

# Ovarian Cancer Histotypes: Report of Statistical Findings

Derek Chiu

2020-11-06

# Contents

<b>Preface</b>	<b>4</b>
<b>1 Introduction</b>	<b>5</b>
<b>2 Methods</b>	<b>6</b>
2.1 Data Processing . . . . .	6
2.2 Normalization Between CodeSets . . . . .	10
2.3 Histotype Classification . . . . .	11
<b>3 Validation</b>	<b>12</b>
3.1 Full Data Distributions . . . . .	12
3.2 Training Set Distributions . . . . .	12
3.3 Normalization . . . . .	16
3.4 Common Sample Distributions . . . . .	26
<b>4 Results</b>	<b>28</b>
4.1 CS1 . . . . .	28
4.2 CS2 . . . . .	31

# List of Tables

2.1 Cohort Distribution amongst CodeSets . . . . .	10
3.1 All CodeSet Histotype Groups . . . . .	12
3.2 All CodeSet Histotypes . . . . .	13
3.3 Common Summary ID CodeSet Histotypes . . . . .	13
3.4 All CodeSet Major Histotypes . . . . .	13
3.5 CS1 Histotypes . . . . .	14
3.6 CS2 Histotypes . . . . .	14
3.7 CS3 Histotypes . . . . .	15
3.8 CS1 Training Set Histotypes . . . . .	15
3.9 CS2 Training Set Histotypes . . . . .	15
3.10 Random1 CS1 vs. CS3 Median Concordance Measures by Histotypes . . . . .	23
3.11 Random1 CS2 vs. CS3 Median Concordance Measures by Histotypes . . . . .	23
3.12 Random3 HGSC CS1 vs. CS3 Median Concordance Measures by Histotypes . . . . .	24
3.13 Random3 HGSC CS2 vs. CS3 Median Concordance Measures by Histotypes . . . . .	24
3.14 Pools Non-Normalized CS2 vs. CS3 Median Concordance Measures by Histotypes . . . . .	26
3.15 Pools Normalized CS2 vs. CS3 Median Concordance Measures by Histotypes . . . . .	26
3.16 All Common Samples Histotype Distribution . . . . .	26
3.17 Distinct Common Samples Histotype Distribution . . . . .	27

# List of Figures

3.1	Random3 Non-Normalized Concordance Measure Distributions . . . . .	18
3.2	Random3 Normalized Concordance Measure Distributions . . . . .	19
3.3	Random2 Non-Normalized Concordance Measure Distributions . . . . .	20
3.4	Random2 Normalized Concordance Measure Distributions . . . . .	21
3.5	Random1 Non-Normalized Concordance Measure Distributions . . . . .	22
3.6	Random1 Normalized Concordance Measure Distributions . . . . .	23
3.7	Random3 HGSC Normalized Concordance Measure Distributions . . . . .	24
3.8	CS2Non vs. CS2Pools Concordance Measure Distributions . . . . .	25
3.9	CS2 Non-Normalized Pools vs. CS3 Concordance Measure Distributions . . . . .	25
3.10	CS2 Normalized Pools vs. CS3 Concordance Measure Distributions . . . . .	25
4.1	CS1 Accuracy . . . . .	28
4.2	CS1 F1-Score . . . . .	29
4.3	CS1 Class-Specific F1-Score . . . . .	30
4.4	CS2 Accuracy . . . . .	31
4.5	CS2 F1-Score . . . . .	32
4.6	CS2 Class-Specific F1-Score . . . . .	33

# Preface

This report of statistical findings describes the classification of ovarian cancer histotypes using data from NanoString CodeSets.

Marina Pavanello conducted the initial exploratory data analysis, Cathy Tang implemented class imbalance techniques, Derek Chiu conducted the normalization and statistical analysis, and Aline Talhouk lead the project.

# 1. Introduction

Ovarian cancer has five major histotypes: high-grade serous carcinoma (HGSC), low-grade serous carcinoma (LGSC), endometrioid carcinoma (ENOC), mucinous carcinoma (MUC), and clear cell carcinoma (CCOC). A common problem with classifying these histotypes is that there is a class imbalance issue. HGSC dominates the distribution, commonly accounting for 70% of cases in many patient cohorts, while the other four histotypes are spread over the rest of the cases.

In the NanoString CodeSets, we also run into a problem with trying to find suitable control pools to normalize the gene expression. For prospective NanoString runs, the pools can be specifically chosen, but for retrospective runs, we have to utilize a combination of common samples and common genes as references for normalization.

The supervised learning is performed under a consensus framework: we consider various classification algorithms and use evaluation metrics to help make decisions of which methods to carry forward for downstream analysis.

## 2. Methods

### 2.1 Data Processing

RNA was extracted from FFPE ovarian carcinoma samples and expression was quantified using NanoString nCounter. Samples were run in three CodeSets. Some samples or pools of samples were repeated across CodeSets for expression normalization. Normalizing CS2 to CS3 can easily follow the [PrOType](#) method for HGSC subtypes because both CodeSets have pool samples. A different technique is implemented when normalizing across CS1, CS2, and CS3 where we use common samples and genes as reference sets.

#### 2.1.1 Raw Data

NanoString CodeSets contained a mix of all probes of interest, six positive controls spiked-in at fixed proportional concentrations (0.125- 128 fM), and eight negative controls (probes without a corresponding target). Gene targets also included 5 housekeeping genes: POLR1B, SDHA, PGK1, ACTB, RPL19. Gene selection was made from top ranked differential gene expression analysis between ovarian cancer histotypes and molecular subtypes of HGSC, as well as containing some genes of interest from unrelated projects. Gene targets in each subsequent CodeSet were re-curated, where non-informative genes were dropped and new potential differentiating genes were added.

There are 3 NanoString CodeSets:

- CS1: OvCa2103\_C953
  - Samples = 412
  - Genes = 275
- CS2: PrOTYPE2\_v2\_C1645
  - Samples = 1223
  - Genes = 384
- CS3: OTTA2014\_C2822
  - Samples = 5424
  - Genes = 532

These datasets contain raw counts extracted straight from NanoString RCC files.

### 2.1.2 Housekeeping Genes

The first normalization step is to normalize all endogenous genes to housekeeping genes (POLR1B, SDHA, PGK1, ACTB, RPL19; reference genes expressed in all cells). We normalize by subtracting the average  $\log_2$  housekeeping gene expression from the  $\log_2$  endogenous gene expression:

$$\log_2 \text{ endogenous expression} - \log_2 \text{ average housekeeping expression} = \text{relative expression}$$

The updated CodeSet dimensions are now:

- CS1: OvCa2103\_C953
  - Samples = 412
  - Genes = 256
- CS2: PrOTYPE2\_v2\_C1645
  - Samples = 1223
  - Genes = 365
- CS3: OTTA2014\_C2822
  - Samples = 5424
  - Genes = 513

The number of genes are reduced by 19: 5 housekeeping, 8 negative, 6 positive (the latter 2 types are not used).

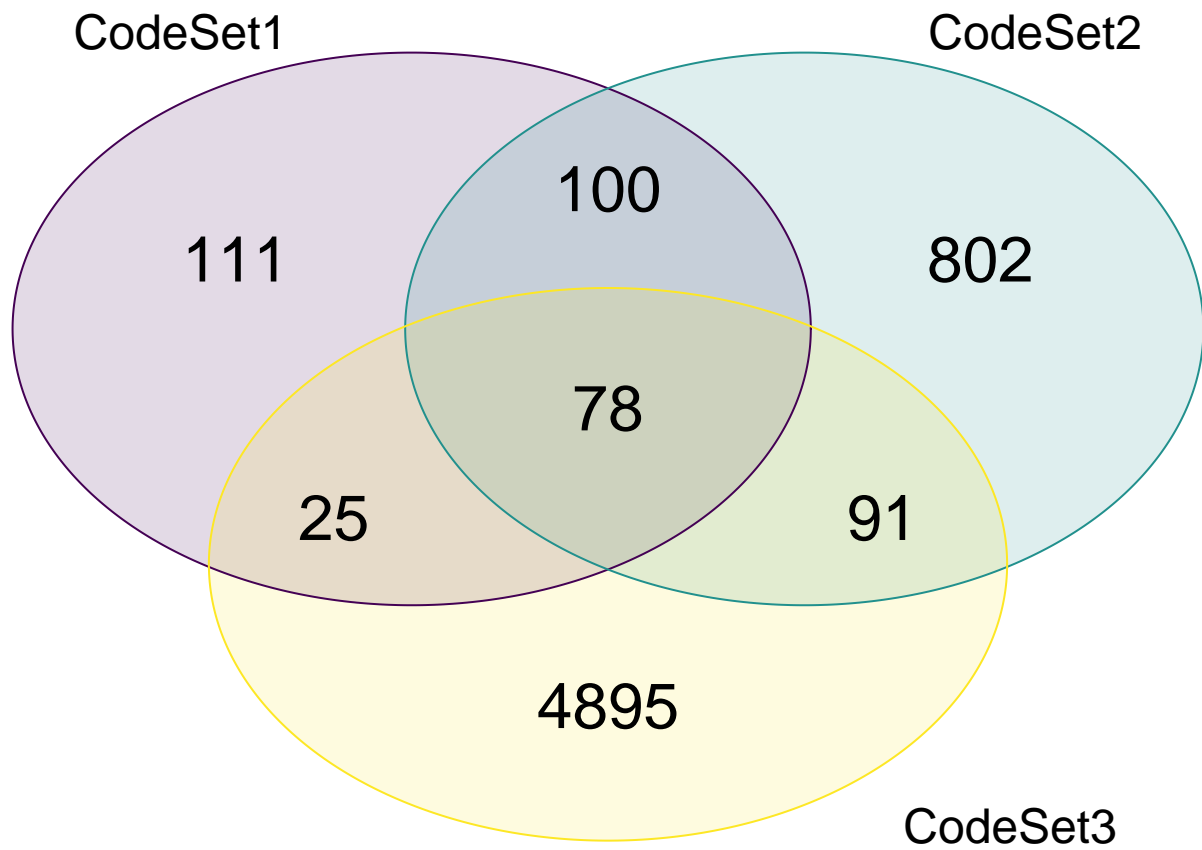
### 2.1.3 Common Samples and Genes

Since the reference pool samples only exist in CS2 and CS3, we need to find an alternative method to normalize all three CodeSets. One method is to select common samples and common genes that exist in all three. We found 72 common genes. Using the `summaryID` identifier, we also found 78 common summary IDs, translating to 320 samples. The number of samples that were matched to each CodeSet differed:

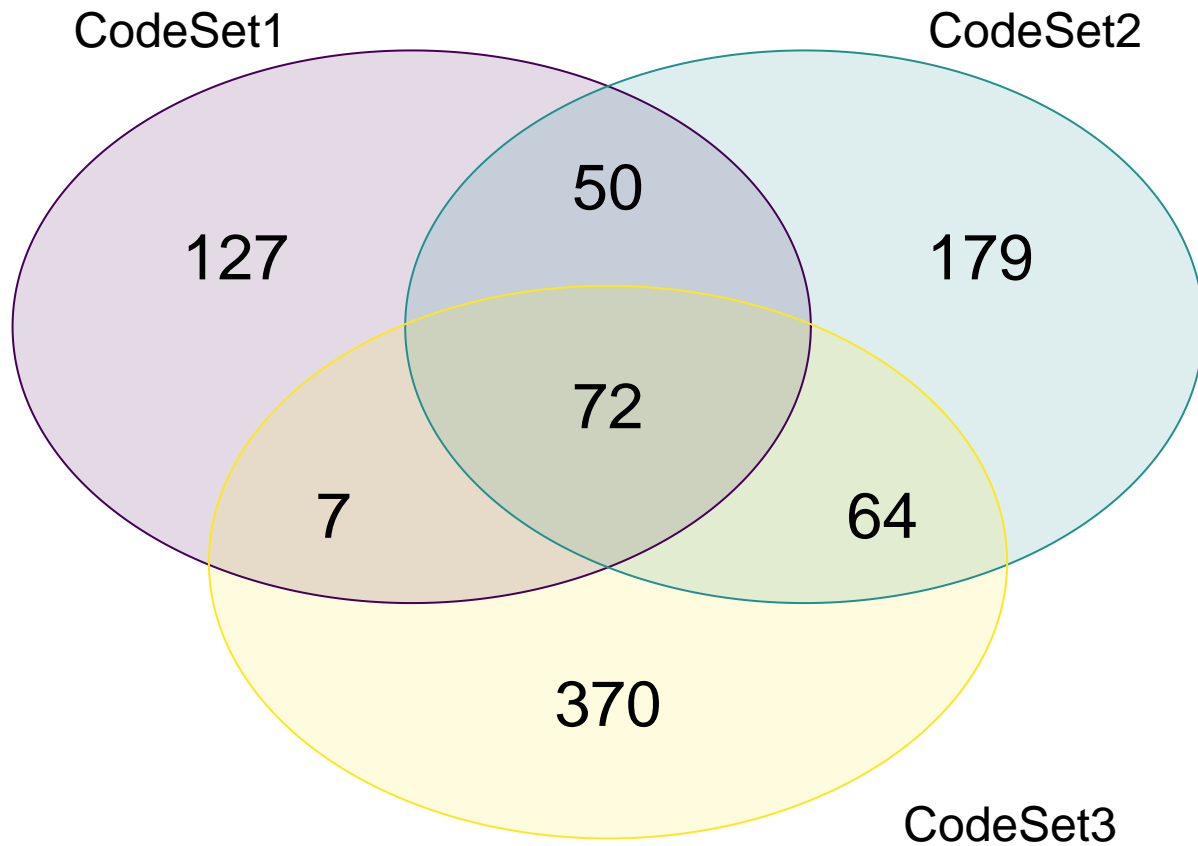
- CS1: OvCa2103\_C953
  - Samples = 93
  - Genes = 72
- CS2: PrOTYPE2\_v2\_C1645
  - Samples = 87
  - Genes = 72
- CS3: OTTA2014\_C2822
  - Samples = 140
  - Genes = 72



2.1.3.1 Overlap of common samples by summary ID



#### 2.1.3.2 Overlap of common genes



\*Excluding housekeeping genes and controls

#### 2.1.4 CS1 Training Set Generation

We use the reference method to normalize CS1 to CS3.

- CS1 reference set: duplicate samples from CS1
  - Samples = 25
  - Genes = 72
- CS3 reference set: corresponding samples in CS3 also found in CS1 reference set
  - Samples = 20
  - Genes = 72
- CS1 validation set: remaining CS1 samples with reference set removed
  - Samples = 387
  - Genes = 72

The final CS1 training set has 304 samples on 72 genes after normalization and keeping only the major histotypes of interest.

Table 2.1: Cohort Distribution amongst CodeSets

cohort	cs1	cs2	cs3
MAYO	6	63	NA
MTL	3	59	NA
OOU	108	43	19
OOUE	32	30	11
VOA	145	122	538
ICON7	NA	416	NA
JAPAN	NA	8	NA
OVAR3	NA	150	NA
POOL-CTRL	NA	12	NA
DOVE4	NA	NA	1160
POOL-1	NA	NA	31
POOL-2	NA	NA	14
POOL-3	NA	NA	13
TNCO	NA	NA	691

### 2.1.5 CS2 Training Set Generation

We use the pool method to normalize CS2 to CS3 so we can be consistent with the PrOType normalization when there are available pools.

- CS2 pools:
  - Samples = 12 (Pool 1 = 4, Pool 2 = 4, Pool 3 = 4)
  - Genes = 365
- CS3 pools:
  - Samples = 22 (Pool 1 = 12, Pool 2 = 5, Pool 3 = 5)
  - Genes = 513
- CS2 validation set: CS2 samples with pools removed
  - Samples = 1214
  - Genes = 365

The final CS2 training set has 945 samples on 136 (common) genes after normalization and keeping only the major histotypes of interest.

### 2.1.6 Cohort Distribution

CodeSets comprised samples from sites collected internationally as shown below. Note that the CS3 pools sample total (n=58) shown here include those that are not used as reference pools, following previous normalization methods. In particular, the distribution of CS3 pools actually used for normalization (n=22) is POOL1 = 12, POOL2 = 5, POOL3 = 5.

## 2.2 Normalization Between CodeSets

After normalization to housekeeping genes and filtering for the five major histotypes of interest, as determined by pathology review and/or IHC, two methods were used to normalize data between CodeSets.

### 2.2.1 Common Samples Method

The common samples method was used to normalize CodeSet1, 2, and 3, where common samples and genes were used as reference sets. Among the samples repeated in all CodeSets we normalized using either: a random set of 3 samples from each major histotype (random3; n=15), a random set of 2 samples from each major histotype (random2; n=10), or a random set of 1 sample from each major histotype (random1; n=5). In each case CodeSet3 expression ( $X_3$ ) was held fixed, while CodeSet1/2 expression ( $X_1$  and  $X_2$ ) were normalized to CodeSet3 by subtracting the average gene expression from the CodeSet1/2 reference set ( $R_1$  or  $R_2$ ) and adding the average gene expression of the CodeSet3 reference set ( $R_3$ ). Alternatively,  $X_1$  (norm) =  $X_1 - R_1 + R_3$  would calibrate CodeSet1 to CodeSet3.

### 2.2.2 Pools Method

The pools method was used to normalize CodeSet2 and CodeSet3. The three reference pools, regularly assayed mixes of samples representing all histotypes, were run in CodeSet2 and CodeSet3 only. CodeSet2 contained 12 reference pool samples (Pool 1 = 4, Pool 2 = 4, Pool 3 = 4) and CodeSet3 contained 22 reference pool samples (Pool 1 = 12, Pool 2 = 5, Pool 3 = 5). Similar to the common samples method, CodeSet2 was normalized to CodeSet3 via:  $X_2$  (norm) =  $X_2 - R_2 + R_3$  where  $R$  is the average expression of the reference pool samples in the respective CodeSet. This method of pool normalization was also used by PrOType to classify HGSC subtypes

### 2.2.3 Concordance Comparison

Concordance between CodeSets using the different normalization strategies was compared in common samples, excluding those used for the normalization, using Pearson's correlation coefficient ( $R^2$ ), coefficient of accuracy (Ca), and Lin's concordance correlation ( $R_c = R^2 \times Ca$ ).

## 2.3 Histotype Classification

We use 6 classification algorithms and 4 subsampling methods across 500 repetitions in the supervised learning framework for CS1 and CS2. The pipeline was run using many SGE batch jobs as a way of parallelization on a CentOS 5 server. Implementations of the techniques below were called from the [splendid](#) package.

- Classifiers:
  - Random Forest
  - Adaboost
  - XGBoost
  - LDA
  - SVM
  - K-Nearest Neighbours
- Subsampling:
  - None
  - Down-sampling
  - Up-sampling
  - SMOTE

## 3. Validation

### 3.1 Full Data Distributions

The histotype distributions on the full data are shown below.

### 3.2 Training Set Distributions

The training set distributions for CS1 and CS2 are shown below.

Table 3.1: All CodeSet Histotype Groups

hist_gr	CS1	CS2	CS3
HGSC	169	757	2453
non-HGSC	196	377	677

Table 3.2: All CodeSet Histotypes

revHist	CS1	CS2	CS3
CARCINOMA-NOS	0	61	23
Carcinoma, NOS	0	0	2
CCOC	57	68	182
CCOC-MCT	0	1	0
Cell-Line	17	48	13
CTRL	0	12	0
ENOC	61	30	272
ENOC-CCOC	0	7	0
ERROR	0	3	0
HGSC	169	757	2453
HGSC-MCT	0	1	0
LGSC	22	29	50
MBOT	0	20	3
MET-NOP	0	21	0
MIXED (ENOC/CCOC)	0	0	1
MIXED (ENOC/LGSC)	0	0	1
MIXED (HGSC/CCOC)	0	0	1
mixed cell	0	0	7
MMMT	0	0	30
MUC	20	61	77
Other (use when 6, 7, or 9 is not distinguished) or unknown if epithelial	0	0	1
Other/Exclude	0	0	8
SBOT	19	10	3
Serous	0	0	2
serous LMP	0	0	1
SQAMOUS	0	1	0
UNK	0	4	0

Table 3.3: Common Summary ID CodeSet Histotypes

revHist	CS1	CS2	CS3
CCOC	3	4	9
Cell-Line	4	5	5
ENOC	4	4	9
HGSC	68	64	98
LGSC	7	5	8
MUC	7	5	11

Table 3.4: All CodeSet Major Histotypes

revHist	CS1	CS2	CS3	CS1_percent	CS2_percent	CS3_percent
CCOC	57	68	182	17.3	7.2	6.0
ENOC	61	30	272	18.5	3.2	9.0
HGSC	169	757	2453	51.4	80.1	80.9
LGSC	22	29	50	6.7	3.1	1.6
MUC	20	61	77	6.1	6.5	2.5

Table 3.5: CS1 Histotypes

CodeSet	revHist	n
CS1	CCOC	57
CS1	Cell-Line	17
CS1	ENOC	61
CS1	HGSC	169
CS1	LGSC	22
CS1	MUC	20
CS1	SBOT	19

Table 3.6: CS2 Histotypes

CodeSet	revHist	n
CS2	CARCINOMA-NOS	61
CS2	CCOC	68
CS2	CCOC-MCT	1
CS2	Cell-Line	48
CS2	CTRL	12
CS2	ENOC	30
CS2	ENOC-CCOC	7
CS2	ERROR	3
CS2	HGSC	757
CS2	HGSC-MCT	1
CS2	LGSC	29
CS2	MBOT	20
CS2	MET-NOP	21
CS2	MUC	61
CS2	SBOT	10
CS2	SQAMOUS	1
CS2	UNK	4

Table 3.7: CS3 Histotypes

CodeSet	revHist	n
CS3	CARCINOMA-NOS	23
CS3	Carcinoma, NOS	2
CS3	CCOC	182
CS3	Cell-Line	13
CS3	ENOC	272
CS3	HGSC	2453
CS3	LGSC	50
CS3	MBOT	3
CS3	MIXED (ENOC/CCOC)	1
CS3	MIXED (ENOC/LGSC)	1
CS3	MIXED (HGSC/CCOC)	1
CS3	mixed cell	7
CS3	MMMT	30
CS3	MUC	77
CS3	Other (use when 6, 7, or 9 is not distinguished) or unknown if epithelial	1
CS3	Other/Exclude	8
CS3	SBOT	3
CS3	Serous	2
CS3	serous LMP	1

Table 3.8: CS1 Training Set Histotypes

histotype	n
CCC	57
ENOCa	59
HGSC	156
LGSC	16
MUC	16

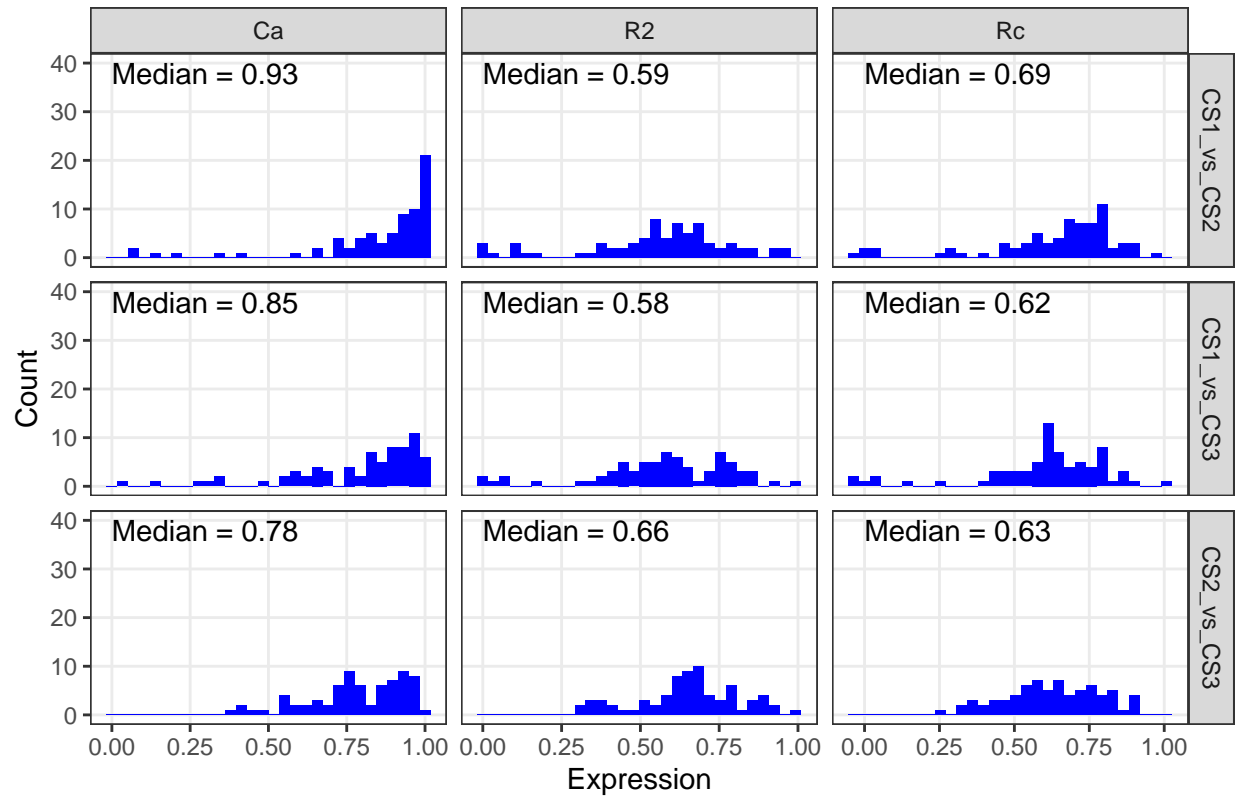
Table 3.9: CS2 Training Set Histotypes

histotype	n
CCOC	68
ENOC	30
HGSC	757
LGSC	29
MUC	61

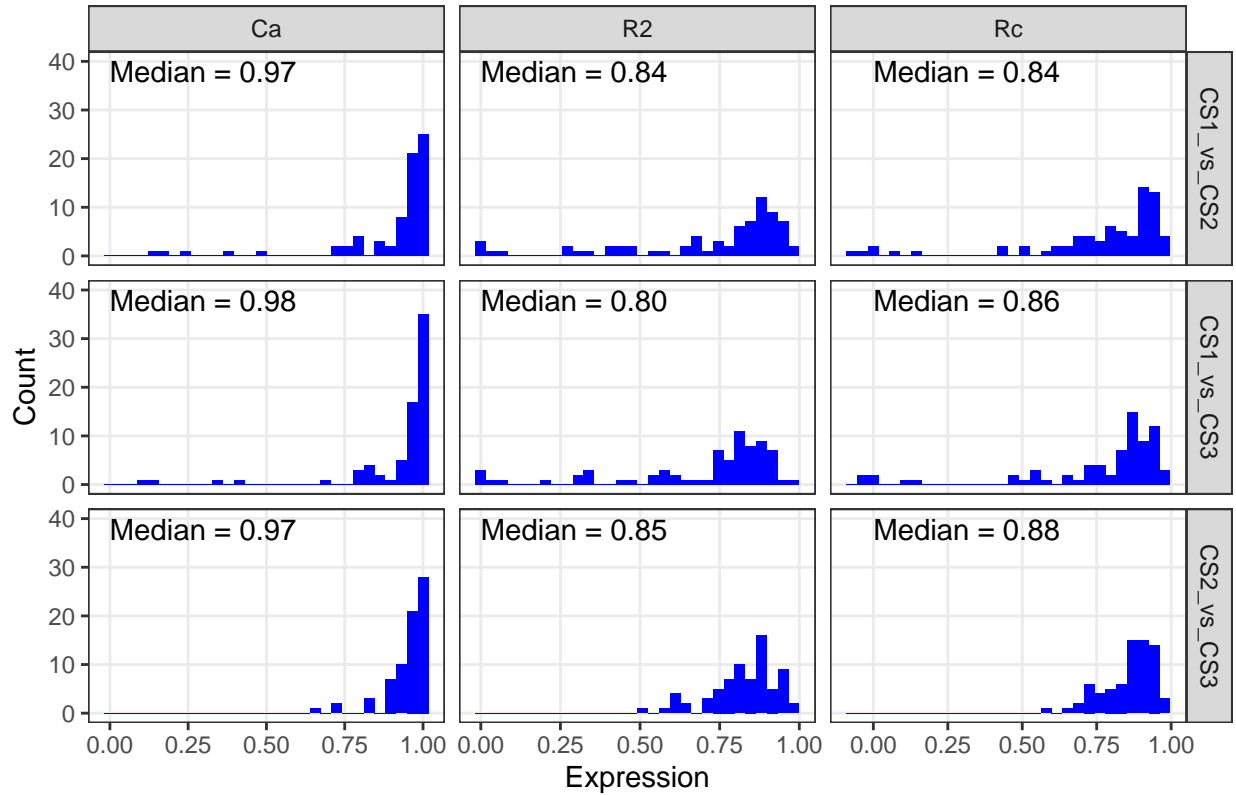


### 3.3 Normalization

Raw Non-Normalized Concordance Measure Distributions



## HK genes Normalized Concordance Measure Distributions



### 3.3.1 Common Samples Method

We employ a new normalization technique using randomly selected samples common to all three CodeSets with a uniform distribution of histotypes as the reference dataset. The number of randomly selected samples ranges from 1-3 per histotype. Hence, the reference dataset has either 5, 10, or 15 samples and we validate on the remaining. Note that ottaID duplicates are collapsed by mean averaging the gene expression.  $n=72$  common samples.

CodeSets 1 and 2 are calibrated to CodeSet3 as follows:

$$X^1(\text{norm}) = X^1 - R^1 + R^3$$

$$X^2(\text{norm}) = X^2 - R^2 + R^3$$

$$X^3(\text{norm}) = X^3$$

### Random3 Non-Normalized Concordance Measure Distributions

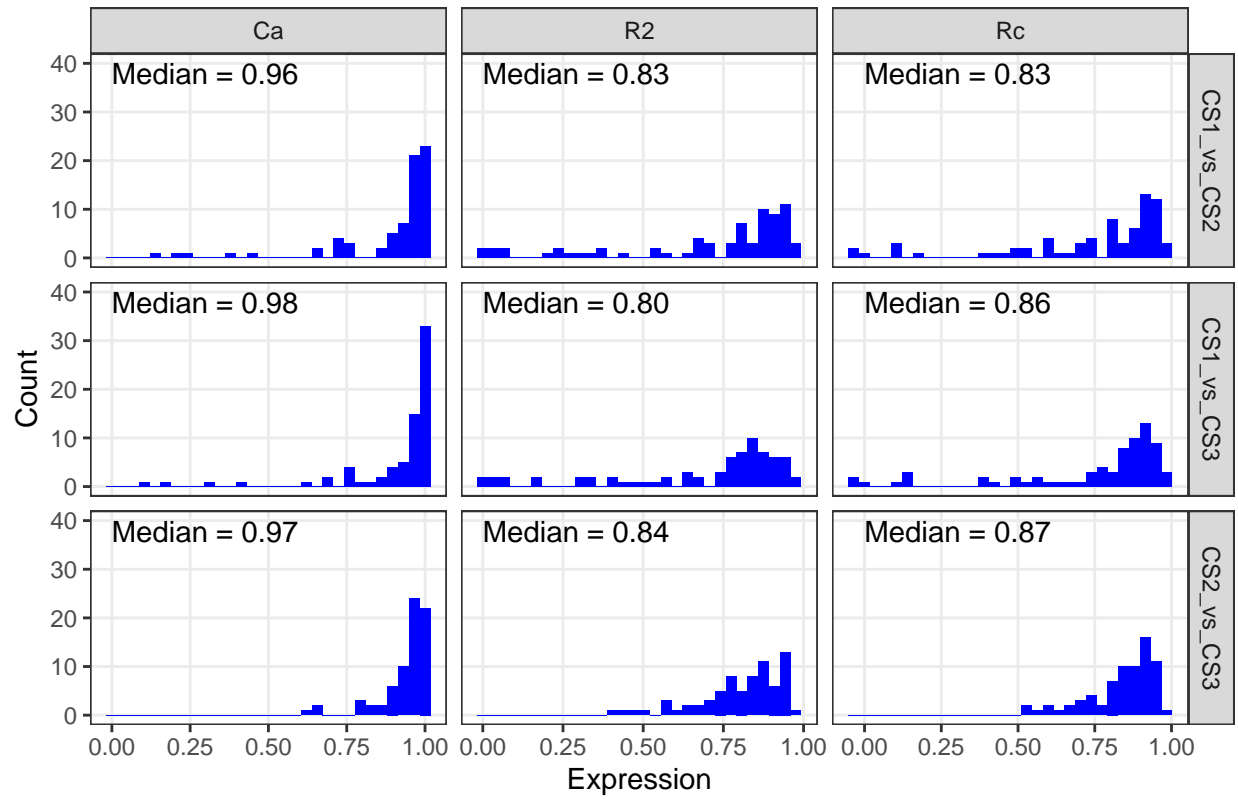


Figure 3.1: Random3 Non-Normalized Concordance Measure Distributions

### Random3 Normalized Concordance Measure Distributions

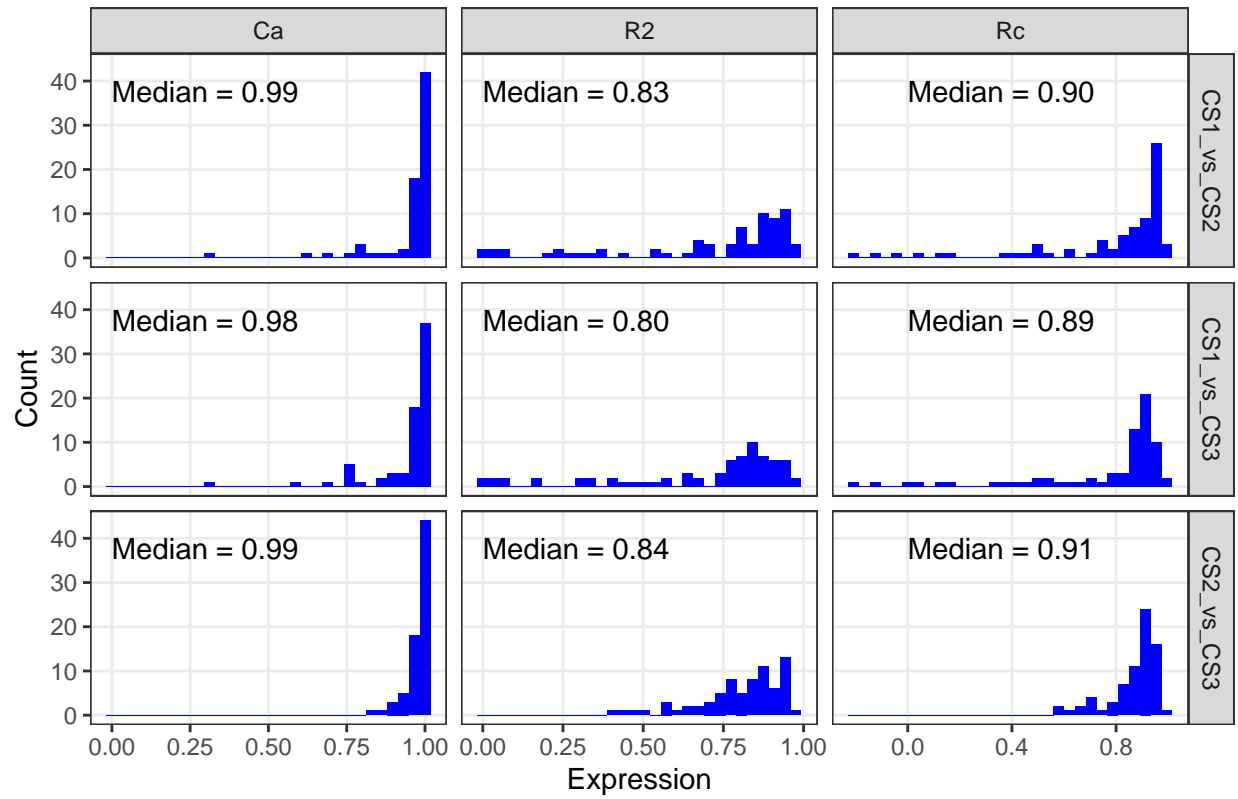


Figure 3.2: Random3 Normalized Concordance Measure Distributions

### Random2 Non-Normalized Concordance Measure Distributions

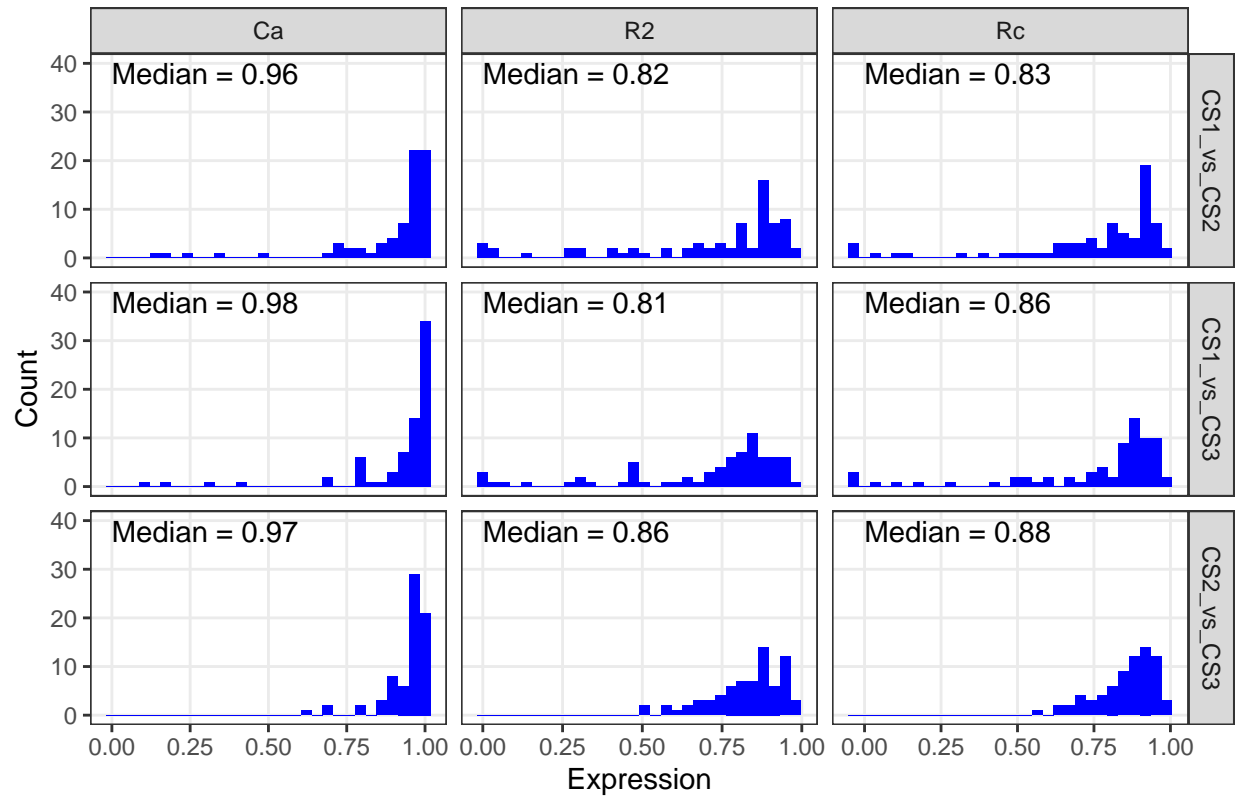


Figure 3.3: Random2 Non-Normalized Concordance Measure Distributions

### Random2 Normalized Concordance Measure Distributions

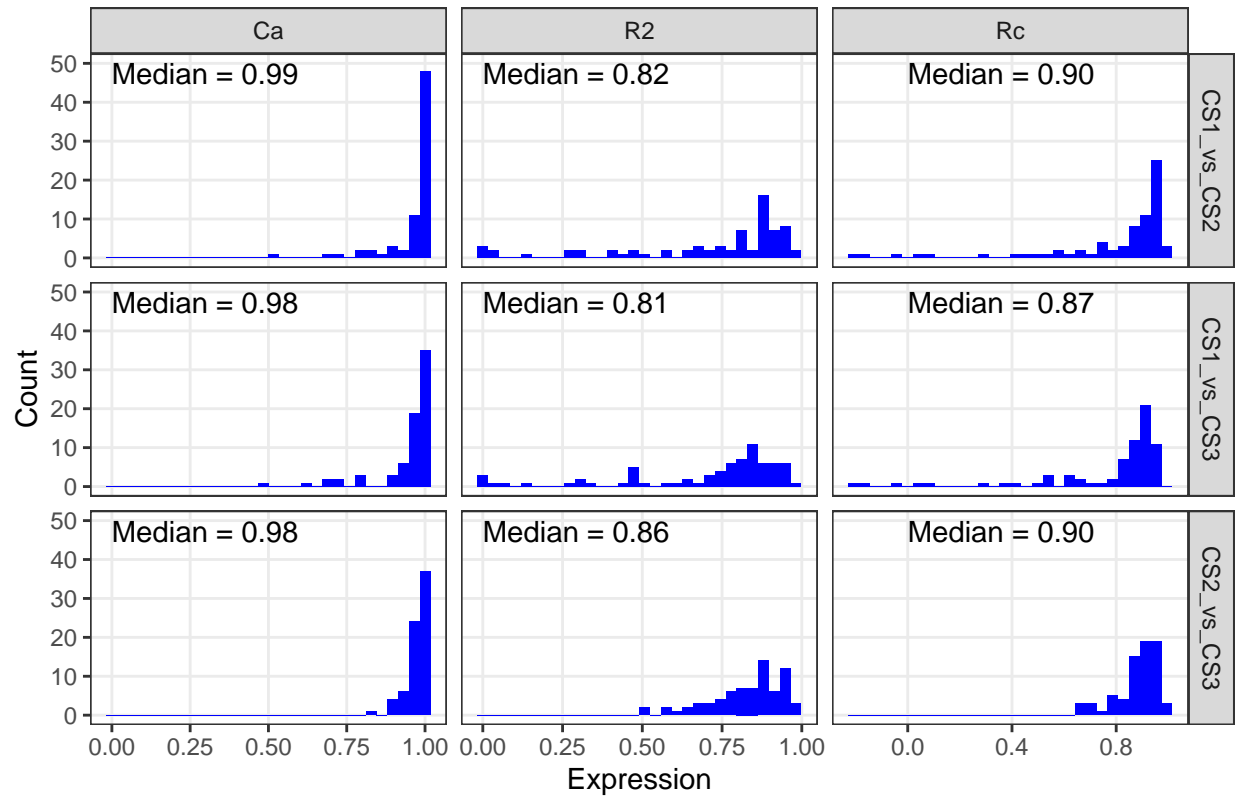


Figure 3.4: Random2 Normalized Concordance Measure Distributions

### Random1 Non-Normalized Concordance Measure Distributions

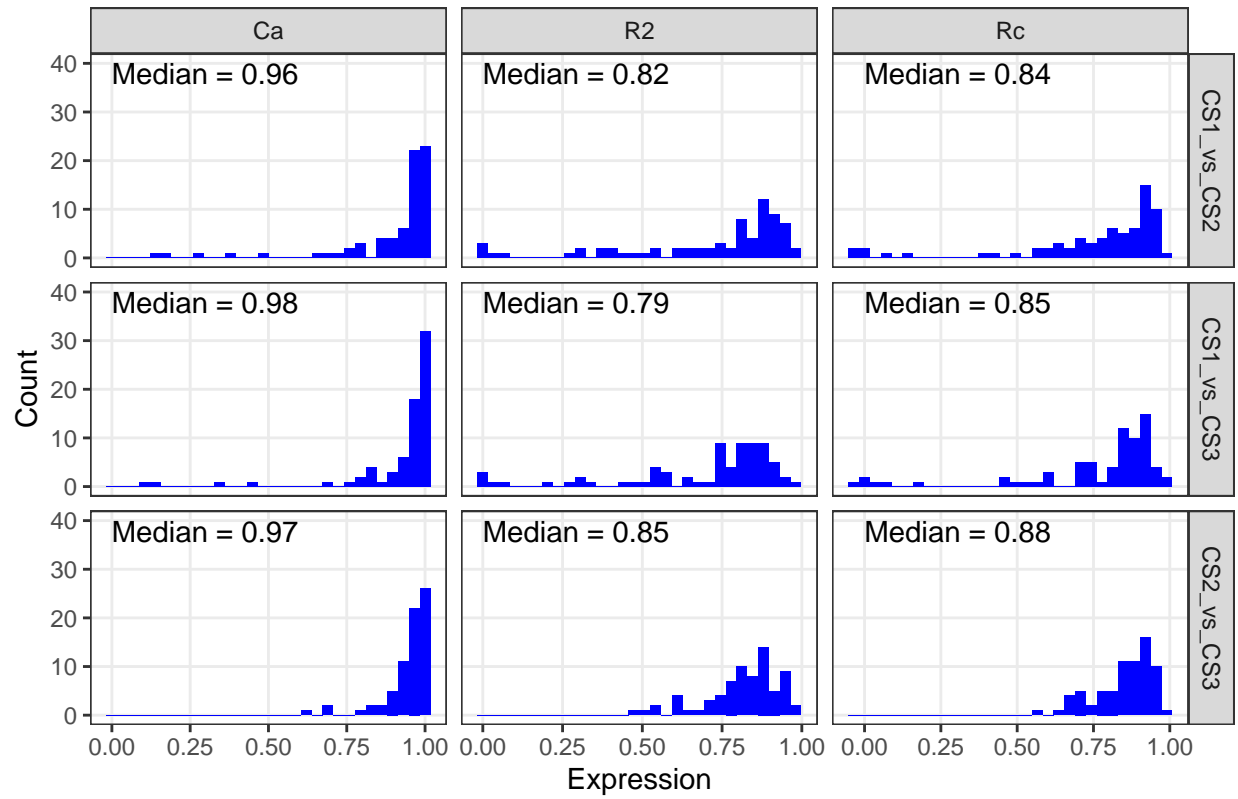


Figure 3.5: Random1 Non-Normalized Concordance Measure Distributions

Table 3.10: Random1 CS1 vs. CS3 Median Concordance Measures by Histotypes

hist	R2-Non	Ca-Non	Rc-Non	R2-Norm	Ca-Norm	Rc-Norm
CCOC	1.00	0.30	0.08	1.00	0.28	0.09
ENOC	1.00	0.54	0.54	1.00	0.61	0.61
HGSC	0.78	0.98	0.85	0.78	0.97	0.86
LGSC	0.97	0.87	0.81	0.97	0.90	0.87
MUC	0.75	0.85	0.68	0.75	0.82	0.64

Table 3.11: Random1 CS2 vs. CS3 Median Concordance Measures by Histotypes

hist	R2-Non	Ca-Non	Rc-Non	R2-Norm	Ca-Norm	Rc-Norm
CCOC	1.00	0.20	0.04	1.00	0.27	0.09
ENOC	1.00	0.67	0.64	1.00	0.64	0.61
HGSC	0.83	0.96	0.86	0.83	0.99	0.89
LGSC	0.99	0.92	0.91	0.99	0.96	0.94
MUC	0.70	0.78	0.55	0.70	0.86	0.57

Random1 Normalized Concordance Measure Distributions

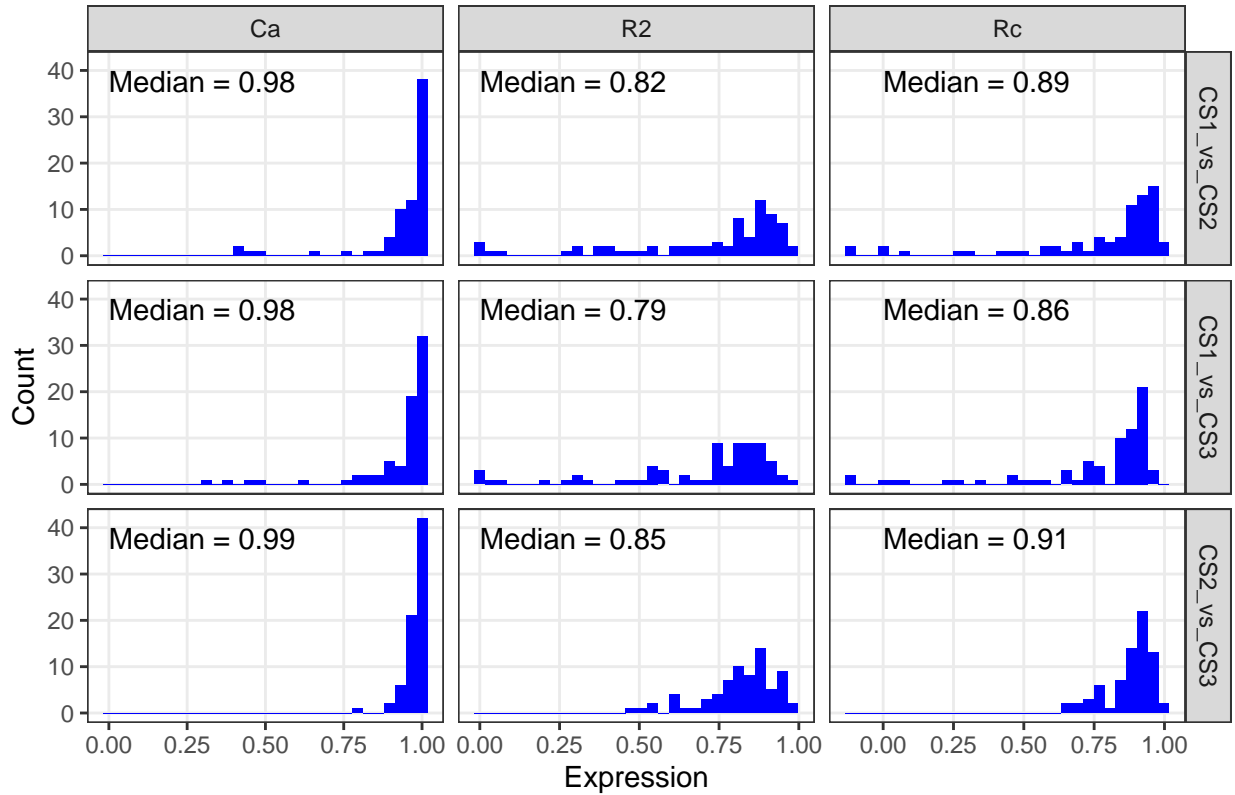


Figure 3.6: Random1 Normalized Concordance Measure Distributions



Table 3.12: Random3 HGSC CS1 vs. CS3 Median Concordance Measures by Histotypes

hist	R2-Non	Ca-Non	Rc-Non	R2-Norm	Ca-Norm	Rc-Norm
CCOC	0.67	0.58	0.29	0.67	0.63	0.26
ENOC	0.89	0.80	0.71	0.89	0.80	0.74
HGSC	0.77	0.98	0.85	0.77	0.99	0.87
LGSC	0.94	0.85	0.79	0.94	0.88	0.83
MUC	0.76	0.93	0.73	0.76	0.93	0.79

Table 3.13: Random3 HGSC CS2 vs. CS3 Median Concordance Measures by Histotypes

hist	R2-Non	Ca-Non	Rc-Non	R2-Norm	Ca-Norm	Rc-Norm
CCOC	0.69	0.54	0.35	0.69	0.64	0.39
ENOC	0.85	0.77	0.66	0.85	0.86	0.76
HGSC	0.82	0.96	0.86	0.82	0.99	0.90
LGSC	0.98	0.95	0.92	0.98	0.94	0.92
MUC	0.77	0.89	0.73	0.77	0.92	0.70

Random3 HGSC Normalized Concordance Measure Distributions

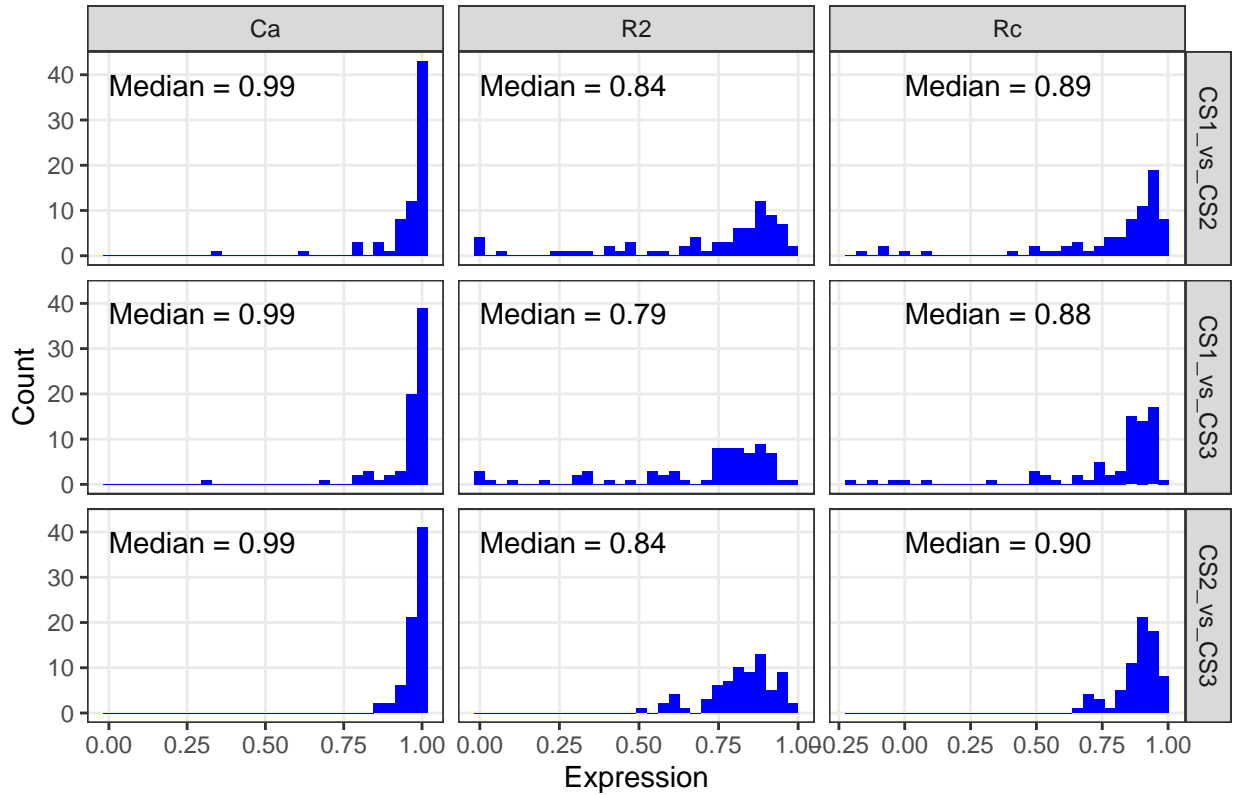


Figure 3.7: Random3 HGSC Normalized Concordance Measure Distributions

### 3.3.2 Pools Method

CodeSet2 contains 12 ref pool samples (Pool 1 = 4, Pool 2 = 4, Pool 3 = 4). CodeSet3 contains 22 ref pool samples (Pool 1 = 12, Pool 2 = 5, Pool 3 = 5). n=86 common samples.

CodeSet2 is calibrated to CodeSet3 as follows:

$$X^2(\text{norm}) = X^2 - R^2 + R^3$$

$$X^3(\text{norm}) = X^3$$

#### CS2Non vs. CS2Pools Concordance Measure Distributions

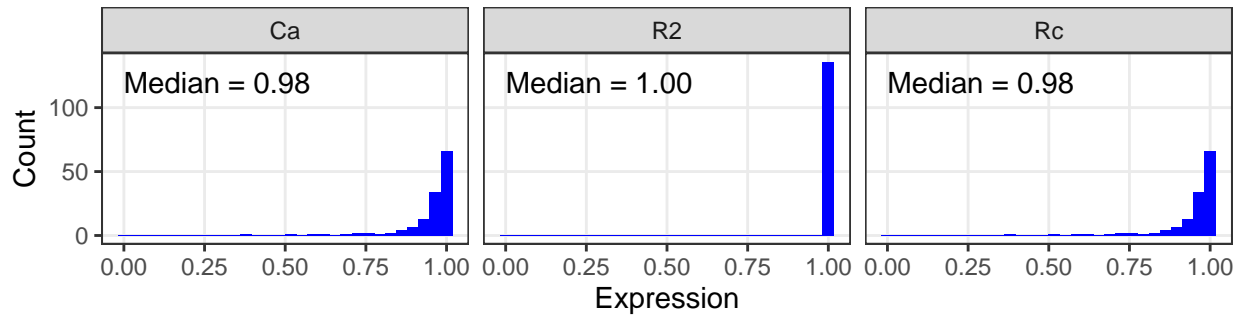


Figure 3.8: CS2Non vs. CS2Pools Concordance Measure Distributions

#### CS2 Non-Normalized Pools vs. CS3 Concordance Measure Distributions

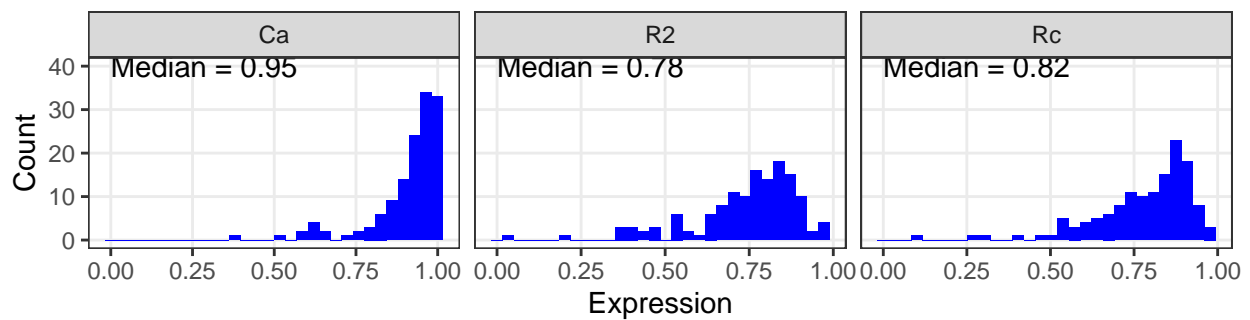


Figure 3.9: CS2 Non-Normalized Pools vs. CS3 Concordance Measure Distributions

#### CS2 Normalized Pools vs. CS3 Concordance Measure Distributions

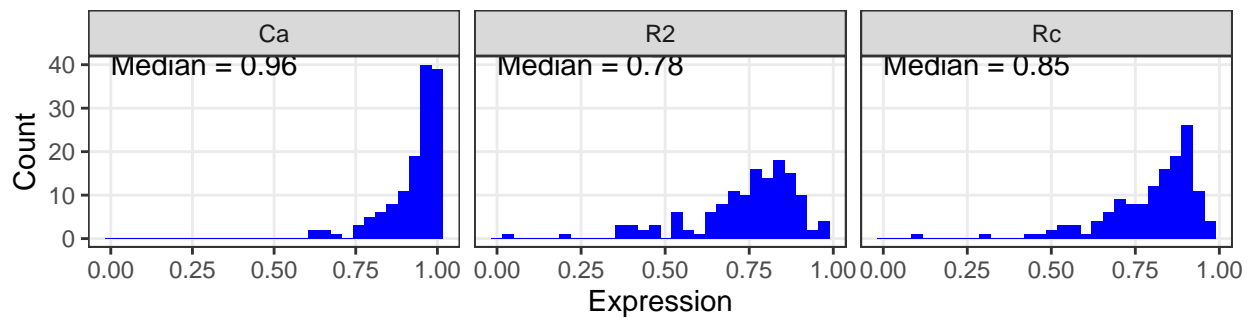


Figure 3.10: CS2 Normalized Pools vs. CS3 Concordance Measure Distributions

Table 3.14: Pools Non-Normalized CS2 vs. CS3 Median Concordance Measures by Histotypes

hist	R2	Ca	Rc
CCOC	0.66	0.51	0.29
ENOC	0.87	0.74	0.64
HGSC	0.77	0.94	0.80
LGSC	0.98	0.94	0.92
MUC	0.75	0.86	0.69

Table 3.15: Pools Normalized CS2 vs. CS3 Median Concordance Measures by Histotypes

hist	R2	Ca	Rc
CCOC	0.66	0.62	0.31
ENOC	0.87	0.74	0.67
HGSC	0.77	0.94	0.82
LGSC	0.98	0.96	0.93
MUC	0.75	0.91	0.70

### 3.4 Common Sample Distributions

Table 3.16: All Common Samples Histotype Distribution

revHist	CS1	CS2	CS3
CCOC	3	4	9
ENOC	4	4	9
HGSC	59	62	95
LGSC	7	5	8
MUC	7	5	11

Table 3.17: Distinct Common Samples Histotype Distribution

revHist	CS1	CS2	CS3
CCOC	3	3	3
ENOC	3	3	3
HGSC	57	57	57
LGSC	4	4	4
MUC	5	5	5

## 4. Results

Here we show internal validation summaries for both CS1 and CS2. The accuracy and F1-scores are the measures of interest. Algorithms are sorted by descending value. The point ranges show the median, 5th and 95th percentiles, coloured by subsampling methods.

### 4.1 CS1

#### 4.1.1 Accuracy

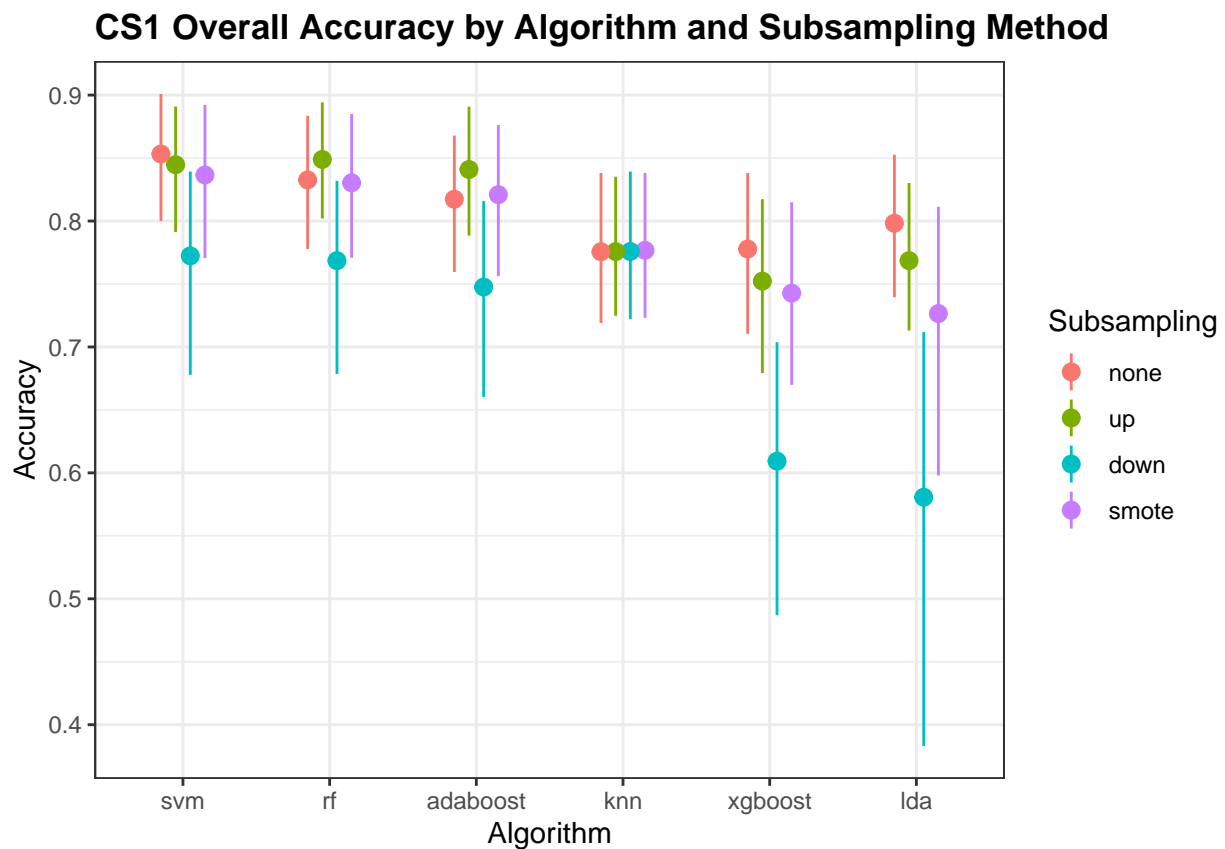


Figure 4.1: CS1 Accuracy

### 4.1.2 F1-Score

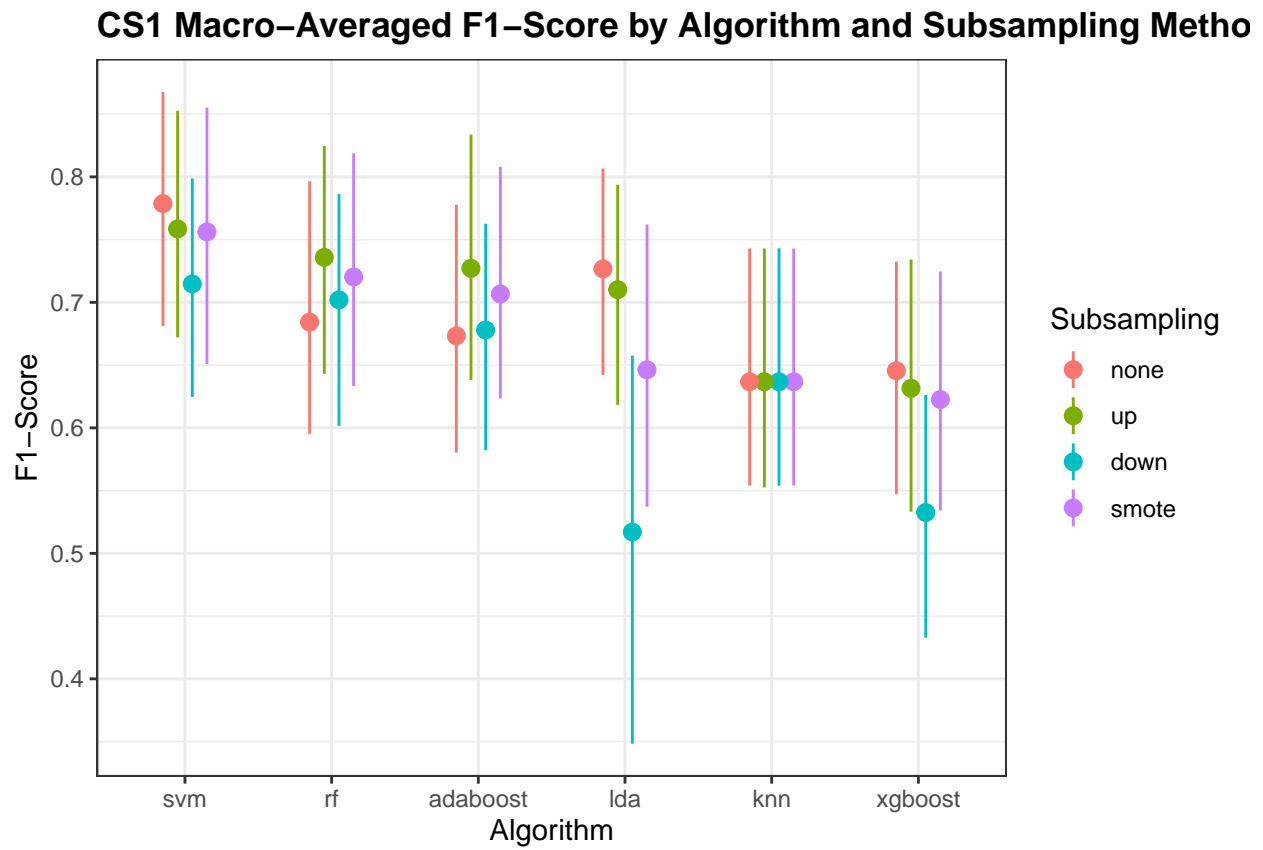


Figure 4.2: CS1 F1-Score

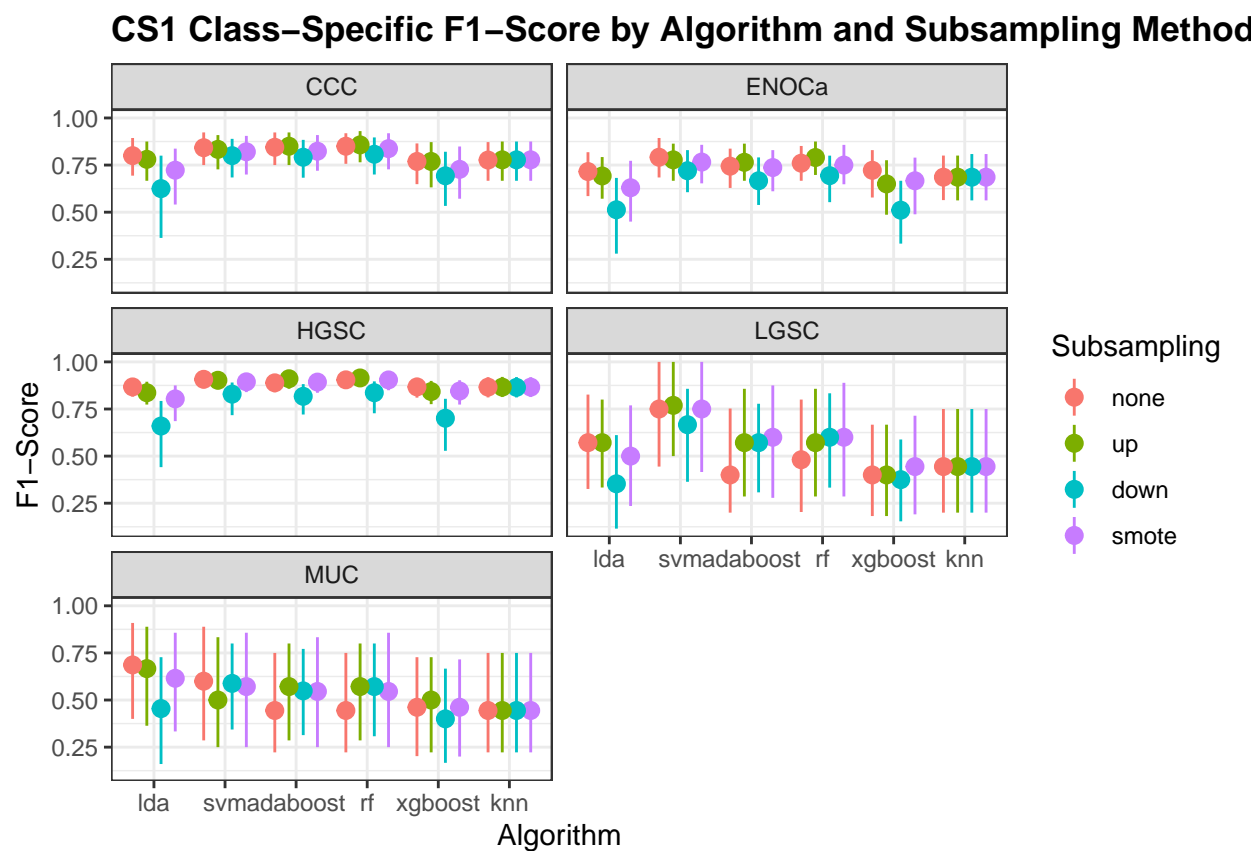


Figure 4.3: CS1 Class-Specific F1-Score

## 4.2 CS2

### 4.2.1 Accuracy

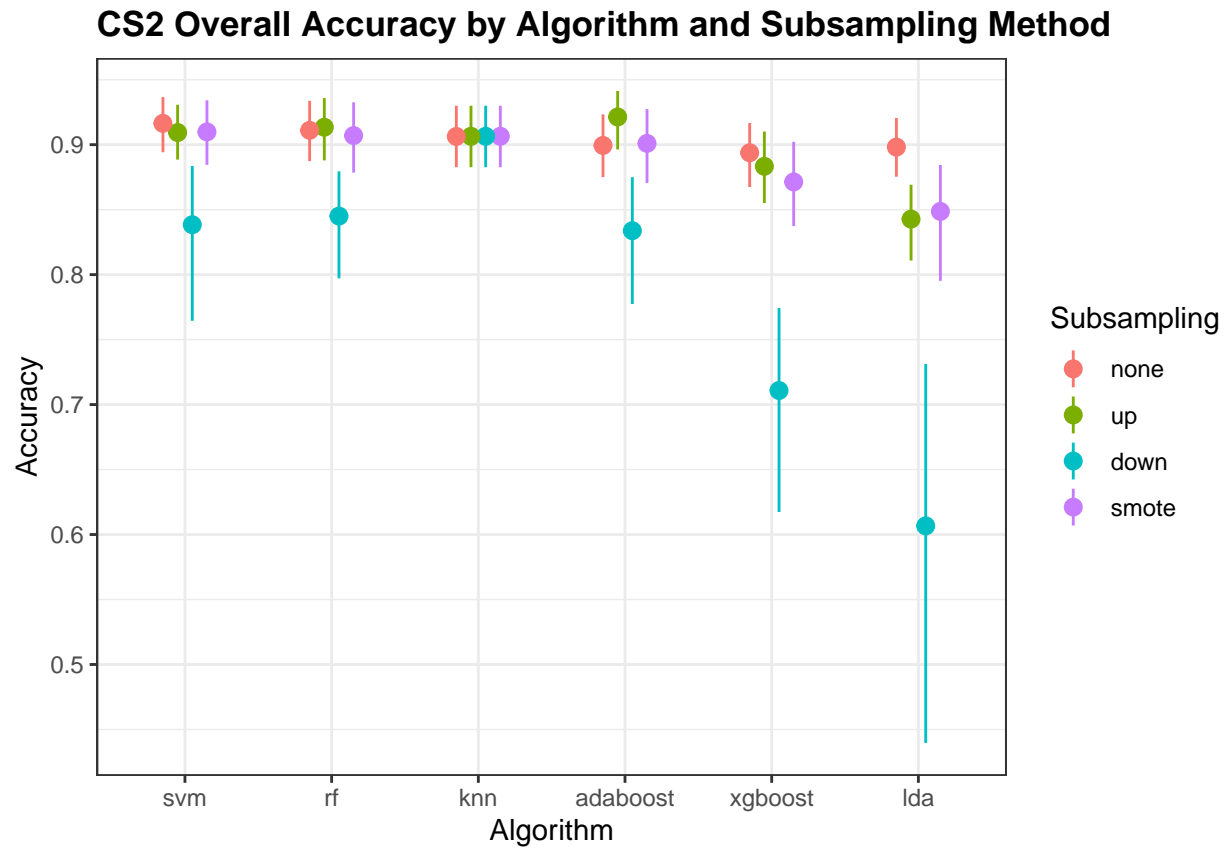


Figure 4.4: CS2 Accuracy



### 4.2.2 F1-Score

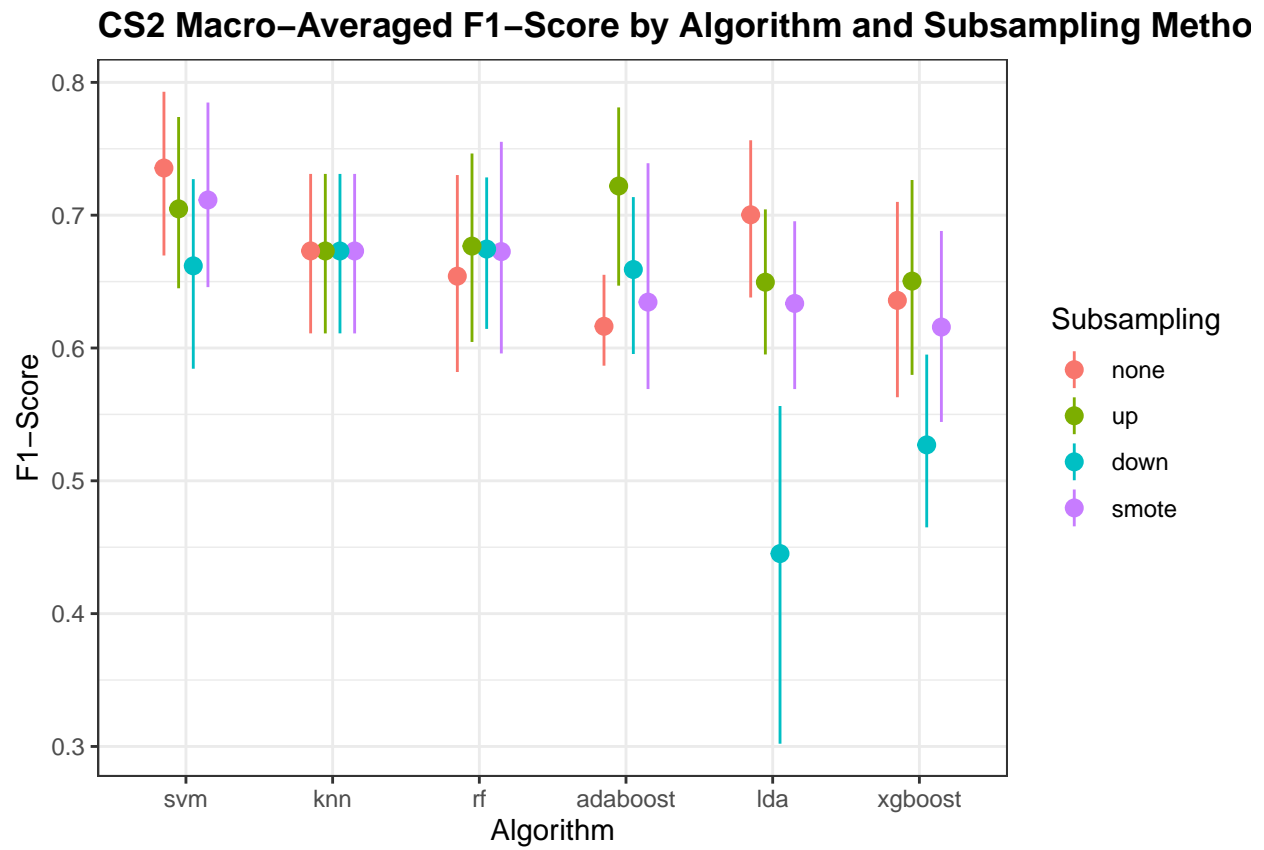


Figure 4.5: CS2 F1-Score

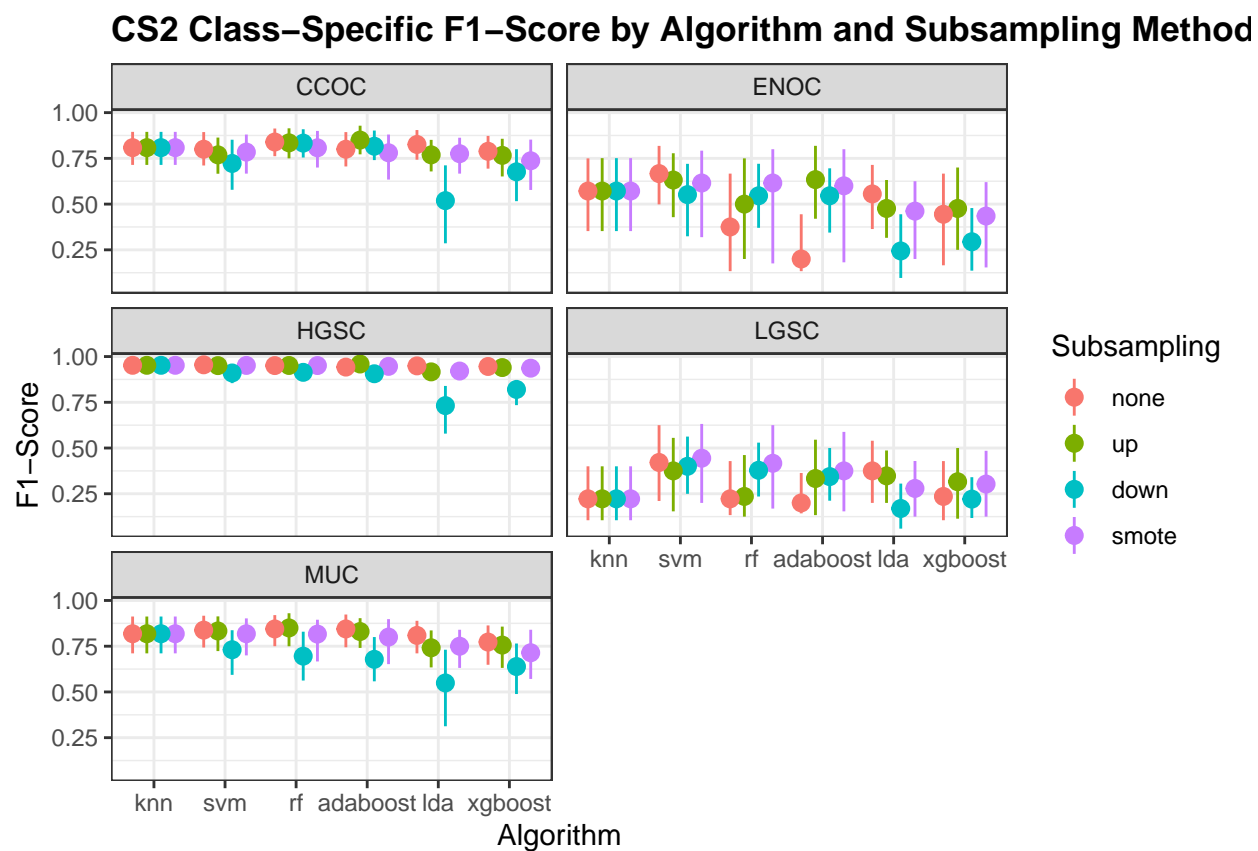


Figure 4.6: CS2 Class-Specific F1-Score