example, are there variables for which the time dependence is different between the treated and the untreated samples.
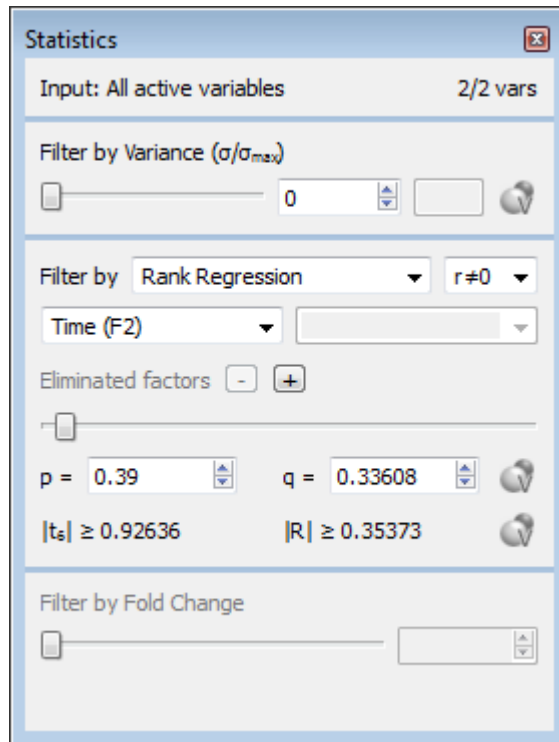
In Qlucore Omics Explorer, to address the first question, do as follows:

• Select a test type and a test direction.

• Choose the annotation corresponding to F2 as the "Selected annotation".

• Eliminate the effect of F1 by adding an eliminated factor and selecting the annotation corresponding to F1.

• Drag the slider to filter the variables by their association to F2 after correcting for F1, or select a desirable significance level (p-value) or false discovery rate (q-value).
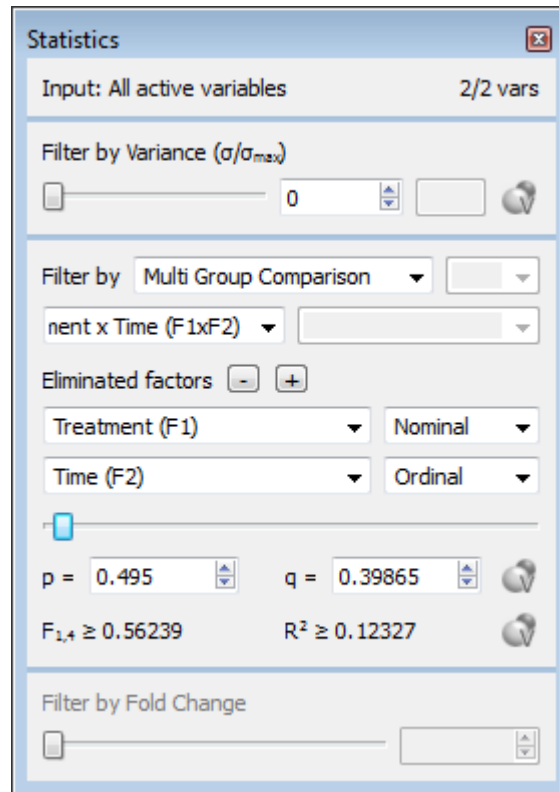
The second question can be addressed as follows:

• Select the active samples to be only those with the specified value of F1 (for example, the ones where F1 has the value T above).

• Then proceed as for the independent measures, single factor design to find variables that are related to F2, given the specified level of F1.

Finally, to find variables that are related to the interaction between F1 and F2, proceed in the following way:

• Select a test type and a test direction.

• Choose the annotation F1xF2 as the "Selected annotation".

• Eliminate the main effects of F1 and F2. For each of the factors, add an eliminated factor and select the corresponding annotation.

• Drag the slider to filter the variables by their association to the interaction between F1 and F2, or select a desirable significance level (p-value) or false discovery rate (q-value).
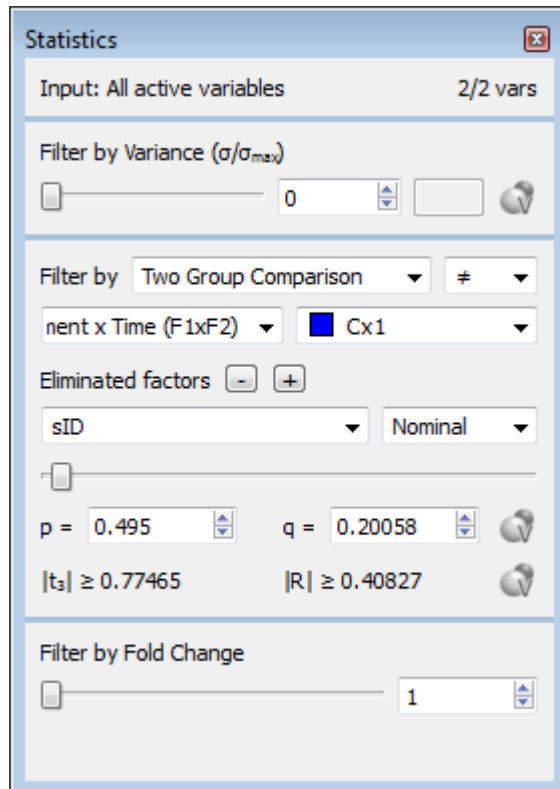
## 4.9 MODEL C: REPEATED MEASURES, SINGLE FACTOR DESIGN. ⓒ

In this ANOVA model, one factor is varied in a collection of samples, and each subject receives all levels of the factor. In this case, by accounting for the subject-specific effects it is possible to obtain a more powerful test than with the independent measures design. This is equivalent to performing a paired test. The data file should contain one annotation (sID) that defines the subject ID, and one annotation that defines the value of the factor (F) for each sample. An example data set can look as follows, where S1,...,S4 denote the four subjects, and C (control) and T (treated) denote the two levels of the factor F. Note that each subject receives both levels of F.

|  | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
|  | Subject (sID) | S1 | S1 | S2 | S2 | S3 | S3 | S4 | S4 |
|  | Treatment (F) | C | T | C | T | C | T | C | T |
| Variable ID |  |  |  |  |  |  |  |  |  |
| v1 |  | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 |  | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

To find variables which are related to the factor F, we could proceed as follows:

• Select a test type and a test direction, depending on the type and the number of levels of F.

• Choose the annotation corresponding to F as the "Selected annotation".

• Add an eliminated factor corresponding to the sID annotation.

• Drag the slider to filter the variables by their association to F, or select a desirable significance level (p-value) or false discovery rate (q-value).



## 4.10 MODEL D: REPEATED MEASURES, FULLY FACTORIAL DESIGN. (D)

In this ANOVA model, several factors (we describe the situation where there are two factors, F1 and F2) are varied in a number of subjects. Each subject provides measurements for all combinations of levels of F1 and F2. The data file should contain one annotation that defines the values for each factor in the samples, and one annotation (sID) corresponding to the subject ID. To study interaction effects, there should also be annotations corresponding to the combinations sIDxF1, sIDxF2 and F1xF2. An example data set can look as follows, where S1 and S2 denote the two subjects, C (control) and T (treated) denote the two levels of F1, and 1 and 2 denote the two levels of F2:

| | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| | Subject (sID) | S1 | S1 | S1 | S1 | S2 | S2 | S2 | S2 |
| | Treatment (F1) | C | C | T | T | C | C | T | T |
| | Time (F2) | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 2 |
| | Subject x Treatment (sIDxF1) | S1xC | S1xC | S1xT | S1xT | S2xC | S2xC | S2xT | S2xT |
| | Subject x Time (sIDxF2) | S1x1 | S1x2 | S1x1 | S1x2 | S2x1 | S2x2 | S2x1 | S2x2 |
| | Treatment x Time (F1xF2) | Cx1 | Cx2 | Tx1 | Tx2 | Cx1 | Cx2 | Tx1 | Tx2 |
| Variable ID | | | | | | | | | |
| v1 | | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 | | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

In this situation, several questions can be asked, for example:

• Are there any variables that are related to F2 given that F1 has a specific value?

• Are there any variables that are related to the interaction between F1 and F2?

To address the first question, select only those samples where F1 has the specified value as the active samples. Then proceed as for the repeated measures single factor design, model C. (C)

The variables that are related to the interaction between F1 and F2 can be found in the following way:

• Select a test type and a test direction, depending on the type and the number of levels of F1xF2.

• Choose the annotation F1xF2 as the "Selected annotation".

• Add eliminated factors corresponding to the annotations sIDxF1 and sIDxF2.

• Drag the slider to filter the variables by their association to F1xF2, or select a desirable significance level (p-value) or false discovery rate (q-value).

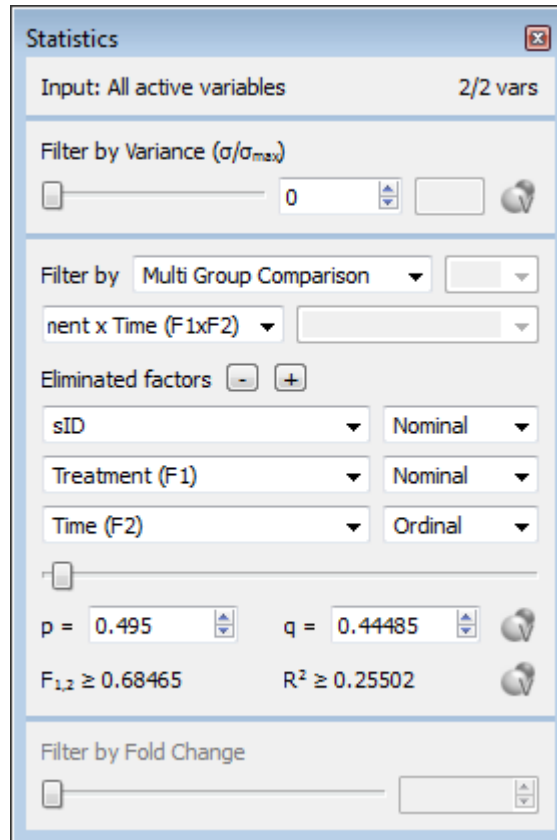## 4.11 MODEL E: REPEATED MEASURES, MIXED DESIGN. (E)

In this ANOVA model, several factors are varied in a number of subjects. Each subject provides measurements for all levels of some factors, but only for one level of other factors. An example data set can look as follows, where S1,...,S4 denote the four subjects, C (control) and T (treated) denote the two levels of F1, and 1 and 2 denote the two levels of F2. Note that each subject provides measurements for only one level of F1 (that is, each subject is assigned either to the control group or to the treated group), but for all levels of F2 (that is, each subject is studied both 1 and 2 weeks after treatment.

| Variable ID | Sample ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| | Subject (sID) | S1 | S1 | S2 | S2 | S3 | S3 | S4 | S4 |
| | Treatment (F1) | C | C | C | C | T | T | T | T |
| | Time (F2) | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 2 |
| | Treatment x Time (F1xF2) | Cx1 | Cx2 | Cx1 | Cx2 | Tx1 | Tx2 | Tx1 | Tx2 |
| v1 | | 0.61 | 0.91 | 0.32 | 0.53 | 0.33 | 0.41 | 0.62 | 0.47 |
| v2 | | 0.37 | 0.37 | 0.43 | 0.41 | 0.08 | 0.05 | 0.35 | 0.59 |

To assess the interaction effect between F1 and F2, proceed as follows:

• Select a test type and a test direction, depending on the type and the number of levels of F1xF2.

• Choose the annotation F1xF2 as the "Selected annotation".

• Add eliminated factors corresponding to the annotations sID, F1 and F2.

• Drag the slider to filter the variables by their association to F1xF2, or select a desirable significance level (p-value) or false discovery rate (q-value).



## 5 EXAMPLE (GDS3521)

To show how the method outlined above works, we apply it to a real data example. This example is based on the data set available from Gene Expression Omnibus (http://www.ncbi.nlm.nih.gov) as data set GDS3521.[1] The data set contains measurements of 45,101 variables in 24 samples. The samples are divided into two groups receiving different treatments (control or hypoxic preconditioning, encoded in the annotation protocol). Furthermore, each sample is studied either 0, 2, 4, or 16 h after treatment (encoded in the annotation time).

_____

[1] See also Thiersch et al (2008): Analysis of the retinal gene expression profile after hypoxic preconditioning identifies candidate genes for neuroprotection. BMC Genomics 9:73.

**5.1 QUESTIONS**

Question 1: Are there any covariates? (1)

In this data set, we do not have any potential confounding factors, such as age or gender, that need to be eliminated. Therefore, the answer is No, and we proceed to Question 2. (2)

Question 2: Independent or repeated measures design? (2)

In this data set, each sample receives only one treatment, and each sample is studied at only one time point. Therefore, we have an independent design, and we proceed to Question 3. (3)

Question 3: How many factors are varied in the experiment? (3)

In this experiment, both the protocol (the treatment) and the time point are varied. Hence, we have **multiple factors,** and we proceed to Model B. (B)

**5.2 MODEL B** (B)

Following the reasoning above, we will use an **independent measures, factorial ANOVA design**. Let F1 correspond to the time annotation (an ordinal variable with four levels), and let F2 correspond to the protocol annotation (a nominal variable with two levels). We will study two aspects of this data set:
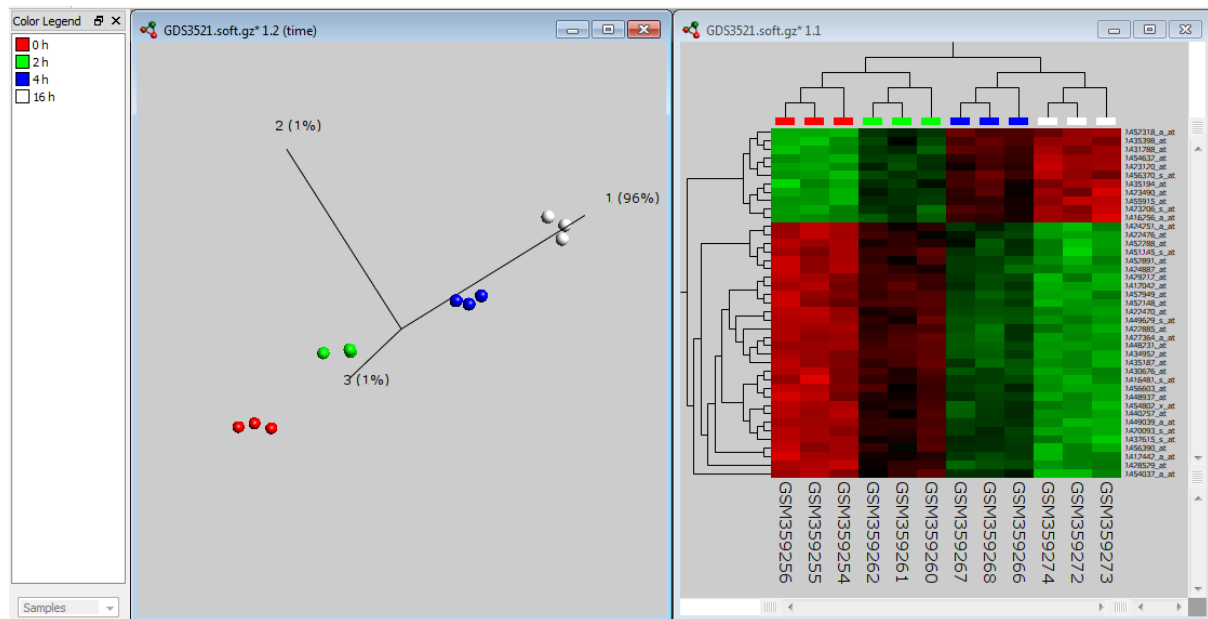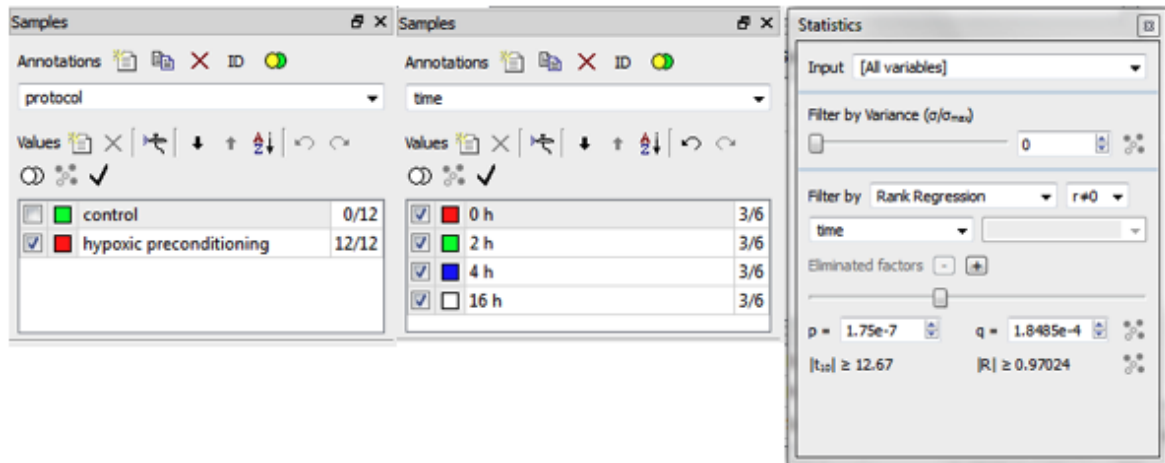
• Are there any variables that are related to the time after treatment for the samples exposed to hypoxic preconditioning? To answer this question, proceed as follows:

o Select the active samples to be only those where protocol has the value hypoxic preconditioning.

o Select Rank Regression as the test type, and $r \neq 0$ as the test direction. Make sure that the values for the annotation time are ordered in increasing order in the Samples window.

o Choose the annotation time as the "Selected annotation".

o Drag the slider to find variables that change with time for the samples exposed to hypoxic preconditioning.

• Are there any variables that are related to the interaction between the protocol and the time after treatment? To address this question, we select all samples as active and proceed as follows.

o We first construct a new annotation (Protocol x Time). Press the New Annotation button to create the annotation, and give it the name Protocol x Time.
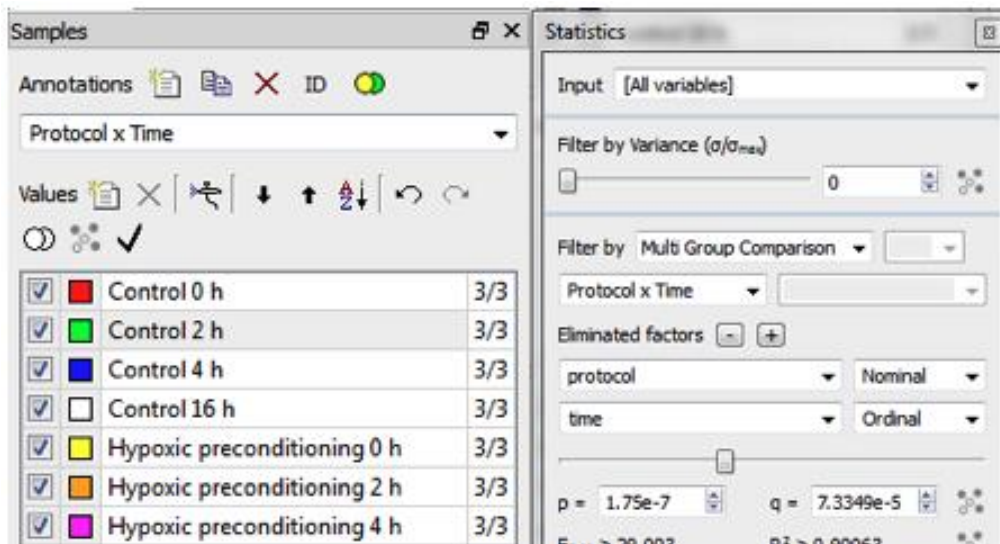
o Repeatedly press the New Value button to create new values for the new annotation. The number of values equals the product of the number of values of protocol and time (that is, 2*4=8 values in this case). Give each value a suitable label.

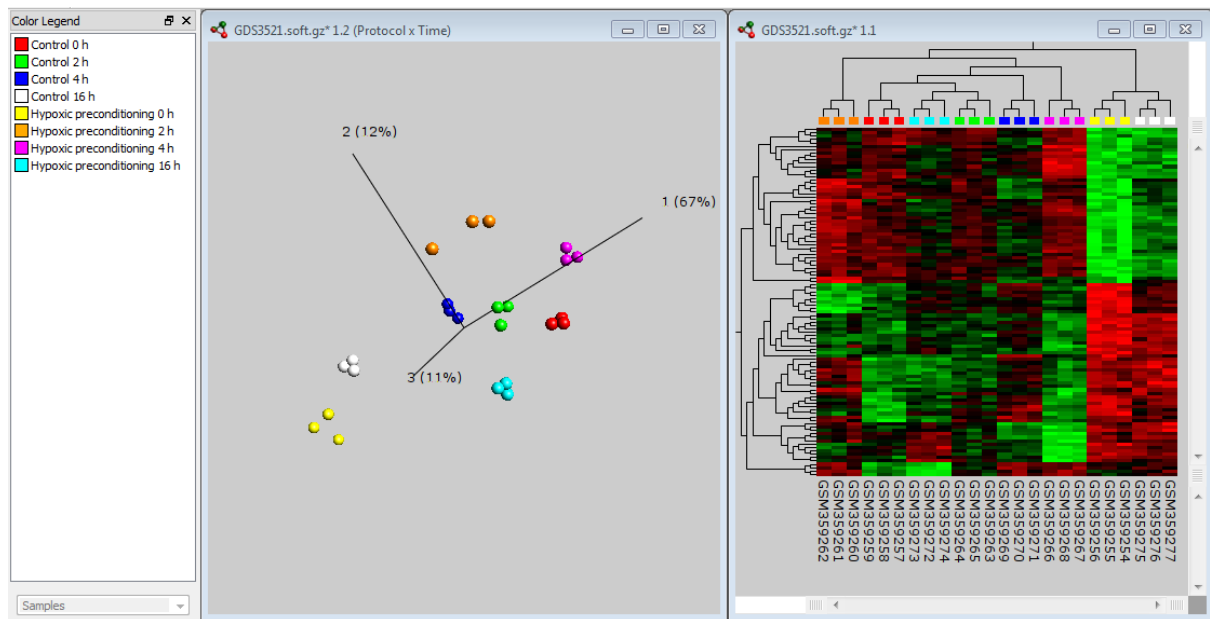o Use the Classify tool to assign each sample the correct value of the new annotation.

o Now, we have the new annotation and we can search for variables that are related to this annotation. The new annotation corresponds to a nominal variable with multiple levels, so we select Multi Group Comparison as the test type. Choose the new annotation as the "Selected annotation".

o Add eliminated factors corresponding to the annotations protocol and time.

o Drag the slider to find variables that depend differently on the time after administration for the different protocols.

**DISCLAIMER**

The contents of this document are subject to revision without notice due to continuous progress in methodology, design, and manufacturing.

Qlucore shall have no liability for any error or damages of any kind resulting from the use of this document.

Qlucore Omics Explorer is only intended for research purposes.