# Sample Pair Mismatch

**Author:** Christophe Lambert and Greta Linse Peterson, Golden Helix, Inc.

## Overview

This script compares genotype calls from NSP and STY files and calculates the correlation between the nearest markers in the two sets. If there is a high correlation, the NSP and STY markers correspond to the same person, otherwise there is a mismatch.

## Recommended Directory Location

Save the script to the following directory:
*..**\Application Data\Golden Helix SVS\UserScripts\Spreadsheet\Numeric\CNV_QA**

**Note:** The **Application Data** folder is a hidden folder on Windows operating systems and its location varies between operating systems. The easiest way to locate this directory on your computer is to open SVS and go to **Tools > Open Folder > UserScripts Folder** and save the script in the **\SVS\Spreadsheet\Numeric\CNV_QA** folder.  If you save the script to the proper folder, it will be accessible from the spreadsheet **Numeric > CNV QA** menu.

## Obtaining the Required Datasets

To compare the NSP and STY genotypes you need the following items:

1. NSP genotypes for all samples with a marker map applied to the spreadsheet.

2. STY genotypes for all samples with a marker map applied to the spreadsheet.

3. A matching spreadsheet with a common sample name as the row labels and NSP file names in the first column and STY file names in the second column. See **Figure 1** for an example of a matching spreadsheet.

**Figure 1: NSP and STY matching spreadsheet.**

Next, the genotypes need to be converted to integers using an additive model.

1) Open the NSP Genotype spreadsheet and convert to integers by going to **Edit > Recode > Recode Genotypes** and select **Encode genotypes numerically based on genetic model: Additive model: DD=2, Dd = 1, dd = 0**.

2) Open the STY Genotype spreadsheet and convert to integers by going to **Edit > Recode > Recode Genotypes** and select **Encode genotypes numerically based on genetic model: Additive model: DD=2, Dd = 1, dd = 0**.

## Using the Script

1) To use the script, open the NSP and STY matching spreadsheet and go to **Numeric > CNV QA > Sample Pair Mismatch.**

2) Choose the columns that contain the NSP and STY CEL names, the numerically-coded NSP and STY spreadsheets, and the maximum distance in base pairs to determine

the markers to compare. The script finds the nearest marker in the STY set less than the threshold for every marker in the NSP set.

3) The NSP and STY correlation will be calculated and output in the NSP STY correlation spreadsheet as a child of the matching spreadsheet. See **Figure 2**.



| | Sample | NSP Name | Sample STY Name | NSP-STY-Cor |
|---|---|---|---|---|
| 1 | CEU_NA06985 | CEU_NA06985_NSP | CEU_NA06985_STY | 0.2935970 |
| 2 | CEU_NA06991 | CEU_NA06991_NSP | CEU_NA06991_STY | 0.2878418 |
| 3 | CEU_NA06993 | CEU_NA06993_NSP | CEU_NA06993_STY | 0.2889659 |
| 4 | CEU_NA06994 | CEU_NA06994_NSP | CEU_NA06994_STY | 0.2966875 |
| 5 | CEU_NA07000 | CEU_NA07000_NSP | CEU_NA07000_STY | 0.2903291 |
| 6 | CEU_NA07019 | CEU_NA07019_NSP | CEU_NA07019_STY | 0.2896972 |
| 7 | CEU_NA07022 | CEU_NA07022_NSP | CEU_NA07022_STY | 0.2935277 |
| 8 | CEU_NA07029 | CEU_NA07029_NSP | CEU_NA07029_STY | 0.2924334 |
| 9 | CEU_NA07034 | CEU_NA07034_NSP | CEU_NA07048_STY | 0.1715511 |
| 10 | CEU_NA07048 | CEU_NA07048_NSP | CEU_NA07034_STY | 0.1671637 |
| 11 | CEU_NA07055 | CEU_NA07055_NSP | CEU_NA07055_STY | 0.299542 |
| 12 | CEU_NA07056 | CEU_NA07056_NSP | CEU_NA07056_STY | 0.2903525 |
| 13 | CEU_NA07345 | CEU_NA07345_NSP | CEU_NA07345_STY | 0.2852376 |
| 14 | CEU_NA07348 | CEU_NA07348_NSP | CEU_NA07348_STY | 0.282452 |
| 15 | CEU_NA07357 | CEU_NA07357_NSP | CEU_NA07357_STY | 0.2961255 |
| 16 | CEU_NA10830 | CEU_NA10830_NSP | CEU_NA10830_STY | 0.2987756 |
| 17 | CEU_NA10831 | CEU_NA10831_NSP | CEU_NA10831_STY | 0.2824123 |
| 18 | CEU_NA10835 | CEU_NA10835_NSP | CEU_NA10835_STY | 0.284601 |
| 19 | CEU_NA10838 | CEU_NA10838_NSP | CEU_NA10838_STY | 0.2855752 |

**Figure 2: NSP STY correlation spreadsheet**

4) Finally, the easiest way to identify pairs that are potential mismatches is to choose **Select > Compare and Activate by Column Agreement** and add the first two columns. Samples that do not match are inactivated. See **Figure 3**. In this case the STY names for CEU_NA07048 and CEU_NA07034 were intentionally switched and you can see that they have the lowest correlation inconsistent with the rest of the correlation values.

5) A histogram is a good visualization tool, but it cannot identify samples that do not fit the distribution. See **Figure 4**.

**Figure 3: NSP STY Correlation sorted by correlation values in ascending order.**
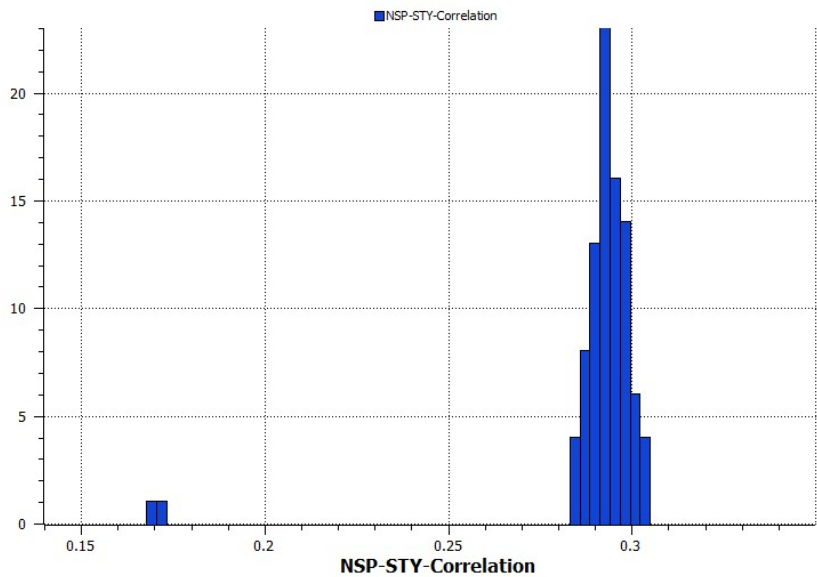


**Figure 4: Histogram of NSP STY Correlation Values showing the two samples that are outliers.**