```python
In [2]: import os
        import flexynesis
        import torch
        import numpy as np
        import seaborn as sns
        import pandas as pd
        import random
        import lightning as pl
```
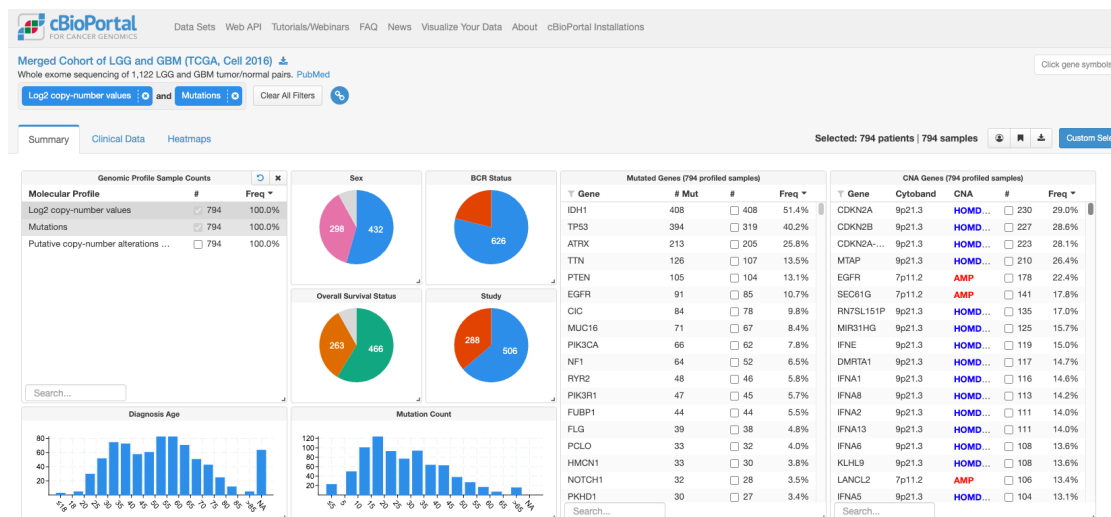
Seed set to 42

# Finding Survival Markers in Lower Grade Glioma (LGG) and Glioblastoma Multiforme (GBM)seed



Here, we demonstrate the capabilities of `flexynesis` on a multi-omic dataset of 506 Brain Lower Grade Glioma (LGG) and 288 Glioblastoma Multiforme (GBM) samples with matching mutation and copy number alteration data downloaded from the cbioportal. The data was split into `train` (70% of the samples) and `test` (30% of the samples) data folders. The data files were processed to follow the same nomenclature.

- `cna.csv` contains "copy number alteration" data
- `mut.csv` contains "mutation" data, which is a binary matrix of genes versus samples.
- `clin.csv` contains "clinical/sample metatada", which is a table of clinical parameters such as age, sex, disease type, histological diagnosis, and overall survival time and status.

## Data Download

The data can be downloaded as follows:

```python
In [3]: if not os.path.exists("lgggbm_tcga_pub_processed"):
```

```
      !wget -O lgggbm_tcga_pub_processed.tgz "https://bimsbstatic.mdc-berli
```

# Importing Train and Test Datasets

We import train and test datasets including mutations and copy number alterations. We rank genes by Laplacian Scores and pick top 10% of the genes, while removing highly redundant genes with a correlation score threshold of 0.8 and a variance threshold of 50%. By setting `concatenate` to `False`, we will be doing an `intermediate` fusion of omic layers.

```
In [4]: data_importer = flexynesis.DataImporter(path ='lgggbm_tcga_pub_processed'
                                     data_types = ['mut', 'cna'], log_
                                     concatenate=False, top_percentile
                                     variance_threshold=0.5)
        train_dataset, test_dataset = data_importer.import_data()
```

```
[INFO] ================= Importing Data =================
[INFO] Validating data folders...

[INFO] ---------------- Reading Data ----------------
[INFO] Importing lgggbm_tcga_pub_processed/train/cna.csv...
[INFO] Importing lgggbm_tcga_pub_processed/train/clin.csv...
[INFO] Importing lgggbm_tcga_pub_processed/train/mut.csv...

[INFO] ---------------- Reading Data ----------------
[INFO] Importing lgggbm_tcga_pub_processed/test/cna.csv...
[INFO] Importing lgggbm_tcga_pub_processed/test/clin.csv...
[INFO] Importing lgggbm_tcga_pub_processed/test/mut.csv...

[INFO] ---------------- Checking for problems with the input data
----------------
[INFO] Data structure is valid with no errors or warnings.

[INFO] ---------------- Processing Data (train) ----------------

[INFO] ---------------- Cleaning Up Data ----------------

[INFO] working on layer:  mut
[INFO] Number of NA values:  0
[INFO] DataFrame mut - Removed 5561 features.

[INFO] working on layer:  cna
[INFO] Number of NA values:  0
[INFO] DataFrame cna - Removed 12375 features.
[INFO] DataFrame mut - Removed 0 samples (0.00%).
[INFO] DataFrame cna - Removed 0 samples (0.00%).
[INFO] Implementing feature selection using laplacian score for layer: mut
with  5503 features  and  556  samples
Calculating Laplacian scores: 100%|
████████████████████████████████████████████████████████████████████
██████████████████████████████████████| 5503/5503 [00:00<00:00, 17621.59it/
s]
Filtering redundant features: 100%|
████████████████████████████████████████████████████████████████████
██████████████████████████████████████| 1000/1000 [00:00<00:00, 16951.61it/
s]
```

```
[INFO] Implementing feature selection using laplacian score for layer: cna
with  12371 features  and  556  samples
Calculating Laplacian scores: 100%|
████████████████████████████████████████████████████████████████████████████
████████████████████████████████████| 12371/12371 [00:00<00:00, 19969.51it/
s]
Filtering redundant features: 100%|
████████████████████████████████████████████████████████████████████████████
████████████████████████████████████| 1237/1237 [00:00<00:00, 275229.65it/
s]
[INFO] ---------------- Processing Data (test) ----------------

[INFO] ---------------- Cleaning Up Data ----------------

[INFO] working on layer:  mut
[INFO] Number of NA values:  0
[INFO] DataFrame mut - Removed 5627 features.

[INFO] working on layer:  cna
[INFO] Number of NA values:  0
[INFO] DataFrame cna - Removed 12382 features.
[INFO] DataFrame mut - Removed 0 samples (0.00%).
[INFO] DataFrame cna - Removed 0 samples (0.00%).

[INFO] ---------------- Harmonizing Data Sets ----------------

[INFO] ---------------- Finished Harmonizing ----------------

[INFO] ---------------- Normalizing Data ----------------

[INFO] ---------------- Normalizing Data ----------------
[INFO] Training Data Stats:  {'feature_count in: cna': 1237, 'feature_coun
t in: mut': 318, 'sample_count': 556}
[INFO] Test Data Stats:  {'feature_count in: cna': 1237, 'feature_count in
: mut': 318, 'sample_count': 238}
[INFO] Merging Feature Logs...
[INFO] Data import successful.
```

# 1. Exploratory Data Analysis

Before building any machine learning models on the data, it is important to first familiarize yourself with the data you are working with. It is important to know the available data matrices, their sizes/shapes, available clinical variables and how they are distributed.

Below you are asked to do simple explorations of the available data.

## 1.1 Print the shapes of the available data matrices

- How many features and samples are available per data type in train/test datasets?

```
In [5]: train_dataset.dat, test_dataset.dat
```

```
Out[5]:  ({'cna': tensor([[-0.1857, -0.7226, -0.8444,  ...,  0.4157, -0.8662,
         -0.8662],
                  [-0.2515, -0.7058, -0.8262,  ...,  0.4130, -0.8484, -0.8484],
                  [-0.2568, -0.7394,  0.8417,  ..., -1.8322,  0.7924,  0.7924],
                  ...,
                  [-0.2252,  3.8013, -0.6767,  ...,  0.6435,  4.1296,  4.1296],
                  [-0.2647, -0.7437, -0.8670,  ..., -1.8981, -0.8885, -0.8885],
                  [-0.1752, -0.7416, -0.8444,  ...,  0.3883, -0.8662, -0.866
         2]]),
            'mut': tensor([[ 0.9822, -0.1485,  1.7075,  ..., -0.0424, -0.0424,
         -0.0601],
                  [ 0.9822, -0.1485,  1.7075,  ..., -0.0424, -0.0424, -0.0601],
                  [ 0.9822, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  ...,
                  [-1.0182, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  [ 0.9822, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  [ 0.9822, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.060
         1]])},
           {'cna': tensor([[-0.2331, -0.7331, -0.8557,  ...,  0.3800, -0.8773,
         -0.8773],
                  [-0.2541, -0.1700, -0.2483,  ...,  0.4377, -0.2799, -0.2799],
                  [ 2.0530, -0.7437, -0.8670,  ...,  0.4377, -0.8885, -0.8885],
                  ...,
                  [-0.2462, -0.7142, -0.8648,  ..., -1.8816, -0.8863, -0.8863],
                  [-0.1910,  0.9163,  0.9913,  ...,  0.4020,  0.8726,  0.8726],
                  [-0.2620,  0.1788,  0.2412,  ...,  0.4184,  0.2016,  0.201
         6]]),
            'mut': tensor([[ 0.9822, -0.1485,  1.7075,  ..., -0.0424, -0.0424,
         -0.0601],
                  [-1.0182, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  [ 0.9822, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  ...,
                  [ 0.9822, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  [-1.0182, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.0601],
                  [-1.0182, -0.1485, -0.5857,  ..., -0.0424, -0.0424, -0.060
         1]])})
```

```
In [6]:  train_dataset.dat['mut'].shape, train_dataset.dat['cna'].shape
```

```
Out[6]:  (torch.Size([556, 318]), torch.Size([556, 1237]))
```

```
In [7]:  test_dataset.dat['mut'].shape, test_dataset.dat['cna'].shape
```

```
Out[7]:  (torch.Size([238, 318]), torch.Size([238, 1237]))
```

```
In [8]:  train_dataset.samples[1:20], train_dataset.features
```

```
Out[8]:  (['TCGA-S9-A6TV',
          'TCGA-HW-8322',
          'TCGA-06-5415',
          'TCGA-VM-A8CB',
          'TCGA-HT-7860',
          'TCGA-HW-7487',
          'TCGA-WH-A86K',
          'TCGA-76-4931',
          'TCGA-28-5204',
          'TCGA-19-5958',
          'TCGA-FG-6690',
          'TCGA-DU-A5TS',
          'TCGA-27-1835',
          'TCGA-DU-8168',
          'TCGA-FG-A6J1',
          'TCGA-DU-A7TA',
          'TCGA-HT-A740',
          'TCGA-06-0169',
          'TCGA-FG-A6J3'],
         {'cna': Index(['SLC30A8', 'ZNF273', 'CLEC5A', 'AGL', 'KCNA5', 'MIR603',
         'SNTB1',
                 'MRPL13', 'MTBP', 'SNORA72|ENSG00000252158.1',
                 ...
                 'CAV1', 'FZD1', 'BCAP29', 'MNX1', 'ADAM22', 'LRP8', 'NOM1', 'RN
         7SL290P',
                 'snoU13|ENSG00000239044.1', 'CPA1'],
               dtype='object', length=1237),
          'mut': Index(['IDH1', 'IDH2', 'ATRX', 'RELN', 'PIK3CA', 'EGFR', 'TP53'
         , 'COL6A3',
                 'TEKT4', 'SVIL',
                 ...
                 'ZNF571', 'PRDM1', 'MYO19', 'ADAMTSL4', 'VDAC3', 'WNT7A', 'ARHG
         AP29',
                 'CTSW', 'PDE10A', 'PTPRS'],
               dtype='object', length=318)})
```

```
In [9]:  train_dataset.samples[1:20], train_dataset.features
```

```
Out[9]:  (['TCGA-S9-A6TV',
          'TCGA-HW-8322',
          'TCGA-06-5415',
          'TCGA-VM-A8CB',
          'TCGA-HT-7860',
          'TCGA-HW-7487',
          'TCGA-WH-A86K',
          'TCGA-76-4931',
          'TCGA-28-5204',
          'TCGA-19-5958',
          'TCGA-FG-6690',
          'TCGA-DU-A5TS',
          'TCGA-27-1835',
          'TCGA-DU-8168',
          'TCGA-FG-A6J1',
          'TCGA-DU-A7TA',
          'TCGA-HT-A740',
          'TCGA-06-0169',
          'TCGA-FG-A6J3'],
         {'cna': Index(['SLC30A8', 'ZNF273', 'CLEC5A', 'AGL', 'KCNA5', 'MIR603',
'SNTB1',
                'MRPL13', 'MTBP', 'SNORA72|ENSG00000252158.1',
                ...
                'CAV1', 'FZD1', 'BCAP29', 'MNX1', 'ADAM22', 'LRP8', 'NOM1', 'RN
7SL290P',
                'snoU13|ENSG00000239044.1', 'CPA1'],
             dtype='object', length=1237),
          'mut': Index(['IDH1', 'IDH2', 'ATRX', 'RELN', 'PIK3CA', 'EGFR', 'TP53'
, 'COL6A3',
                'TEKT4', 'SVIL',
                ...
                'ZNF571', 'PRDM1', 'MYO19', 'ADAMTSL4', 'VDAC3', 'WNT7A', 'ARHG
AP29',
                'CTSW', 'PDE10A', 'PTPRS'],
             dtype='object', length=318)})
```

## 1.2 Explore sample annotations

- What are the available clinical variables? Are they available in both train and test datasets? (See .ann)

```
In [10]:  train_dataset.ann, test_dataset.ann
```

```
Out[10]:  ({'AGE': tensor([nan, 50., 39., 60., 33., 60., 39., 65., 70., 72., 56.,
         70., 42., 53.,
                 55., 44., 32., 34., 68., 52., nan, nan, 31., 46., 41., 59., na
         n, 54.,
                 63., 49., 56., 56., 33., 57., 55., 72., 64., 53., 62., 41., 53
         ., 51.,
                 nan, 52., 39., 61., 60., 21., 24., 79., 58., nan, 53., 74., 49
         ., 54.,
                 nan, 23., 38., 36., 45., 60., 46., 17., nan, 75., 53., 43., 64
         ., 67.,
                 26., 26., 73., 39., 33., 58., 43., 34., nan, 65., 28., 64., 63
         ., 75.,
                 65., nan, 68., 75., 33., 30., nan, 57., 47., 59., 38., 28., 69
         ., 44.,
                 40., 57., 21., 67., 36., 48., 37., 47., 47., 26., 43., 48., 24
         ., 51.,
                 76., 44., 31., 30., 56., 54., nan, 66., 59., 64., 78., 65., 51
         ., nan,
                 47., 51., 49., 42., 47., 40., 58., 70., 66., 72., 50., 29., 81
         ., 52.,
                 66., 74., 77., nan, 64., 72., 52., 49., 51., nan, 20., 76., 36
         ., 61.,
                 30., 48., 48., 63., 74., 70., 37., 48., 25., 61., 76., 53., 40
         ., 42.,
                 41., 59., 68., 58., 51., 53., 44., 44., 81., 81., 63., 52., 66
         ., 34.,
                 47., 35., 62., 61., 39., 86., 33., 71., 63., 84., 77., 56., 53
         ., 64.,
                 52., nan, 30., 30., 80., 36., 53., 57., 41., 37., 59., 68., 60
         ., 50.,
                 27., 76., 67., 70., 32., 37., 44., 43., 51., 59., 46., 60., 88
         ., 68.,
                 58., 30., nan, 45., 49., 25., 78., 60., 36., 43., 60., nan, 81
         ., 70.,
                 40., 54., 76., 14., 20., 76., 37., 45., 55., 71., 70., 38., 65
         ., 38.,
                 nan, 26., 32., 27., 35., 34., 50., 65., 67., 34., 58., 40., 39
         ., 59.,
                 49., 36., 47., 70., nan, 73., 54., 18., 47., 60., 53., 22., 31
         ., 60.,
                 31., 71., 62., 33., 29., 43., 23., 55., 53., 29., 69., 66., 40
         ., 47.,
                 42., 51., 65., 78., 85., 30., 38., 72., 36., 59., 46., 31., 51
         ., 40.,
                 nan, 30., 33., 37., 60., 72., 25., nan, 25., 54., 57., 46., na
         n, 55.,
                 26., 29., 69., 30., 23., 42., 27., 57., 57., 64., 47., 63., 74
         ., 32.,
                 40., 67., 45., 82., 44., 48., 74., 60., 39., 32., 67., 85., 30
         ., 75.,
                 41., 38., 31., 35., 57., 29., 40., 38., 35., 58., 30., 49., 51
         ., 37.,
                 62., 32., 77., 67., 54., 74., 33., 66., 57., 58., 61., 43., 21
         ., 31.,
                 52., 69., 31., 65., 25., 64., 60., 35., 74., 52., 32., 59., 49
         ., 62.,
                 72., 34., 38., nan, 25., 50., 27., 37., 46., 21., nan, 45., 52
         ., 42.,
                 74., 66., 47., 52., 38., 29., 29., 34., 22., 72., 31., 63., 46
         ., 58.,
```

```
                38., 52., 62., 30., 66., 28., nan, 38., 64., 53., 63., 51., 78
        ., 43.,
                53., 67., 29., 63., nan, 49., 34., 36., 27., nan, 49., 67., 69
        ., 60.,
                48., 21., 78., 74., 66., 33., 56., 41., 50., 81., 54., 47., 37
        ., 62.,
                75., 59., 39., 52., 60., 28., 34., 44., 26., nan, 61., 33., 50
        ., 59.,
                73., 69., 72., 56., 48., 60., 54., 55., 48., 76., 65., 40., 39
        ., 34.,
                43., 56., 53., 55., 33., 63., 36., 61., 54., 67., nan, 33., 31
        ., 35.,
                74., 29., 25., 51., nan, 62., 43., 51., 67., 81., 54., nan, 56
        ., 35.,
                24., 28., 32., 53., 50., 69., 43., 45., 65., 58., 58., 67., 48
        ., 62.,
                30., 40., nan, 76., 54., 26., 62., nan, 42., 41., 39., 45., 31
        ., 40.,
                32., nan, 21., 30., 73., 54., 68., 63., 38., 43.], dtype=torc
        h.float64),
          'OS_MONTHS': tensor([        nan, 1.8800e+01, 1.7800e+01, 8.5000e+00,
        1.0000e-01, 5.0000e-01,
                1.3400e+01, 5.3000e+00, 9.2000e+00, 1.4900e+01, 5.4000e+00,
        2.5000e+01,
                1.5100e+01, 2.1300e+01, 1.4200e+01, 4.1000e+00, 7.7800e+01,
        1.0000e-01,
                3.3000e+00, 1.0600e+01,         nan,         nan, 1.0000e-01,
        9.4000e+00,
                1.9200e+01, 7.9000e+00,         nan, 2.8000e+00, 2.3700e+01,
        1.3000e+01,
                8.0000e+00, 4.2300e+01, 4.3700e+01, 1.6700e+01, 7.4000e+00,
        4.5000e+00,
                1.6000e+00, 6.2000e+01, 3.9500e+01, 9.0700e+01, 4.8000e+00,
        3.0000e+00,
                        nan, 1.0370e+02, 4.0100e+01, 1.4600e+01, 2.7000e+00,
        3.0000e-01,
                4.0000e-01, 2.0900e+01, 1.4900e+01,         nan, 2.7000e+00,
        1.8000e+01,
                6.2000e+00, 1.4700e+01,         nan, 1.1400e+02, 1.2700e+01,
        4.7900e+01,
                1.4300e+01, 1.9100e+01, 1.7400e+01, 5.0100e+01,         nan,
        2.7000e+00,
                1.0000e-01, 1.5300e+01, 8.3000e+00, 2.4200e+01, 2.5000e+00,
        1.3000e+00,
                2.2900e+01, 3.6800e+01, 8.1900e+01, 3.7000e+00, 9.0000e-01,
        7.2900e+01,
                        nan, 1.4700e+01, 1.5300e+01, 3.0000e-01, 5.2000e+00,
        3.6000e+00,
                4.3000e+00,         nan, 1.2700e+01, 9.5000e+00, 4.4000e+01,
        7.1000e+00,
                        nan, 4.7000e+00, 1.5300e+01, 1.8800e+01, 0.0000e+00,
        7.8000e+00,
                8.0000e+00, 2.0000e-01, 3.4000e+00, 8.0000e+00, 4.7000e+00,
        7.0000e-01,
                5.7600e+01, 5.9300e+01, 8.6000e+00, 1.1800e+01, 1.8700e+01,
        2.7600e+01,
                1.9000e+01, 1.5600e+01, 1.0000e-01, 1.0500e+01, 1.1000e+00,
        9.7000e+00,
                1.7900e+01, 9.4000e+00, 7.6000e+00, 6.5700e+01,         nan,
        3.7000e+00,
```

```
              9.3000e+00, 1.3800e+01, 1.2000e+00, 7.0000e-01, 1.9900e+01,
       nan,
              1.1700e+01, 1.4300e+01, 2.9600e+01, 4.0000e+00, 1.5700e+01,
       4.4400e+01,
              1.5000e+00, 1.2000e+01, 3.4000e+00, 6.3000e+00, 3.3100e+01,
       2.3600e+01,
              1.3600e+01, 2.3000e+01, 5.8000e+00, 2.6600e+01, 5.3000e+00,
       nan,
              8.3000e+00, 8.5000e+00, 7.0000e+00, 3.8000e+00, 2.1300e+01,
       nan,
              6.7000e+00, 8.0000e+00, 1.8000e+00, 3.2000e+00, 5.4000e+00,
       6.0300e+01,
              1.0400e+01, 1.8400e+01, 2.8900e+01, 1.8700e+01, 7.5200e+01,
       6.2100e+01,
              3.8000e+00, 2.9900e+01, 1.5200e+01, 9.3000e+00, 8.8800e+01,
       6.1000e+00,
              4.2000e+00, 2.0000e-01, 4.5000e+00, 1.5000e+01, 2.0000e-01,
       1.0000e+01,
              2.4200e+01, 4.3900e+01, 2.7000e+00, 2.2000e+00, 1.2200e+01,
       0.0000e+00,
              1.6600e+01, 5.0800e+01, 2.7000e+00, 3.0000e-01, 5.7900e+01,
       2.2000e+01,
              2.0000e-01, 2.0000e-01, 2.0000e+01, 2.8500e+01, 2.0000e-01,
       6.1000e+00,
              9.8000e+00, 1.2600e+01, 9.4300e+01, 1.6000e+00, 1.3600e+01,
       nan,
              1.9400e+01, 3.1400e+01, 8.8000e+00, 7.0000e+00, 1.2000e+00,
       3.0700e+01,
              2.8000e+00, 2.9900e+01, 1.4900e+01, 5.9000e+00, 2.1000e+01,
       1.0000e-01,
              6.0000e-01, 4.8000e+00, 3.4000e+00, 4.9000e+00, 2.0000e-01,
       1.0000e-01,
              1.1300e+01, 8.2000e+00, 1.7500e+01, 1.8800e+01, 1.2200e+01,
       1.5900e+01,
              1.0000e+00, 4.0000e+00, 1.4300e+01, 5.0000e+00,         nan,
       4.3000e+00,
              1.0000e+00, 1.0510e+02, 1.3000e+00, 1.4900e+01, 4.0000e+00,
       3.0000e+00,
              4.2000e+00,         nan, 1.4400e+01, 5.3000e+00, 1.6800e+01,
       3.6000e+00,
              4.8000e+00, 1.5410e+02, 2.8500e+01, 3.6000e+00, 6.3500e+01,
       5.3000e+00,
              7.5000e+00, 2.0000e-01, 5.0000e+00, 9.0000e+00, 4.4000e+00,
       1.9900e+01,
                     nan, 6.5000e+00, 2.5000e+00, 7.2000e+00, 3.9100e+01,
       2.1700e+01,
              4.7000e+00, 1.2800e+01, 7.4000e+00, 1.8220e+02, 1.8300e+01,
       1.1900e+01,
              2.6300e+01, 8.8000e+00, 2.9000e+01, 6.4000e+00, 6.3000e+00,
       3.0200e+01,
                     nan, 1.8000e+01, 4.7700e+01, 6.0000e+01, 4.5000e+01,
       1.1500e+01,
              2.6000e+00, 1.1300e+01, 7.0500e+01, 2.4200e+01, 8.7400e+01,
       1.2300e+01,
              2.5000e+00, 5.7000e+00, 7.6000e+00, 2.3300e+01, 4.0900e+01,
       5.7000e+00,
              1.5400e+01, 2.0700e+01, 2.2300e+01, 2.2700e+01, 1.4300e+01,
       1.4300e+01,
              1.6000e+01, 2.1400e+01, 1.3600e+01, 4.5000e+00, 3.9000e+00,
       1.4900e+01,
```

```
           1.3000e+01, 9.8000e+00, 3.9000e+00, 1.9600e+01, 1.3300e+01,
2.0100e+01,
           4.3000e+00, 2.8300e+01,        nan, 1.1500e+01, 1.7700e+01,
1.8900e+01,
           5.4000e+00, 8.0000e+00, 2.0000e+00,        nan, 5.6100e+01,
1.1700e+01,
           7.2000e+00, 1.0900e+01,        nan, 5.0000e-01, 2.0000e-01,
4.2700e+01,
           2.2000e+00, 5.3000e+00, 2.4000e+01, 7.8200e+01, 8.3700e+01,
5.1000e+00,
           7.9000e+00, 5.2000e+00, 8.5000e+00, 7.4000e+00, 9.4000e+00,
2.3500e+01,
           1.6600e+01, 1.3300e+01, 1.0300e+01, 8.9000e+00, 7.1000e+00,
9.5500e+01,
           3.3000e+00, 1.6900e+01, 3.5900e+01, 4.6600e+01, 4.2000e+01,
3.1000e+00,
           2.5300e+01, 1.2500e+01, 1.3000e+01, 2.9000e+00, 3.7000e+00,
0.0000e+00,
           8.4000e+00, 1.4500e+01, 1.0000e-01, 9.7800e+01, 4.0100e+01,
6.2300e+01,
           1.0300e+01, 1.7200e+01, 1.2400e+01, 8.8000e+00, 1.5900e+01,
1.3820e+02,
           4.8000e+00, 8.8000e+00, 1.1200e+01, 5.5000e+00, 1.0500e+01,
0.0000e+00,
           1.1400e+01, 1.7800e+01, 2.4100e+01, 3.1300e+01, 2.9800e+01,
1.3420e+02,
           1.4000e+00, 1.0400e+01, 4.9000e+01, 1.3800e+01, 1.3100e+01,
1.6200e+01,
           3.4900e+01, 1.6430e+02, 5.0000e+00, 2.4700e+01, 4.3900e+01,
4.4500e+01,
           6.8000e+00, 1.1000e+01, 3.5000e+00, 1.8000e+01, 2.7800e+01,
nan,
           2.6800e+01, 5.8000e+00, 1.8800e+01, 5.3000e+00, 2.8100e+01,
1.0000e-01,
                  nan, 4.6000e+01, 1.0800e+01, 2.0800e+01, 6.5000e+00,
1.2000e+00,
           6.6000e+00, 4.2000e+01, 6.7400e+01, 1.3400e+01, 1.7000e+01,
6.6000e+00,
           8.5000e+00, 8.9000e+00, 2.0600e+01, 5.3000e+00, 1.6400e+01,
6.0000e+00,
           4.6000e+00, 9.4000e+01, 3.6000e+00, 1.5500e+01, 5.2100e+01,
3.6800e+01,
                  nan, 2.2000e+00, 3.9000e+00, 9.4000e+00, 4.7600e+01,
1.2600e+01,
           1.5300e+01, 0.0000e+00, 2.5900e+01, 7.7000e+00, 4.4600e+01,
1.6100e+01,
                  nan, 7.8000e+00, 2.0000e-01, 1.4700e+01, 3.8200e+01,
nan,
           2.0100e+01, 2.0400e+01, 1.5500e+01, 6.9000e+00, 6.0000e+00,
2.4000e+00,
           6.7000e+00, 7.0000e+00, 1.9500e+01, 3.0100e+01, 7.0600e+01,
1.4900e+01,
           5.9000e+00, 2.8000e+00, 1.0000e-01, 1.0300e+01, 3.0100e+01,
1.0400e+01,
           2.0000e-01, 8.9000e+00, 5.9100e+01, 7.5000e+00, 1.1600e+01,
5.0600e+01,
           7.0000e+00, 3.1600e+01, 3.2400e+01,        nan, 9.0000e+00,
6.2900e+01,
           2.3000e+01, 2.0000e+00, 1.0800e+01, 1.9500e+01, 7.4000e+00,
3.2000e+00,
```

```
             5.0000e+00, 2.9000e+00, 6.4000e+00, 2.0000e-01, 5.0700e+01,
7.6000e+00,
             1.3100e+01, 1.2000e+01, 2.3000e+00, 4.6000e+01, 1.5700e+01,
4.0500e+01,
             1.6200e+01, 2.0500e+01, 1.4200e+01, 2.5400e+01, 2.1600e+01,
9.2000e+00,
             1.2900e+01, 1.0500e+01,        nan, 7.5100e+01, 6.0000e+00,
2.0700e+01,
             6.0000e+00, 2.3200e+01, 1.0000e-01, 8.3000e+00,        nan,
1.8900e+01,
             1.7500e+01, 1.2200e+01, 1.1000e+01, 4.4000e+00, 2.0000e-01,
nan,
             5.4000e+00, 1.3700e+01, 1.0900e+01, 8.7000e+01, 8.0000e-01,
3.2600e+01,
             2.2900e+01, 7.3000e+00, 4.8000e+00, 7.5000e+01, 1.4900e+01,
5.4000e+00,
             1.7700e+01, 1.7400e+01, 1.0000e-01, 2.4200e+01, 1.4400e+01,
1.6500e+01,
                    nan, 4.1000e+00, 2.6700e+01, 5.2000e+00, 1.0600e+01,
nan,
             1.1800e+01, 3.2300e+01, 2.3600e+01, 6.1000e+00, 2.2500e+01,
5.0000e+00,
             7.2900e+01,        nan, 9.4000e+00, 6.7000e+00, 1.0800e+01,
1.3900e+01,
             1.0500e+01, 4.0300e+01, 5.7000e+00, 3.4000e+00], dtype=torch.f
loat64),
   'OS_STATUS': tensor([nan, 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0.,
0., 1., 0., 0., 0., 0.,
             1., 0., nan, nan, 0., 0., 0., 1., nan, 1., 1., 1., 0., 0., 0.,
1., 0., 1.,
             0., 1., 0., 0., 0., 0., nan, 0., 0., 1., 0., 0., 0., 0., 0., n
an, 1., 1.,
             0., 0., nan, 1., 0., 0., 0., 0., 0., 0., nan, 1., 0., 0., 0.,
1., 0., 0.,
             0., 1., 0., 1., 1., 0., nan, 0., 0., 0., 1., 1., nan, nan, 1.,
1., 1., 0.,
             nan, 0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 1., 0., 0., 0., 1
., 0., 0.,
             1., 0., 0., 0., 1., 0., 0., 0., 1., 1., nan, 1., 0., 1., 0., 0
., 1., nan,
             0., 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 1., 0., 1., 1.
, 0., nan,
             1., 0., 1., 0., 1., nan, 0., 1., 0., 1., 0., 0., 0., 1., 1., 0
., 0., 1.,
             0., 0., 1., 0., 0., 0., 0., 0., 1., 1., 1., 0., 0., 1., 1., 1.
, 1., 0.,
             1., 1., 0., 0., 1., 0., 0., 1., 0., 0., 1., 1., 1., 1., 0., 1.
, 0., nan,
             0., 0., 1., 0., 0., 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1.
, 0., 0.,
             0., 0., 0., 0., 0., 1., 1., 1., 0., 0., nan, 0., 1., 1., 1., 1
., 0., 0.,
             0., nan, 1., 1., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 0., 0
., 1., 1.,
             nan, 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 1., 0., 1., 1., 0
., 0., 1.,
             nan, 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 0., 1., 0., 0., 1
., 0., 1.,
             0., 0., 0., 1., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 0., 1.
, 1., 0.,
```

```
             0., 0., nan, 1., 0., 0., 1., 0., 0., nan, 0., 1., 0., 1., nan,
        1., 0., 0.,
             1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 0., 0., 0., 1., 1., 0.
        , 0., 1.,
             1., 1., 0., 1., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
             0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 0., 1., 1., 1., 0.
        , 0., 1.,
             1., 1., 1., 1., 0., 0., 1., 0., 1., 1., 0., 0., 0., 1., 1., 0.
        , 0., nan,
             0., 0., 0., 0., 0., 0., nan, 1., 0., 0., 1., 1., 1., 0., 1., 0
        ., 0., 0.,
             0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1., nan, 0., 0., 0
        ., 1., 0.,
             1., 0., 1., 1., 0., 1., nan, 0., 0., 1., 0., nan, 0., 1., 0.,
        0., 0., 0.,
             0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 1.
        , 1., 0.,
             0., 1., 1., nan, 0., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0
        ., 0., 1.,
             1., 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 1., 0., nan, 0
        ., 0., 0.,
             0., 0., 0., 1., nan, 1., 0., 0., 0., 1., 0., nan, 0., 0., 0.,
        0., 0., 0.,
             0., 0., 0., 1., 0., 1., 1., 0., 0., 0., 0., 0., nan, 1., 1., 0
        ., 1., nan,
             0., 1., 1., 0., 0., 0., 0., nan, 0., 0., 1., 1., 0., 0., 0.,
        0.],
          dtype=torch.float64),
   'KARNOFSKY_PERFORMANCE_SCORE': tensor([ nan,  90.,  nan, 100.,  nan,
        nan,  nan,  nan,  80.,  80.,  80.,  90.,
             100.,  80.,  70.,  70.,  80.,  nan,  80.,  90.,  nan,  nan,  n
        an,  80.,
              nan,  80.,  nan, 100.,  70.,  60.,  50., 100.,  90.,  80.,  n
        an,  nan,
              nan,  90., 100.,  90.,  nan,  90.,  nan, 100., 100.,  80.,  n
        an,  nan,
              80.,  80.,  60.,  nan,  80.,  80.,  90.,  90.,  nan,  90.,  n
        an,  nan,
             100.,  90.,  90.,  nan,  nan,  60.,  nan,  80.,  40.,  80.,  n
        an,  70.,
              80.,  90.,  nan,  nan,  80.,  nan,  nan,  70.,  nan,  nan,  n
        an,  60.,
              nan,  nan,  20.,  60.,  90.,  nan,  nan,  nan,  90.,  nan,  n
        an,  90.,
              nan,  nan,  60.,  90.,  90.,  nan,  90.,  nan,  70., 100., 10
        0.,  nan,
              90., 100.,  nan,  90.,  40.,  nan,  90.,  nan,  80.,  nan,  n
        an,  nan,
              90.,  60.,  nan,  20.,  90.,  nan,  nan,  nan,  nan, 100.,  8
        0.,  70.,
              nan,  60., 100.,  nan,  90.,  nan,  nan,  80.,  60., 100.,  8
        0.,  nan,
              80.,  80.,  nan,  80.,  50.,  nan,  80.,  80.,  nan,  70.,  6
        0.,  90.,
              80.,  nan,  60.,  80.,  80., 100.,  nan,  90.,  80.,  80.,  n
        an, 100.,
              90.,  90.,  80.,  nan,  nan,  90.,  90.,  90.,  nan,  nan,  n
        an,  90.,
              80.,  nan,  nan,  nan,  80.,  nan,  60.,  90., 100., 100.,  n
```

```
an,  nan,
           80.,  nan, 100.,  40.,  80.,  nan,  90.,  nan,  nan,  nan,  n
an,  90.,
           90.,  nan,  80.,  nan,  80.,  nan,  90., 100., 100.,  40.,  n
an,  nan,
           80.,  60.,  90.,  80.,  nan,  60.,  nan,  nan, 100.,  nan,  n
an,  nan,
           nan,  80.,  60.,  nan, 100., 100.,  nan,  nan,  80.,  40.,  n
an, 100.,
          100.,  nan,  nan,  80.,  nan,  60.,  80.,  nan,  nan,  90.,  8
0.,  90.,
           nan,  nan,  90.,  90.,  nan,  nan, 100.,  80.,  nan,  70., 10
0.,  60.,
          100.,  80.,  nan,  nan,  90.,  nan,  nan,  60.,  nan,  nan,  n
an,  nan,
           60., 100., 100.,  nan,  nan,  80.,  80., 100., 100.,  80.,  8
0.,  nan,
           80.,  90., 100.,  80.,  nan,  nan,  70.,  nan,  90.,  80.,  n
an,  70.,
           90.,  80.,  nan, 100.,  80.,  90.,  90.,  90.,  nan,  nan,  n
an, 100.,
           nan,  60.,  80.,  nan,  nan,  80.,  60.,  80.,  nan,  nan,  n
an,  nan,
           60.,  nan, 100., 100.,  80.,  80.,  80.,  80.,  70., 100.,  6
0.,  90.,
          100.,  nan,  nan,  nan,  nan, 100.,  80.,  nan,  90.,  70., 10
0.,  80.,
          100.,  60., 100.,  nan,  nan,  60.,  90.,  nan,  nan,  nan,  n
an, 100.,
           nan,  90.,  90.,  90.,  80.,  nan,  80.,  60.,  40.,  60.,  n
an,  90.,
           80.,  80.,  80.,  80.,  nan, 100.,  80.,  90.,  nan,  80.,  n
an,  80.,
           80.,  90.,  80.,  80., 100.,  nan, 100.,  90.,  60.,  nan,  9
0.,  nan,
           90.,  80.,  90.,  80., 100.,  nan,  nan,  nan,  80.,  nan,  n
an,  nan,
           40.,  90.,  80., 100.,  nan,  nan,  90.,  nan,  nan,  80., 10
0.,  nan,
           nan,  nan,  50.,  nan,  nan,  90.,  nan,  90.,  nan,  nan,  6
0.,  90.,
           60.,  90., 100.,  90.,  nan,  80.,  nan,  80.,  nan, 100., 10
0.,  nan,
           nan,  nan, 100.,  90.,  nan,  nan,  nan,  60.,  80.,  nan, 10
0.,  90.,
          100.,  80.,  nan,  nan,  nan,  80.,  nan,  60.,  nan,  60.,  9
0.,  nan,
          100.,  90.,  90.,  nan,  80.,  nan,  80.,  80.,  80.,  nan,  n
an,  80.,
           70.,  40.,  80.,  nan, 100.,  nan, 100.,  90.,  nan,  90.,  8
0.,  80.,
           nan,  nan,  70.,  60.,  nan,  80.,  80.,  90.,  nan,  nan,  9
0.,  nan,
           90.,  nan,  nan,  80.,  nan,  60.,  nan,  nan,  70.,  nan,  6
0.,  nan,
           80.,  nan,  90., 100.,  nan,  nan,  nan,  80.,  90., 100., 10
0.,  80.,
           90.,  70.,  nan,  80.,  80., 100.,  nan,  nan,  nan,  90.,  6
0.,  nan,
          100.,  90.,  nan, 100., 100.,  70.,  nan,  nan, 100.,  80.,  n
```

```
         an, 100.,
                nan,  80., 100.,  90.], dtype=torch.float64),
      'STUDY': tensor([0., 0., 0., 1., 0., 0., 0., 0., 1., 1., 1., 0., 0., 1
 ., 0., 0., 0., 0.,
                1., 0., 0., 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 0., 1.
 , 0., 1.,
                0., 0., 0., 0., 1., 0., 1., 0., 0., 0., 0., 0., 1., 1., 0., 0.
 , 1., 1.,
                0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 1., 1., 1.
 , 0., 0.,
                0., 0., 0., 1., 1., 0., 0., 0., 0., 0., 1., 1., 1., 0., 1., 1.
 , 0., 0.,
                0., 0., 1., 0., 0., 0., 0., 0., 1., 0., 1., 1., 0., 0., 0., 1.
 , 0., 0.,
                0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 1., 1., 1.
 , 0., 0.,
                1., 0., 0., 0., 0., 0., 1., 1., 0., 1., 0., 0., 1., 1., 0., 1.
 , 1., 0.,
                1., 1., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 1., 0., 1., 0.
 , 0., 0.,
                0., 0., 1., 1., 0., 0., 0., 0., 1., 0., 1., 0., 0., 0., 1., 1.
 , 1., 0.,
                1., 0., 0., 0., 0., 0., 0., 1., 0., 1., 1., 1., 1., 1., 0., 1.
 , 1., 0.,
                0., 0., 1., 0., 1., 0., 0., 0., 1., 1., 1., 1., 0., 1., 1., 1.
 , 0., 0.,
                0., 0., 0., 0., 0., 1., 1., 1., 1., 0., 0., 0., 1., 0., 1., 1.
 , 0., 0.,
                0., 0., 1., 0., 0., 0., 1., 0., 0., 1., 0., 1., 1., 0., 0., 0.
 , 1., 0.,
                0., 0., 0., 0., 0., 0., 1., 1., 1., 0., 1., 1., 0., 0., 1., 0.
 , 0., 0.,
                0., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 1., 1., 0., 0., 0.
 , 1., 1.,
                1., 0., 1., 1., 0., 0., 0., 0., 0., 1., 1., 0., 0., 1., 0., 1.
 , 1., 0.,
                1., 0., 0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 1., 1., 0., 1.
 , 0., 0.,
                1., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1., 0., 1., 1., 0., 1.
 , 0., 0.,
                1., 1., 0., 1., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0.
 , 0., 0.,
                0., 0., 0., 0., 1., 0., 1., 1., 1., 1., 0., 0., 1., 1., 1., 1.
 , 0., 0.,
                1., 0., 0., 1., 1., 1., 1., 0., 1., 1., 0., 0., 0., 1., 1., 0.
 , 0., 0.,
                0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 0.
 , 0., 0.,
                0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1.
 , 1., 0.,
                1., 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 0., 1., 0., 0.
 , 1., 0.,
                1., 1., 1., 0., 0., 0., 1., 1., 0., 0., 0., 1., 0., 1., 0., 0.
 , 0., 0.,
                0., 0., 0., 0., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 0., 0.
 , 0., 1.,
                1., 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., 1., 1., 1., 0., 0.
 , 0., 0.,
                0., 0., 0., 1., 0., 1., 0., 0., 0., 1., 1., 0., 1., 0., 0., 0.
 , 0., 0.,
```

              0., 1., 0., 0., 1., 1., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0.
        , 1., 0.,
              1., 0., 1., 1., 0., 0., 0., 0., 1., 0., 1., 1., 0., 1., 0.,
        0.],
            dtype=torch.float64),
     'BCR_STATUS': tensor([1., 1., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0.,
        0., 0., 0., 0., 0., 0.,
              0., 0., 1., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 0., 1., 0., 0., 0., 1.
        , 0., 0.,
              0., 0., 1., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 0., 0., 0.
        , 0., 1.,
              1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0.
        , 0., 0.,
              1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0.
        , 0., 0.,
              0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0.
        , 0., 0.,
              0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 1., 0.
        , 0., 1.,
              0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1.
        , 0., 1.,
              0., 1., 0., 0., 1., 1., 1., 1., 0., 0., 0., 0., 1., 0., 0., 0.
        , 0., 1.,
              0., 0., 0., 1., 0., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              1., 0., 0., 1., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0.
        , 0., 0.,
              0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0.
        , 1., 0.,
              0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 1., 0., 1.
        , 0., 0.,
              1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 1., 0.,
              1., 0., 0., 0., 1., 0., 0., 1., 1., 0., 0., 0., 0., 0., 1., 0.
        , 0., 0.,
              0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 1., 1., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 0., 1., 0., 0., 0., 1., 0., 1., 0., 0., 1., 0., 0., 0., 0.
        , 0., 1.,
              0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0.
        , 0., 0.,
              0., 1., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0.
        , 0., 1.,
              1., 0., 0., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 1., 1., 0., 0.
        , 0., 0.,
              0., 1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 0., 1.,
              0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0.
        , 0., 0.,
              0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0.
        , 1., 0.,
              0., 1., 1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0.

```
                , 0., 0.,
                  0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1.
                , 0., 1.,
                  0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0.
                , 0., 1.,
                  0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 1.,
                0.],
                  dtype=torch.float64),
         'HISTOLOGICAL_DIAGNOSIS': tensor([nan, 2., 3., 1., 3., 0., 3., 0., 1.,
        1., 1., 3., 3., 1., 3., 3., 3., 0.,
                  1., 0., nan, nan, 0., 3., 0., 1., nan, 1., 2., 1., 3., 3., 3.,
        1., 2., 1.,
                  3., 3., 0., 3., 1., 2., nan, 3., 0., 2., 3., 0., 1., 1., 3., n
        an, 1., 1.,
                  0., 3., nan, 3., 0., 2., 2., 3., 0., 3., nan, 1., 0., 1., 1.,
        1., 3., 3.,
                  2., 2., 0., 1., 1., 3., nan, 0., 3., 3., 1., 1., 1., nan, 1.,
        1., 0., 0.,
                  nan, 0., 1., 2., 3., 2., 0., 0., 1., 3., 1., 1., 2., 3., 3., 1
        ., 0., 2.,
                  3., 3., 0., 3., 1., 2., 0., 3., 1., 3., nan, 3., 0., 1., 1., 1
        ., 0., nan,
                  1., 3., 3., 0., 0., 3., 1., 1., 3., 1., 3., 2., 1., 1., 0., 1.
        , 1., nan,
                  1., 1., 0., 3., 3., nan, 0., 1., 0., 0., 1., 3., 1., 3., 1., 0
        ., 2., 0.,
                  0., 3., 1., 1., 3., 2., 0., 2., 1., 0., 1., 3., 2., 0., 1., 1.
        , 1., 3.,
                  1., 0., 2., 0., 0., 2., 0., 1., 0., 1., 1., 1., 1., 1., 3., 1.
        , 1., nan,
                  0., 2., 1., 0., 1., 3., 3., 3., 1., 1., 1., 1., 2., 1., 1., 1.
        , 2., 3.,
                  0., 3., 3., 0., 0., 1., 1., 1., 1., 0., nan, 3., 1., 2., 1., 1
        ., 0., 2.,
                  3., nan, 1., 0., 3., 0., 1., 2., 3., 1., 2., 1., 1., 3., 3., 0
        ., 1., 2.,
                  nan, 0., 2., 0., 3., 0., 1., 1., 1., 3., 1., 1., 2., 0., 1., 0
        ., 2., 3.,
                  nan, 3., 3., 2., 3., 0., 1., 0., 0., 1., 3., 1., 1., 0., 0., 2
        ., 1., 1.,
                  1., 2., 1., 1., 3., 3., 2., 3., 2., 1., 1., 0., 3., 1., 2., 1.
        , 1., 2.,
                  1., 3., nan, 3., 1., 2., 1., 1., 3., nan, 2., 1., 1., 1., nan,
        1., 3., 0.,
                  1., 2., 0., 3., 0., 0., 1., 1., 1., 1., 1., 0., 1., 1., 0., 1.
        , 0., 3.,
                  1., 1., 0., 1., 2., 1., 2., 1., 3., 3., 2., 0., 0., 3., 0., 2.
        , 3., 3.,
                  2., 3., 0., 3., 1., 3., 1., 1., 1., 1., 2., 0., 1., 1., 1., 1.
        , 3., 3.,
                  1., 0., 2., 1., 1., 1., 1., 3., 1., 1., 0., 3., 2., 1., 1., 0.
        , 2., nan,
                  3., 1., 2., 0., 3., 0., nan, 3., 1., 0., 0., 3., 1., 0., 0., 2
        ., 2., 2.,
                  0., 1., 3., 2., 0., 2., 3., 3., 0., 3., 3., 2., nan, 2., 3., 1
        ., 1., 2.,
                  1., 2., 3., 0., 2., 1., nan, 1., 3., 1., 3., nan, 2., 1., 2.,
        3., 1., 0.,
                  1., 1., 1., 2., 2., 2., 1., 1., 3., 0., 0., 1., 3., 1., 3., 2.
        , 2., 0.,
```

```
              0., 0., 2., nan, 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 0., 0
., 0., 1.,
              1., 0., 0., 2., 0., 1., 3., 1., 0., 1., 0., 1., 1., 1., nan, 3
., 2., 0.,
              3., 3., 2., 1., nan, 1., 2., 0., 3., 1., 1., nan, 1., 3., 0.,
0., 0., 3.,
              3., 1., 3., 3., 1., 1., 2., 2., 3., 0., 2., 3., nan, 1., 3., 2
., 1., nan,
              1., 0., 1., 1., 3., 0., 2., nan, 1., 3., 1., 1., 3., 1., 3.,
0.],
            dtype=torch.float64),
    'SEX': tensor([nan, 1., 1., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 0.
, 0., 0., 1., 1.,
              1., 0., nan, nan, 0., 1., 1., 1., nan, 1., 0., 1., 1., 1., 0.,
1., 0., 0.,
              0., 0., 1., 0., 0., 1., nan, 1., 1., 0., 1., 1., 0., 0., 0., n
an, 1., 1.,
              1., 1., nan, 0., 1., 1., 0., 0., 1., 1., nan, 0., 1., 1., 1.,
0., 0., 1.,
              0., 0., 1., 1., 1., 1., nan, 0., 1., 0., 1., 0., 1., nan, 1.,
1., 1., 1.,
              nan, 1., 1., 1., 1., 0., 1., 1., 1., 0., 0., 1., 1., 0., 0., 1
., 1., 1.,
              1., 0., 1., 0., 1., 0., 1., 0., 1., 1., nan, 1., 0., 0., 0., 1
., 0., nan,
              0., 0., 0., 1., 1., 1., 0., 1., 1., 0., 0., 1., 0., 1., 1., 1.
, 0., nan,
              0., 0., 1., 0., 1., nan, 0., 0., 0., 0., 1., 0., 1., 1., 0., 1
., 0., 0.,
              0., 0., 1., 0., 1., 1., 1., 0., 1., 0., 1., 0., 0., 0., 0., 1.
, 1., 0.,
              1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 0., 1., 1., 0., 1.
, 1., nan,
              1., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 0., 1., 1., 1., 1.
, 1., 0.,
              0., 1., 0., 0., 0., 0., 1., 1., 1., 0., nan, 0., 1., 1., 1., 0
., 0., 1.,
              1., nan, 1., 0., 1., 0., 1., 1., 0., 0., 0., 1., 0., 0., 1., 0
., 1., 0.,
              nan, 1., 0., 0., 1., 0., 1., 1., 1., 0., 1., 1., 1., 1., 0., 1
., 1., 1.,
              nan, 0., 1., 0., 1., 0., 0., 0., 1., 1., 0., 1., 0., 0., 1., 1
., 1., 0.,
              1., 1., 1., 1., 1., 0., 0., 1., 0., 0., 1., 1., 0., 1., 1., 0.
, 1., 1.,
              1., 1., nan, 1., 1., 0., 1., 0., 1., nan, 0., 0., 1., 1., nan,
0., 1., 0.,
              1., 1., 0., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 0., 1.
, 1., 0.,
              1., 1., 0., 0., 1., 0., 1., 0., 1., 1., 1., 1., 0., 0., 1., 0.
, 1., 1.,
              1., 1., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 1., 0., 0., 1.
, 1., 1.,
              1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 0., 0., 0., 0.
, 0., nan,
              1., 1., 0., 0., 1., 1., nan, 1., 0., 0., 0., 0., 0., 1., 1., 1
., 0., 0.,
              1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., nan, 1., 1., 1
., 0., 1.,
              0., 1., 1., 1., 0., 1., nan, 0., 1., 0., 0., nan, 0., 1., 0.,
```

              1., 1., 1.,
              1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1.
, 0., 1.,
              0., 0., 0., nan, 1., 1., 0., 1., 1., 0., 0., 0., 1., 1., 0., 0
., 1., 1.,
              1., 1., 1., 1., 0., 1., 1., 1., 1., 0., 1., 0., 1., 1., nan, 1
., 0., 1.,
              0., 1., 1., 1., nan, 1., 0., 0., 0., 0., 1., nan, 1., 1., 0.,
0., 0., 0.,
              0., 0., 0., 1., 0., 1., 1., 1., 0., 1., 1., 1., nan, 0., 1., 1
., 1., nan,
              1., 0., 0., 0., 0., 1., 1., nan, 0., 1., 1., 1., 1., 1., 0.,
1.],
          dtype=torch.float64)},
 {'AGE': tensor([nan, 49., 35., 68., 48., 62., 63., 34., 44., 52., 44.,
38., 61., 49.,
              79., 37., 63., 39., 35., 71., 63., 39., 38., 51., nan, 55., 52
., nan,
              30., 36., 52., 19., nan, 35., 74., 56., 83., 54., 63., 32., 53
., 56.,
              53., 29., 41., 62., nan, 86., 55., nan, 60., 58., 69., 48., na
n, 40.,
              nan, 20., 41., 62., 65., 74., 57., 53., 30., 59., 23., 33., na
n, 41.,
              28., 40., 70., 57., 43., 54., 47., 66., 73., 33., nan, 74., 73
., 55.,
              65., nan, 38., 29., 52., 58., 24., nan, 54., 50., 50., nan, 47
., nan,
              35., 52., nan, 33., 79., 37., 36., 31., 39., 36., 30., 57., 59
., 72.,
              30., 66., 76., 34., 36., nan, 62., 36., 30., 66., 87., 30., 21
., 66.,
              83., 48., 31., 41., nan, 64., 34., 53., 20., 61., nan, 28., 24
., 41.,
              73., 32., 63., 58., 33., 56., nan, 54., 65., 34., 50., 42., 48
., 59.,
              52., 59., 37., nan, 44., 61., 59., 62., 43., 62., 78., 33., 61
., 72.,
              60., 43., 89., 43., 35., 44., nan, 38., 23., nan, 41., nan, na
n, 52.,
              63., 53., nan, 64., 22., 54., 55., 44., 32., 31., 38., 31., 75
., 33.,
              76., 59., 50., 73., 25., 31., 29., nan, 63., 38., nan, 27., 47
., 57.,
              24., 34., 66., 78., 66., 35., 62., 59., 57., 31., 73., 56., 58
., 59.,
              45., 60., 24., 58., 59., 78., 58., 39., 39., nan, 50., 54., na
n, 61.],
          dtype=torch.float64),
   'OS_MONTHS': tensor([       nan, 2.5300e+01, 8.0000e-01, 1.2300e+01,
1.4600e+01, 3.4300e+01,
              1.0000e-01, 1.0400e+01, 2.7000e+01, 1.7500e+01, 7.7000e+00,
1.0000e-01,
              1.0400e+01, 7.8000e+00, 3.1000e+00, 1.8100e+01, 1.7500e+01,
6.8000e+00,
              7.7900e+01, 1.1300e+01, 2.6700e+01, 2.1600e+01, 1.3070e+02,
5.8000e+00,
                       nan, 9.7000e+00, 1.5600e+01,        nan, 1.6000e+01,
8.3600e+01,
              1.0900e+01, 1.2000e+01,        nan, 2.3200e+01, 1.0000e-01,

```
          4.7000e+00,
          4.9000e+00, 4.1000e+00, 1.7100e+01, 2.7000e+01, 1.0800e+01,
   3.0000e-01,
          1.2600e+01, 1.0000e-01, 1.6000e+00, 1.5100e+01,         nan,
   6.9000e+00,
          1.5400e+01,         nan, 5.7000e+00, 9.3000e+00, 5.4000e+00,
   1.7900e+01,
                  nan, 1.7400e+01,         nan, 1.0000e-01, 4.8400e+01,
   1.0300e+01,
          2.3000e+00, 5.2000e+00, 5.1000e+00, 1.4200e+01, 9.4000e+00,
   1.5400e+01,
          2.0000e-01, 1.4000e+01,         nan, 1.8400e+01, 9.9000e+00,
   2.1500e+01,
          7.6000e+00, 1.9600e+01, 4.9400e+01, 6.2000e+00, 1.6100e+01,
   1.5000e+01,
          1.7500e+01, 4.1000e+00,         nan, 7.2000e+00, 2.1000e+00,
   6.0000e-01,
          1.6000e+00,         nan, 6.9000e+01, 3.0700e+01, 1.1300e+01,
   4.7000e+00,
          2.0000e-01,         nan, 3.2000e+00, 2.7000e+00, 1.0600e+01,
   nan,
          7.0000e+00,         nan, 1.2440e+02, 3.7000e+00,         nan,
   3.0000e+01,
          2.9000e+00, 1.4200e+01, 1.4000e+01, 1.9600e+01, 7.6000e+00,
   7.6000e+00,
          6.6000e+00, 5.8700e+01, 9.9000e+00, 1.2000e+00, 1.7300e+01,
   2.5000e+00,
          4.7000e+00, 7.9900e+01, 4.6000e+00,         nan, 5.0000e-01,
   6.4000e+00,
          5.9000e+00, 5.4000e+00, 1.1400e+01, 3.3600e+01, 5.0500e+01,
   8.0000e+00,
          1.8000e+00, 4.2000e+00, 1.8900e+01, 3.7100e+01,         nan,
   1.2700e+01,
          7.2500e+01, 2.3000e+00, 1.5400e+01, 6.7000e+00,         nan,
   1.6500e+01,
          1.6100e+01, 1.5610e+02, 7.8000e+00, 1.5100e+01, 3.2000e+00,
   3.9000e+00,
          7.9000e+00, 8.0000e+00,         nan, 1.3500e+01, 7.5000e+00,
   3.3000e+01,
          1.0200e+01, 4.1100e+01, 2.3800e+01, 1.6800e+01, 2.0000e-01,
   1.5000e+01,
          7.4000e+00,         nan, 2.2600e+01, 7.8200e+01, 3.1700e+01,
   1.5800e+01,
          2.1000e+01, 2.4000e+00, 1.9000e+00, 1.0000e-01, 8.0000e-01,
   8.9000e+00,
          3.4900e+01, 3.6500e+01, 9.0000e-01, 6.6000e+00, 9.4500e+01,
   1.5000e+00,
                  nan, 1.8900e+01, 2.5500e+01,         nan, 5.0100e+01,
   nan,
                  nan, 8.1000e+00, 1.3900e+01, 2.2400e+01,         nan,
   1.2050e+02,
          2.0000e-01, 3.7000e+01, 1.6200e+01, 2.5000e+00, 3.1800e+01,
   7.5100e+01,
          3.2000e+00, 1.9800e+01, 1.3600e+01, 3.5000e+00, 7.4000e+00,
   1.1500e+01,
          1.0790e+02, 1.9000e+00, 7.7000e+00, 4.0800e+01, 3.0300e+01,
   nan,
          4.6000e+00, 8.3000e+00,         nan, 0.0000e+00, 6.4600e+01,
   1.2900e+01,
          9.4600e+01, 1.4500e+01, 4.8000e+00, 2.7000e+00, 2.5000e+00,
```

```
1.8200e+01,
          4.2700e+01, 1.1500e+01, 3.4000e+00, 4.7000e+00, 3.7000e+00,
1.0900e+01,
          7.9000e+00, 6.8000e+00, 4.7000e+00, 5.0000e+00, 1.3370e+02,
6.0000e-01,
          5.4200e+01, 1.4700e+01, 2.0800e+01, 1.9200e+01, 1.2600e+01,
nan,
          2.2500e+01, 3.4300e+01,         nan, 5.7000e+00], dtype=torch.f
loat64),
   'OS_STATUS': tensor([nan, 1., 0., 1., 0., 0., 0., 0., 0., 1., 0., 0.,
1., 0., 1., 0., 0., 0.,
          0., 1., 1., 0., 1., 1., nan, 0., 0., nan, 0., 0., 0., 0., nan,
0., 0., 0.,
          1., 0., 1., 0., 0., 0., 0., 0., 0., 0., nan, 1., 1., nan, 1.,
0., 1., 0.,
          nan, 0., nan, 0., 0., 1., 0., 0., 1., 0., 0., 1., 0., 0., nan,
0., 0., 0.,
          1., 0., 0., 0., 0., 1., 0., 0., nan, 0., 1., 0., 0., nan, 0.,
0., 0., 1.,
          0., nan, 1., 1., 0., nan, 0., nan, 0., 0., nan, 1., 1., 0., 1.
, 0., 0., 0.,
          0., 0., 1., 1., 0., 1., 1., 1., 0., nan, 0., 0., 0., 0., 1., 1
., 1., 1.,
          1., 0., 0., 0., nan, 1., 0., 1., 0., 1., nan, 0., 0., 0., 0.,
0., 1., 0.,
          1., 0., nan, 0., 0., 0., 0., 1., 0., 1., 0., 1., 0., nan, 0.,
0., 0., 1.,
          1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 1., nan, 1., 1., n
an, 1., nan,
          nan, 0., 1., 1., nan, 1., 0., 0., 0., 1., 0., 1., 0., 1., 1.,
0., 1., 0.,
          0., 0., 0., 1., 0., nan, 0., 0., nan, 0., 0., 1., 0., 1., 1.,
1., 0., 0.,
          0., 1., 0., 0., 1., 0., 1., 0., 0., 0., 1., 0., 0., 1., 1., 0.
, 0., nan,
          0., 0., nan, 0.], dtype=torch.float64),
   'KARNOFSKY_PERFORMANCE_SCORE': tensor([ nan,  nan,  70.,  nan,  nan,
nan,  nan,  80.,  nan,  60.,  80.,  80.,
           80., 100.,  nan,  50.,  80.,  nan,  90.,  70.,  nan, 100.,  9
0.,  nan,
           nan,  40.,  70.,  nan,  nan,  nan, 100.,  90.,  nan,  nan,  n
an, 100.,
           60.,  90.,  80.,  nan,  nan,  80.,  nan,  nan,  80., 100.,  n
an,  60.,
           nan,  nan,  80.,  80.,  60.,  90.,  nan,  90.,  nan,  nan, 10
0.,  80.,
           80.,  60.,  80.,  70., 100.,  nan, 100.,  90.,  nan,  nan, 10
0.,  nan,
           80., 100.,  90.,  nan,  60., 100.,  90.,  80.,  nan,  80.,  n
an,  nan,
           60.,  nan, 100.,  90., 100.,  nan,  nan,  nan,  80.,  nan,  9
0.,  nan,
           90.,  nan, 100.,  nan,  nan,  80.,  40.,  90.,  80.,  nan,  9
0.,  80.,
           nan,  80., 100.,  nan,  nan,  80.,  80., 100.,  80.,  nan,  n
an,  90.,
           nan,  nan,  nan, 100.,  nan,  nan,  nan,  90.,  80.,  nan,  n
an,  nan,
          100.,  nan, 100.,  80.,  nan,  nan,  90.,  nan,  nan,  80.,  4
0.,  90.,
```

```
              60., 100.,  nan,  70.,  80.,  nan,  nan,  90.,  90.,  nan,  n
an,  nan,
              nan,  nan, 100.,  nan,  nan,  80.,  90.,  nan,  nan,  nan,  4
0.,  80.,
              80.,  nan,  60.,  nan,  90.,  nan,  nan,  80.,  90.,  nan, 10
0.,  nan,
              nan, 100.,  nan,  nan,  nan,  nan,  nan, 100.,  nan,  nan,  n
an,  90.,
              90., 100.,  nan, 100.,  nan,  nan,  nan,  40.,  nan,  90.,  9
0.,  nan,
              40.,  nan,  nan,  nan, 100.,  60., 100.,  80.,  60.,  60.,  n
an, 100.,
              nan,  80., 100.,  90.,  nan,  90., 100., 100.,  80.,  nan, 10
0.,  nan,
              90.,  80.,  nan,  nan,  80.,  nan, 100., 100.,  nan,  80.],
          dtype=torch.float64),
   'STUDY': tensor([0., 1., 0., 1., 0., 0., 1., 1., 0., 1., 1., 0., 1., 1
., 1., 0., 0., 1.,
              0., 1., 1., 0., 0., 1., 0., 1., 1., 0., 0., 0., 0., 0., 0., 0.
, 0., 0.,
              1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 1., 1.
, 1., 0.,
              0., 0., 0., 0., 0., 1., 1., 0., 1., 1., 0., 1., 0., 0., 0., 0.
, 0., 0.,
              1., 0., 0., 0., 0., 1., 1., 0., 0., 1., 1., 0., 1., 0., 0., 0.
, 1., 1.,
              0., 0., 1., 1., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 1., 0.
, 0., 0.,
              0., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 0., 0., 1., 0., 1.
, 1., 0.,
              1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 1., 0.
, 1., 0.,
              0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0.
, 0., 1.,
              0., 0., 1., 0., 1., 1., 1., 0., 1., 0., 0., 1., 0., 0., 0., 0.
, 0., 1.,
              0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 0., 0., 0., 1., 1., 0.
, 1., 1.,
              0., 0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 1., 0., 1., 1., 1.
, 1., 0.,
              0., 1., 0., 0., 1., 0., 0., 1., 0., 1., 0., 0., 0., 1., 1., 0.
, 0., 0.,
              0., 0., 0., 0.], dtype=torch.float64),
   'BCR_STATUS': tensor([1., 0., 1., 0., 0., 0., 0., 0., 1., 0., 0., 1.,
0., 0., 0., 0., 1., 0.,
              0., 0., 0., 1., 0., 0., 1., 0., 0., 1., 0., 1., 1., 1., 1., 1.
, 0., 1.,
              0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 1., 0., 0.
, 0., 0.,
              1., 0., 1., 0., 1., 0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 0.
, 1., 0.,
              0., 1., 1., 0., 0., 0., 0., 1., 1., 0., 0., 0., 0., 0., 1., 0.
, 0., 0.,
              0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.
, 0., 1.,
              0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 1., 0., 0., 0.
, 0., 1.,
              0., 1., 1., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1.
, 0., 1.,
              1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0.
```

```
        , 0., 0.,
                0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 1.
        , 0., 1.,
                1., 1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 1., 0., 0., 1.
        , 0., 0.,
                0., 1., 0., 1., 0., 1., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0.
        , 0., 1.,
                0., 0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.
        , 1., 1.,
                1., 1., 0., 0.], dtype=torch.float64),
        'HISTOLOGICAL_DIAGNOSIS': tensor([nan, 1., 3., 1., 3., 3., 1., 1., 2.,
        1., 1., 2., 1., 1., 1., 0., 3., 1.,
                2., 1., 1., 0., 2., 1., nan, 1., 1., nan, 0., 0., 0., 2., nan,
        3., 3., 3.,
                1., 3., 1., 0., 2., 3., 2., 2., 3., 3., nan, 1., 1., nan, 1.,
        1., 1., 0.,
                nan, 2., nan, 3., 0., 1., 1., 2., 1., 1., 3., 1., 2., 0., nan,
        0., 0., 2.,
                1., 0., 0., 3., 3., 1., 1., 3., nan, 1., 1., 0., 1., nan, 0.,
        0., 1., 1.,
                0., nan, 1., 1., 0., nan, 2., nan, 2., 2., nan, 1., 1., 0., 1.
        , 0., 0., 0.,
                0., 1., 1., 1., 0., 1., 1., 0., 3., nan, 3., 0., 0., 1., 3., 1
        ., 1., 0.,
                1., 0., 0., 3., nan, 0., 3., 1., 2., 2., nan, 2., 3., 3., 1.,
        0., 1., 2.,
                0., 3., nan, 3., 1., 3., 3., 0., 0., 2., 0., 1., 0., nan, 0.,
        2., 3., 1.,
                3., 0., 1., 3., 1., 1., 1., 2., 1., 3., 3., 1., nan, 0., 0., n
        an, 0., nan,
                nan, 3., 1., 3., nan, 1., 2., 1., 0., 1., 3., 2., 2., 1., 1.,
        3., 1., 1.,
                0., 3., 2., 2., 0., nan, 1., 2., nan, 3., 3., 1., 0., 1., 1.,
        1., 1., 2.,
                2., 1., 2., 0., 1., 0., 0., 1., 0., 1., 3., 2., 3., 1., 1., 0.
        , 2., nan,
                3., 3., nan, 0.], dtype=torch.float64),
        'SEX': tensor([nan, 1., 0., 1., 1., 1., 0., 1., 1., 0., 0., 0., 1., 0.
        , 1., 1., 1., 1.,
                0., 1., 0., 0., 0., 1., nan, 0., 1., nan, 0., 1., 1., 1., nan,
        0., 1., 0.,
                1., 1., 1., 1., 0., 1., 0., 1., 0., 0., nan, 1., 1., nan, 1.,
        1., 0., 1.,
                nan, 0., nan, 0., 0., 1., 0., 0., 1., 1., 0., 1., 0., 0., nan,
        0., 1., 0.,
                1., 1., 0., 1., 1., 1., 0., 1., nan, 1., 0., 1., 0., nan, 1.,
        1., 0., 1.,
                0., nan, 1., 1., 1., nan, 0., nan, 1., 1., nan, 0., 0., 1., 1.
        , 1., 0., 0.,
                1., 0., 0., 1., 0., 0., 0., 1., 1., nan, 1., 0., 1., 1., 1., 1
        ., 1., 1.,
                1., 0., 1., 0., nan, 1., 1., 0., 1., 0., nan, 1., 0., 0., 0.,
        1., 0., 0.,
                0., 0., nan, 1., 1., 0., 1., 1., 1., 0., 0., 0., 0., nan, 1.,
        1., 1., 1.,
                1., 1., 0., 0., 1., 0., 1., 1., 1., 1., 0., 0., nan, 0., 1., n
        an, 1., nan,
                nan, 0., 1., 1., nan, 0., 1., 1., 1., 1., 1., 0., 1., 1., 0.,
        1., 0., 0.,
                1., 0., 1., 1., 1., nan, 1., 1., nan, 1., 0., 1., 1., 0., 1.,
```

```
      1., 0., 0.,
               1., 0., 0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 0., 1., 0., 1.
       , 1., nan,
               1., 0., nan, 0.], dtype=torch.float64)})
```

In [11]:
```python
from IPython.display import display

# Get sorted lists of keys (clinical variable names) for both datasets
train_keys = sorted(train_dataset.ann.keys())
test_keys = sorted(test_dataset.ann.keys())

# Combine keys from both train and test (in case one is missing any)
all_keys = sorted(set(train_keys) | set(test_keys))

# Create a list to store table rows
rows = []
for key in all_keys:
    train_avail = "Yes" if key in train_dataset.ann else "No"
    test_avail = "Yes" if key in test_dataset.ann else "No"
    # Optionally include additional information like the tensor shape
    train_shape = train_dataset.ann[key].shape if key in train_dataset.an
    test_shape = test_dataset.ann[key].shape if key in test_dataset.ann e
    rows.append([key, train_avail, test_avail, train_shape, test_shape])

# Build a DataFrame from the rows
df = pd.DataFrame(rows, columns=["Clinical Variable", "Train", "Test", "T

# Display the DataFrame in the notebook
display(df)
```

| | Clinical Variable | Train | Test | Train Shape | Test Shape |
|---|---|---|---|---|---|
| **0** | AGE | Yes | Yes | (556,) | (238,) |
| **1** | BCR_STATUS | Yes | Yes | (556,) | (238,) |
| **2** | HISTOLOGICAL_DIAGNOSIS | Yes | Yes | (556,) | (238,) |
| **3** | KARNOFSKY_PERFORMANCE_SCORE | Yes | Yes | (556,) | (238,) |
| **4** | OS_MONTHS | Yes | Yes | (556,) | (238,) |
| **5** | OS_STATUS | Yes | Yes | (556,) | (238,) |
| **6** | SEX | Yes | Yes | (556,) | (238,) |
| **7** | STUDY | Yes | Yes | (556,) | (238,) |

- Make a histogram plot of the follow up times in months (OS_MONTHS) (use sns.histplot)

In [12]:
```python
import matplotlib.pyplot as plt

data1 = train_dataset.ann["OS_MONTHS"]
data2 = test_dataset.ann["OS_MONTHS"]

# Create a histogram plot of OS_MONTHS for both train and test data
plt.figure(figsize=(10, 6))
sns.histplot(data1, bins=30, kde=False, label="Train")
sns.histplot(data2, bins=30, kde=False, label="Test")
```

```
plt.xlabel("Follow-up Time (Months)")
plt.ylabel("Frequency")
plt.title("Histogram of Follow-up Times (OS_MONTHS)")
plt.legend()
plt.show()
```



- Make a histogram of the age distribution of the patients in the training data; facet the histogram by "SEX" variable (see flexynesis.utils.plot_boxplot)

In [13]:
```
cat_x = train_dataset.ann["SEX"]
num_y = train_dataset.ann["AGE"]
flexynesis.utils.plot_boxplot(cat_x, num_y, title_x = 'Categories', title
```

/home/thesamurai/micromamba/envs/flexynesisenv/lib/python3.11/site-package
s/flexynesis/utils.py:155: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be remove
d in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for
the same effect.

- Make a summary of all available clinical variables (see flexynesis.print_summary_stats)

In [14]: 
```python
flexynesis.print_summary_stats(train_dataset)
```

```
Summary for variable: AGE
Numerical Variable Summary: Median = 51.0, Mean = 50.33781190019194
------
Summary for variable: OS_MONTHS
Numerical Variable Summary: Median = 11.6, Mean = 19.090978886756236
------
Summary for variable: OS_STATUS
Numerical Variable Summary: Median = 0.0, Mean = 0.36153846153846153
------
Summary for variable: KARNOFSKY_PERFORMANCE_SCORE
Numerical Variable Summary: Median = 80.0, Mean = 82.45454545454545
------
Summary for variable: STUDY
Categorical Variable Summary:
  Label: Brain Lower Grade Glioma, Count: 353
  Label: Glioblastoma multiforme, Count: 203
------
Summary for variable: BCR_STATUS
Categorical Variable Summary:
  Label: IGC, Count: 454
  Label: NCH, Count: 102
------
Summary for variable: HISTOLOGICAL_DIAGNOSIS
Categorical Variable Summary:
  Label: astrocytoma, Count: 115
  Label: glioblastoma, Count: 201
  Label: oligoastrocytoma, Count: 79
  Label: oligodendroglioma, Count: 126
  Label: nan, Count: 35
------
Summary for variable: SEX
Categorical Variable Summary:
  Label: Female, Count: 209
  Label: Male, Count: 312
  Label: nan, Count: 35
------
```

- Notice that the categorical variables such as "SEX", "STUDY", "HISTOLOGICAL_DIAGNOSIS" are encoded numerically in the "dataset.ann" objects. Use dataset.label_mappings to map the STUDY variable to their original labels. Print the top 10 values in dataset.ann['STUDY'] and the mapped label values.

```
In [15]:  train_dataset.label_mappings
```

```
Out[15]:  {'STUDY': {0: 'Brain Lower Grade Glioma', 1: 'Glioblastoma multiforme'},
           'BCR_STATUS': {0: 'IGC', 1: 'NCH'},
           'HISTOLOGICAL_DIAGNOSIS': {0: 'astrocytoma',
            1: 'glioblastoma',
            2: 'oligoastrocytoma',
            3: 'oligodendroglioma',
            4: 'nan'},
           'SEX': {0: 'Female', 1: 'Male', 2: 'nan'}}
```

```
In [16]:  test_dataset.label_mappings
```

```
Out[16]: {'STUDY': {0: 'Brain Lower Grade Glioma', 1: 'Glioblastoma multiforme'},
          'BCR_STATUS': {0: 'IGC', 1: 'NCH'},
          'HISTOLOGICAL_DIAGNOSIS': {0: 'astrocytoma',
           1: 'glioblastoma',
           2: 'oligoastrocytoma',
           3: 'oligodendroglioma',
           4: nan},
          'SEX': {0: 'Female', 1: 'Male', 2: nan}}
```

```python
In [17]: import math
         study_vals = train_dataset.ann["STUDY"][:10]
         mapping = train_dataset.label_mappings["STUDY"]
         df = pd.DataFrame({
             "STUDY": study_vals.tolist(),
             "Mapped Label": [
                 mapping[int(x.item())] if not math.isnan(x.item()) else "NaN"
                 for x in study_vals
             ]
         })
         display(df)
```

|   | STUDY | Mapped Label |
|---|-------|--------------|
| **0** | 0.0 | Brain Lower Grade Glioma |
| **1** | 0.0 | Brain Lower Grade Glioma |
| **2** | 0.0 | Brain Lower Grade Glioma |
| **3** | 1.0 | Glioblastoma multiforme |
| **4** | 0.0 | Brain Lower Grade Glioma |
| **5** | 0.0 | Brain Lower Grade Glioma |
| **6** | 0.0 | Brain Lower Grade Glioma |
| **7** | 0.0 | Brain Lower Grade Glioma |
| **8** | 1.0 | Glioblastoma multiforme |
| **9** | 1.0 | Glioblastoma multiforme |

- Now, let's explore the data matrices. Make a PCA plot of the mutation data matrix and color the samples by "HISTOLOGICAL_DIAGNOSIS". See flexynesis.plot_dim_reduced function

First create a pandas data frame with the data matrix of interest with feature and sample names

> df = pd.DataFrame(train_dataset.dat['cna'], index = train_dataset.samples, columns= train_dataset.features['cna'])

Check the data frame contents

> df.head()

Make a PCA plot of CNA values using the labels from the STUDY variable

**Note**: if you couldn't map the labels above, you can also use train_dataset.dat['STUDY'] as labels

In [18]: 
```python
df = pd.DataFrame(train_dataset.dat['cna'], index = train_dataset.samples
```

In [19]: 
```python
df.head()
```

Out[19]:

| | SLC30A8 | ZNF273 | CLEC5A | AGL | KCNA5 | MIR603 | SNTI |
|---|---|---|---|---|---|---|---|
| TCGA-P5-A735 | -0.185735 | -0.722640 | -0.844373 | 0.389339 | -0.215060 | 0.479218 | -0.1935 |
| TCGA-S9-A6TV | -0.251503 | -0.705830 | -0.826243 | 0.386623 | -0.153882 | 0.544251 | -0.2603 |
| TCGA-HW-8322 | -0.256765 | -0.739450 | 0.841715 | -1.835223 | -0.193781 | 0.575646 | -0.2657 |
| TCGA-06-5415 | -0.238350 | 0.788125 | 0.785059 | 0.381190 | 1.359602 | -0.740713 | -0.2470 |
| TCGA-VM-A8CB | -0.262026 | -0.743652 | -0.867035 | -2.188329 | -0.148562 | 0.582374 | -0.2710 |

5 rows × 1237 columns

In [20]: 
```python
ds = train_dataset
```

In [21]: 
```python
f = 'STUDY'
labels = [ds.label_mappings[f][x] for x in ds.ann[f].numpy()]
```

In [22]: 
```python
flexynesis.plot_dim_reduced(df, labels, color_type = 'categorical', metho
```



PCA Scatter Plot with Colored Labels

Labels
- Brain Lower Grade Glioma
- Glioblastoma multiforme

- (Optional exercise ideas):
  - Make a PCA plot coloring the samples by HISTOLOGICAL_DIAGNOSIS, GENDER, or any other clinical variable
  - Repeat the same exercise on the mutation data matrix.

In [23]:
```python
from flexynesis import plot_dim_reduced

# List available clinical variables in the training dataset (from .ann)
clinical_vars = list(train_dataset.ann.keys())

# Prepare the data matrices for CNA and MUT
ds = train_dataset
df_cna = pd.DataFrame(ds.dat["cna"], index=ds.samples, columns=ds.feature
df_mut = pd.DataFrame(ds.dat["mut"], index=ds.samples, columns=ds.feature

def get_labels(variable):
    """
    For a given clinical variable, returns a tuple:
      (labels, color_type)
    If a mapping exists, labels are mapped (categorical);
    otherwise, raw values are returned (numerical).
    """
    vals = ds.ann[variable].numpy()
    if variable in ds.label_mappings:
        labels = [ds.label_mappings[variable][int(x.item())] if not math.
                  for x in ds.ann[variable]]
        color_type = "categorical"
    else:
        labels = vals  # Use raw numeric values
        color_type = "numerical"
    return labels, color_type
```

In [24]:
```python
# Loop over each clinical variable and generate PCA plots for both CNA an
for var in clinical_vars:
    labels, color_type = get_labels(var)

    # Plot PCA for CNA data
    print(f"PCA plot for CNA data colored by: {var}")
    fig1 = plot_dim_reduced(df_cna, labels=labels, color_type=color_type,
    # If fig is returned, display it:
    if fig1 is not None:
        fig1.show()
    else:
        plt1.show()

    # Plot PCA for MUT data
    print(f"PCA plot for MUT data colored by: {var}")
    fig2 = plot_dim_reduced(df_mut, labels=labels, color_type=color_type,
    # If fig is returned, display it:
    if fig2 is not None:
        fig2.show()
    else:
        plt2.show()
```

PCA plot for CNA data colored by: AGE

## PCA Scatter Plot with Numerical Labels



PCA plot for MUT data colored by: AGE

## PCA Scatter Plot with Numerical Labels



PCA plot for CNA data colored by: OS_MONTHS

## PCA Scatter Plot with Numerical Labels



PCA plot for MUT data colored by: OS_MONTHS

## PCA Scatter Plot with Numerical Labels



PCA plot for CNA data colored by: OS_STATUS

## PCA Scatter Plot with Numerical Labels



PCA plot for MUT data colored by: OS_STATUS

## PCA Scatter Plot with Numerical Labels



PCA plot for CNA data colored by: KARNOFSKY_PERFORMANCE_SCORE

## PCA Scatter Plot with Numerical Labels



PCA plot for MUT data colored by: KARNOFSKY_PERFORMANCE_SCORE

## PCA Scatter Plot with Numerical Labels



PCA plot for CNA data colored by: STUDY

## PCA Scatter Plot with Colored Labels



**Labels**
- Brain Lower Grade Glioma
- Glioblastoma multiforme

PCA plot for MUT data colored by: STUDY

## PCA Scatter Plot with Colored Labels



**Labels**
- Brain Lower Grade Glioma
- Glioblastoma multiforme

PCA plot for CNA data colored by: BCR_STATUS

PCA plot for MUT data colored by: BCR_STATUS



PCA plot for CNA data colored by: HISTOLOGICAL_DIAGNOSIS

## PCA Scatter Plot with Colored Labels



PCA plot for MUT data colored by: HISTOLOGICAL_DIAGNOSIS

## PCA Scatter Plot with Colored Labels



PCA plot for CNA data colored by: SEX

## PCA Scatter Plot with Colored Labels



PCA plot for MUT data colored by: SEX

## PCA Scatter Plot with Colored Labels



# 2. Training a single model using manually set hyperparameters

Now that we have familiarized ourselves with the dataset at hand, we can start building

models.

First we will do a single model training by manually setting hyperparameters. Based on the model performance, we will try modifying individual hyperparameters and build more and more models and see if we can improve model performance.

We will need to define the following components for starting a model training:

```
1. Split the train_dataset into train/validation
components
2. Define data loaders for both train and validation
splits
3. Define a pytorch-lightning trainer
4. Define a model with hyperparameters
5. Fit the model
```

In [25]:
```python
# randomly assign 80% of samples for training, 20% for validation
train_indices = random.sample(range(0, len(train_dataset)), int(len(train
val_indices = list(set(range(len(train_dataset))) - set(train_indices))
train_subset = train_dataset.subset(train_indices)
val_subset = train_dataset.subset(val_indices)

# define data loaders for train/validation splits
from torch.utils.data import DataLoader
train_loader = DataLoader(train_subset, batch_size=32, shuffle=True)
val_loader = DataLoader(val_subset, batch_size=32, shuffle = False)
```

Now, we need to define a model with manually set hyperparameters and a lightning-trainer fit the model.

**Notice**: Notice the callback we are passing to the trainer which enables us to plot the loss values as the training progresses.

In [26]:
```python
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 32, 'hidden_dim_factor': 2, 'lr': 0.001, 'super
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                     callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
```

`Trainer.fit` stopped: `max_epochs=20` reached.

While we can observe how well the model training went based on the "loss" values, we can also evaluate the model performance on test dataset

```
In [27]:   # evaluate the model performance on predicting the target variable
           flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```

Out[27]:

|   | method | var | variable_type | metric | value |
|---|--------|-----|---------------|--------|-------|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.792810 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.810673 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.587156 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.855286 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.750372 |

## 2.1 Exercise

- Now, repeat the above model training and evaluation by manually changing the hyperparameters (Try at least 5 different combinations)
- See if you can find a better hyperparameter combination that yields a better classification performance than the initial setup we provided.
- See the default hyperparameter ranges we use for Flexynesis here: https://github.com/BIMSBbioinfo/flexynesis/blob/69b92ca9370551e9fcc82a756cb42c72bef4a4b1/flexynesis/config.py#L7, but feel free to try outside these ranges too.

- Also try to observe the impact of the changing parameters on how the train/ validation loss curves change.

```
myparams = {'latent_dim': XX, 'hidden_dim_factor': XX,
'lr': XX, 'supervisor_hidden_dim': XX, 'epochs': XX}

model = flexynesis.DirectPred(config = myparams, dataset =
train_dataset, target_variables=['STUDY'])

trainer = pl.Trainer(max_epochs=myparams['epochs'],
default_root_dir="./", logger=False,
enable_checkpointing=False,

callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])

trainer.fit(model, train_loader, val_loader)

flexynesis.evaluate_wrapper("DirectPred",
model.predict(test_dataset), test_dataset)
```

In [28]:
```python
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 64, 'hidden_dim_factor': 0.2345, 'lr': 0.00123,
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                        callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```

HPO Step=1 out of 1
(latent_dim=64, hidden_dim_factor=0.2345, lr=0.00123, supervisor_hidden_dim=16, epochs=76)



`Trainer.fit` stopped: `max_epochs=76` reached.

Out[28]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.752288 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.782020 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.519515 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.855671 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.742671 |

In [29]:
```
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 128, 'hidden_dim_factor': 0.3465, 'lr': 0.01234
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                     callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```



HPO Step=1 out of 1
(latent_dim=128, hidden_dim_factor=0.3465, lr=0.01234, supervisor_hidden_dim=20, epochs=80)

`Trainer.fit` stopped: `max_epochs=80` reached.

Out[29]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.770588 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.786511 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.536965 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.859516 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.752307 |

In [30]:
```
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 96, 'hidden_dim_factor': 0.5674, 'lr': 0.00012,
```

```
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                        callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```
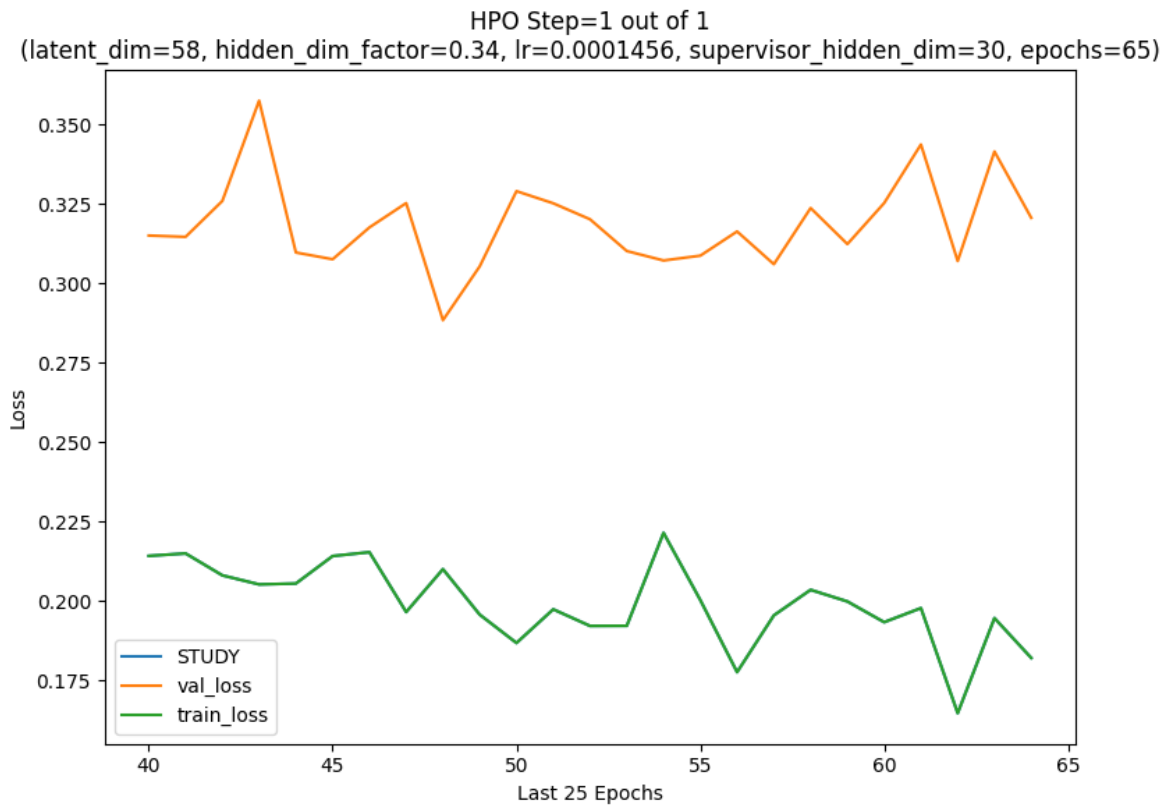
HPO Step=1 out of 1
(latent_dim=96, hidden_dim_factor=0.5674, lr=0.00012, supervisor_hidden_dim=24, epochs=25)



`Trainer.fit` stopped: `max_epochs=25` reached.

Out[30]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.788889 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.789152 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.553191 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.861669 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.762400 |

In [31]:
```
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 87, 'hidden_dim_factor': 0.4789, 'lr': 0.00134,
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                        callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```

HPO Step=1 out of 1
(latent_dim=87, hidden_dim_factor=0.4789, lr=0.00134, supervisor_hidden_dim=8, epochs=45)



`Trainer.fit` stopped: `max_epochs=45` reached.

Out[31]:

|   | method | var | variable_type | metric | value |
|---|--------|-----|---------------|--------|-------|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.760784 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.787405 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.532710 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.860054 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.742641 |

In [32]:
```python
# Define a model with manually set hyperparameters for the DirectPred mod
myparams = {'latent_dim': 58, 'hidden_dim_factor': 0.34, 'lr': 0.0001456,
model = flexynesis.DirectPred(config = myparams, dataset = train_dataset,
trainer = pl.Trainer(max_epochs=myparams['epochs'], default_root_dir="./"
                     callbacks=[flexynesis.LiveLossPlot(myparams, 1, 1)])
# Fit the model
trainer.fit(model, train_loader, val_loader)
flexynesis.evaluate_wrapper("DirectPred", model.predict(test_dataset), te
```

HPO Step=1 out of 1
(latent_dim=58, hidden_dim_factor=0.34, lr=0.0001456, supervisor_hidden_dim=30, epochs=65)

`Trainer.fit` stopped: `max_epochs=65` reached.

Out[32]:

|   | method | var | variable_type | metric | value |
|---|--------|-----|---------------|--------|-------|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.780392 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.798319 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.560784 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.869589 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.788800 |

**Warning!!**: In reality, we don't select the best models based on performance on the test dataset.

The best model is selected based on the validation loss value, where the model parameters that yields the lowest validation loss is selected to be the best model.

The validation dataset which we use to compute the validation loss is basically a subset of the training dataset.

# 3. Automating the Hyperparameter Optimisation Procedure

What we did in the above section was to set random hyperparameters, build a model, evaluate the model and try different hyperparameters based on our previous model performance. However, this process can be quite time consuming and arbitrary. This process can be automated using a Bayesian approach, where the model training is sequentially done for a number of hyperparameter optimisation iterations.

Now, we are ready to do a model training using hyperparameter optimisation.

- `model_class` : We pick `DirectPred` (a fully connected network) for now.
- `config_name` : We use the default/built-in hyperparameter search space for `DirectPred` class.
- `target_variables` : 'STUDY' variable contains the type of disease
- `n_iter` : We do 5 iterations of hyperparameter optimisation. For demonstration purposes, we set it to a small number.
- `plot_losses` : We want to visualize how the training progresses.
- `early_stop_patience` : If a training does not show any signs of improving the performance on the `validation` part of the `train_dataset` for at least 10 epochs, we stop the training. This not only significantly decreases the amount spent on training by avoiding unnecessary continuation of unpromising training runs, but also helps avoid over-fitting the network on the training data.

**Note 1**: Notice how the hyperparameters using in different HPO steps change at each iteration.

**Note 2**: Also notice that we are running the model for more epochs (500 by default) however, by using "early_stop_patience=10", we avoid lengthy training when validation performance is not improving.

**Note 3**: Try to follow the the loss curves and the used hyperparameters. See if you can spot which combination yields the lowest/best loss values.

**Warning!!**: In reality we need to set `n_iter` to higher values so that the optimizer can collect enough data points to learn trends in the parameter space.

```python
In [33]:  # Define a tuner; See n_iter is the number of
tuner = flexynesis.HyperparameterTuning(train_dataset,
                                        model_class = flexynesis.DirectPr
                                        config_name = "DirectPred",
                                        target_variables = ['STUDY'],
                                        n_iter=5,
                                        plot_losses=True,
                                        early_stop_patience=10)
### Perform Training
model, best_params = tuner.perform_tuning()
```

HPO Step=5 out of 5
(latent_dim=108, hidden_dim_factor=0.4784475635799447, lr=0.0019705139841641947, supervisor_hidden_dim=28, epochs=500, batch_size=128)



Validation: |

…

| Validate metric | DataLoader 0 |
|:---:|:---:|
| STUDY<br>val_loss | 0.41549500823020935<br>0.41549500823020935 |

Tuning Progress: 100%|

███████████████████████████████████████████████████████

████████████████████████| 5/5 [00:53<00:00, 10.75s/it, Iteration=5, B
est Loss=0.392]
[INFO] current best val loss: 0.39201247692108154; best params: {'latent_d
im': np.int64(88), 'hidden_dim_factor': 0.22455613724931636, 'lr': 0.00162
16946392416981, 'supervisor_hidden_dim': np.int64(32), 'epochs': 500, 'bat
ch_size': np.int64(64)} since 2 hpo iterations

In [34]: ## See which hyperparameter combination was the best
best_params

Out[34]: {'latent_dim': np.int64(88),
 'hidden_dim_factor': 0.22455613724931636,
 'lr': 0.0016216946392416981,
 'supervisor_hidden_dim': np.int64(32),
 'epochs': 15,
 'batch_size': np.int64(64)}

In [35]: ## Evaluate the model and visualising the results
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model.pred

Out[35]:

|   | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | STUDY | categorical | balanced_acc | 0.795425 |
| 1 | DirectPred | STUDY | categorical | f1_score | 0.811167 |
| 2 | DirectPred | STUDY | categorical | kappa | 0.589309 |
| 3 | DirectPred | STUDY | categorical | average_auroc | 0.869281 |
| 4 | DirectPred | STUDY | categorical | average_aupr | 0.769059 |

Let's extract the sample embeddings and make a PCA plot and color by the target
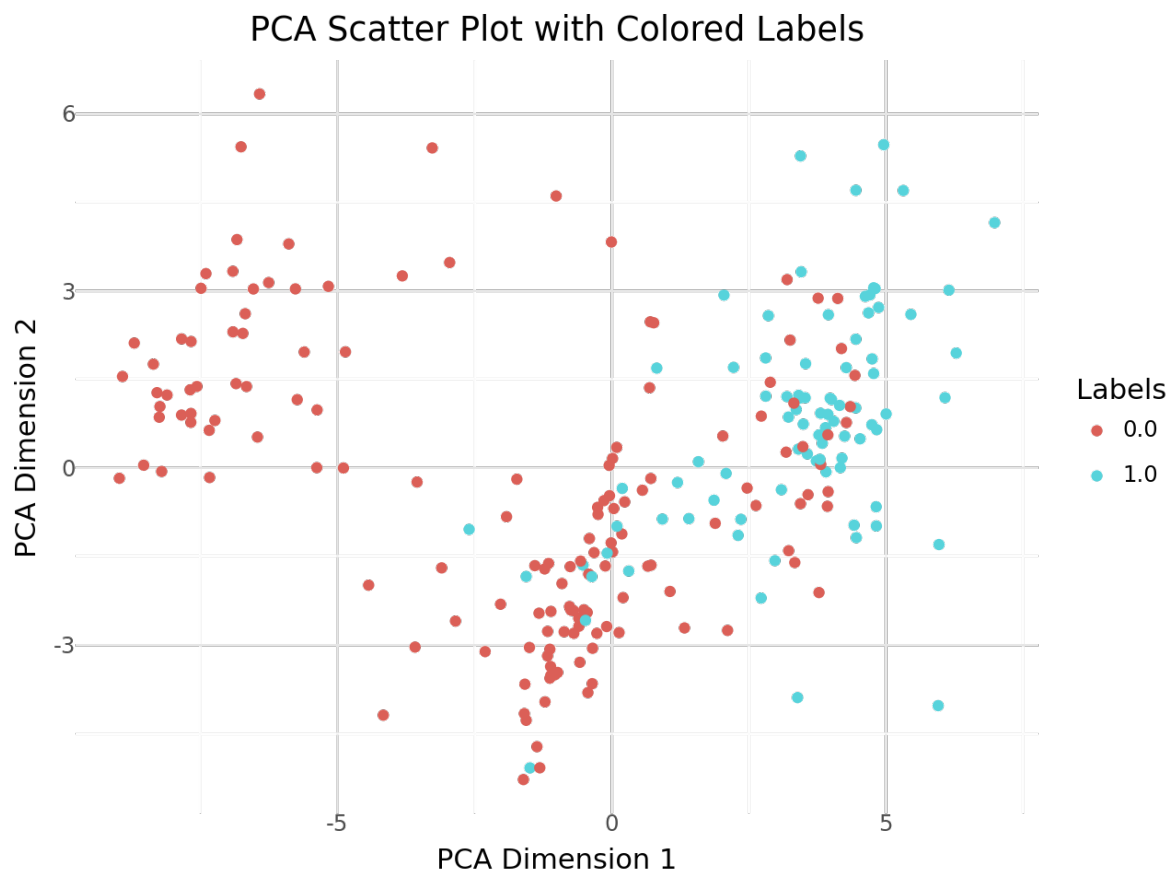
variable

```
In [36]: train_embeddings = model.transform(train_dataset)
         flexynesis.plot_dim_reduced(train_embeddings, train_dataset.ann['STUDY'])
```

PCA Scatter Plot with Colored Labels



Repeat the same for the test dataset: extract sample embeddings for test dataset samples and make a PCA plot, colored by "STUDY" variable

```
In [37]: test_embeddings = model.transform(test_dataset)
         flexynesis.plot_dim_reduced(test_embeddings, test_dataset.ann['STUDY'])
```

## PCA Scatter Plot with Colored Labels



## 3.1 Exercises

**Exercise 1**:

Look up what Harrell's C-index means and write down a simple description of what it measures.

Harrell's C-index is a metric used to evaluate the performance of survival models, similar in spirit to the area under the ROC curve.

It measures how well the model can correctly rank individuals based on their predicted risk or survival times. In simple terms, a high C-index indicates that, in most pairs of patients, the one predicted to have a higher risk (or shorter survival) indeed experiences the event (e.g., death) before the other.

**Exercise 2**:

Now, you build a model using hyperparameter tuning (run at least 10 HPO steps) to predict the survival outcomes of patients. Evaluate the final model on test dataset, which computes the "C-index".

Feel free to cheat from the tutorial available here: https://github.com/BIMSBbioinfo/flexynesis/blob/main/examples/tutorials/survival_subtypes_LGG_GBM.ipynb See how "OS_STATUS" and "OS_MONTHS" were used.

```
In [38]:  HPO_ITER = 10
```
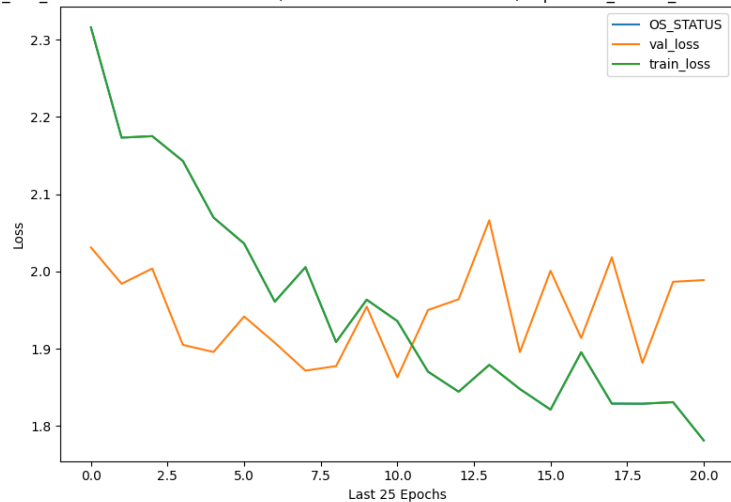
```
In [39]: tuner = flexynesis.HyperparameterTuning(train_dataset,
                                    model_class = flexynesis.DirectPr
                                    config_name = "DirectPred",
                                    surv_event_var="OS_STATUS",
                                    surv_time_var="OS_MONTHS",
                                    target_variables = [],
                                    n_iter=HPO_ITER,
                                    plot_losses=True,
                                    early_stop_patience=10)
         model, best_params = tuner.perform_tuning()
```

HPO Step=10 out of 10
(latent_dim=73, hidden_dim_factor=0.4460402072493238, lr=0.0001097797668246611, supervisor_hidden_dim=12, epochs=500, batch_size=32)



```
Validation: |
…
```

| Validate metric | DataLoader 0 |
|:---:|:---:|
| OS_STATUS | 1.9887990669240672 |
| val_loss | 1.9887990669240672 |

```
Tuning Progress: 100%|
```

```
                                                    | 10/10 [02:07<00:00, 12.73s/it, Iteration=10,
Best Loss=1.99]
[INFO] current best val loss: 1.9883608039858038; best params: {'latent_di
m': np.int64(85), 'hidden_dim_factor': 0.2819923823871969, 'lr': 0.0036760
824635741246, 'supervisor_hidden_dim': np.int64(22), 'epochs': 500, 'batch
_size': np.int64(32)} since 5 hpo iterations
```

```
In [40]: best_params
```

```
Out[40]: {'latent_dim': np.int64(85),
          'hidden_dim_factor': 0.2819923823871969,
          'lr': 0.0036760824635741246,
          'supervisor_hidden_dim': np.int64(22),
          'epochs': 17,
          'batch_size': np.int64(32)}
```

In [41]: ```
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model.pred
                                 surv_event_var=model.surv_event_var, surv_tim
```

Out[41]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| **0** | DirectPred | OS_STATUS | numerical | cindex | 0.671416 |

**Exercise 3:**

Again build a model using hyperparameter tuning to predict survival outcomes (as in Exercise 1), however, this time use additional clinical variables as targets.

```
    flexynesis.HyperparameterTuning(train_dataset,
                                         model_class =
    flexynesis.DirectPred,
                                         config_name =
    "DirectPred",

    surv_event_var="OS_STATUS",
                                         surv_time_var="OS_MONTHS",
                                         target_variables = [], =>
    What other variables can you use here? Try "AGE" and/or
    "HISTOLOGICAL_DIAGNOSIS" and see the model performance
                                         ...
```

**See if you can get a better C-index using additional target variables.**

In [42]:
```python
import itertools


# Exclude survival outcome variables from the targets.
excluded_vars = ["OS_MONTHS", "OS_STATUS"]
# All clinical variables available in train_dataset.ann (excluding surviv
clinical_targets = [var for var in train_dataset.ann.keys() if var not in

results = []  # to collect tuning results (model names and best parameter

def tune_model(target_vars):
    print(f"Tuning model for target variable(s): {target_vars}")
    tuner = flexynesis.HyperparameterTuning(
        train_dataset,
        model_class=flexynesis.DirectPred,
        config_name="DirectPred",
        surv_event_var="OS_STATUS",
        surv_time_var="OS_MONTHS",
        target_variables=target_vars,
        n_iter=HPO_ITER,
        plot_losses=True,
        early_stop_patience=10
    )
    model, best_params = tuner.perform_tuning()
    return best_params

# Generate combinations for r=1 (single) and r=2 (pair) only.
for r in range(1, 3):  # Only consider combinations of size 1 and 2.
    for combo in itertools.combinations(clinical_targets, r):
```

```
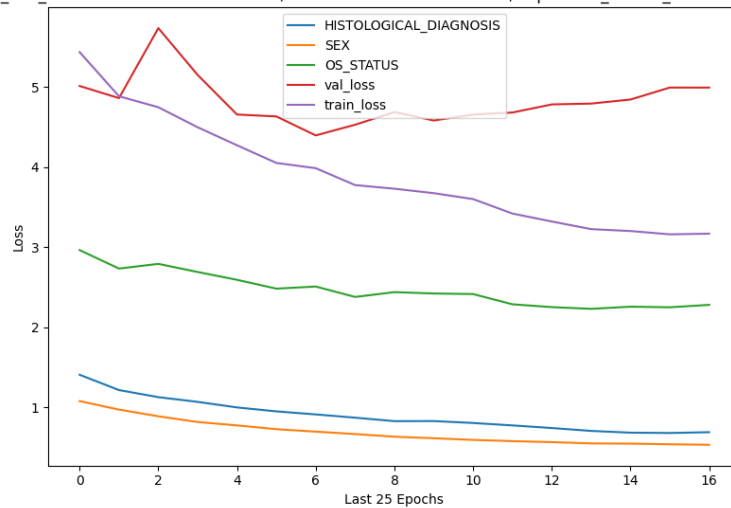        target_list = list(combo)
        best_params = tune_model(target_list)
        model_name = " + ".join(combo)
        results.append({
            "Model Name": model_name,
            "Best Parameters": best_params
        })

# Create and display the summary table with model names and their best pa
results_df = pd.DataFrame(results)
display(results_df)
```

HPO Step=10 out of 10
(latent_dim=44, hidden_dim_factor=0.4285500868749189, lr=0.002232991402453379, supervisor_hidden_dim=15, epochs=500, batch_size=64)



Validation: |
…

| Validate metric | DataLoader 0 |
|---|---|
| HISTOLOGICAL_DIAGNOSIS | 1.018038272857666 |
| OS_STATUS | 2.755630630630631 |
| SEX | 1.221535325050354 |
| val_loss | 4.995204375670837 |

Tuning Progress: 100%|

██████████████████████████████████████████████████████████

██████████████████████████████████████| 10/10 [03:39<00:00, 21.90s/it, Iteration=10,
Best Loss=3.65]
[INFO] current best val loss: 3.6494312341268125; best params: {'latent_di
m': np.int64(101), 'hidden_dim_factor': 0.3586039542763725, 'lr': 0.009336
67617651351, 'supervisor_hidden_dim': np.int64(30), 'epochs': 500, 'batch_
size': np.int64(32)} since 5 hpo iterations

| | Model Name | Best Parameters |
|---|---|---|
| 0 | AGE | {'latent_dim': 57, 'hidden_dim_factor': 0.2973... |
| 1 | KARNOFSKY_PERFORMANCE_SCORE | {'latent_dim': 43, 'hidden_dim_factor': 0.4154... |
| 2 | STUDY | {'latent_dim': 126, 'hidden_dim_factor': 0.315... |
| 3 | BCR_STATUS | {'latent_dim': 80, 'hidden_dim_factor': 0.3594... |
| 4 | HISTOLOGICAL_DIAGNOSIS | {'latent_dim': 70, 'hidden_dim_factor': 0.4874... |
| 5 | SEX | {'latent_dim': 95, 'hidden_dim_factor': 0.4662... |
| 6 | AGE + KARNOFSKY_PERFORMANCE_SCORE | {'latent_dim': 117, 'hidden_dim_factor': 0.357... |
| 7 | AGE + STUDY | {'latent_dim': 58, 'hidden_dim_factor': 0.4929... |
| 8 | AGE + BCR_STATUS | {'latent_dim': 27, 'hidden_dim_factor': 0.2155... |
| 9 | AGE + HISTOLOGICAL_DIAGNOSIS | {'latent_dim': 22, 'hidden_dim_factor': 0.3603... |
| 10 | AGE + SEX | {'latent_dim': 87, 'hidden_dim_factor': 0.3369... |
| 11 | KARNOFSKY_PERFORMANCE_SCORE + STUDY | {'latent_dim': 81, 'hidden_dim_factor': 0.4019... |
| 12 | KARNOFSKY_PERFORMANCE_SCORE + BCR_STATUS | {'latent_dim': 63, 'hidden_dim_factor': 0.4461... |
| 13 | KARNOFSKY_PERFORMANCE_SCORE + HISTOLOGICAL_DIA... | {'latent_dim': 104, 'hidden_dim_factor': 0.467... |
| 14 | KARNOFSKY_PERFORMANCE_SCORE + SEX | {'latent_dim': 110, 'hidden_dim_factor': 0.327... |
| 15 | STUDY + BCR_STATUS | {'latent_dim': 91, 'hidden_dim_factor': 0.2373... |
| 16 | STUDY + HISTOLOGICAL_DIAGNOSIS | {'latent_dim': 60, 'hidden_dim_factor': 0.4758... |
| 17 | STUDY + SEX | {'latent_dim': 34, 'hidden_dim_factor': 0.3133... |
| 18 | BCR_STATUS + HISTOLOGICAL_DIAGNOSIS | {'latent_dim': 111, 'hidden_dim_factor': 0.449... |
| 19 | BCR_STATUS + SEX | {'latent_dim': 61, 'hidden_dim_factor': 0.4455... |
| 20 | HISTOLOGICAL_DIAGNOSIS + SEX | {'latent_dim': 101, 'hidden_dim_factor': 0.358... |

```
In [43]: # Extract best parameters into a new DataFrame.
         param_rows = []
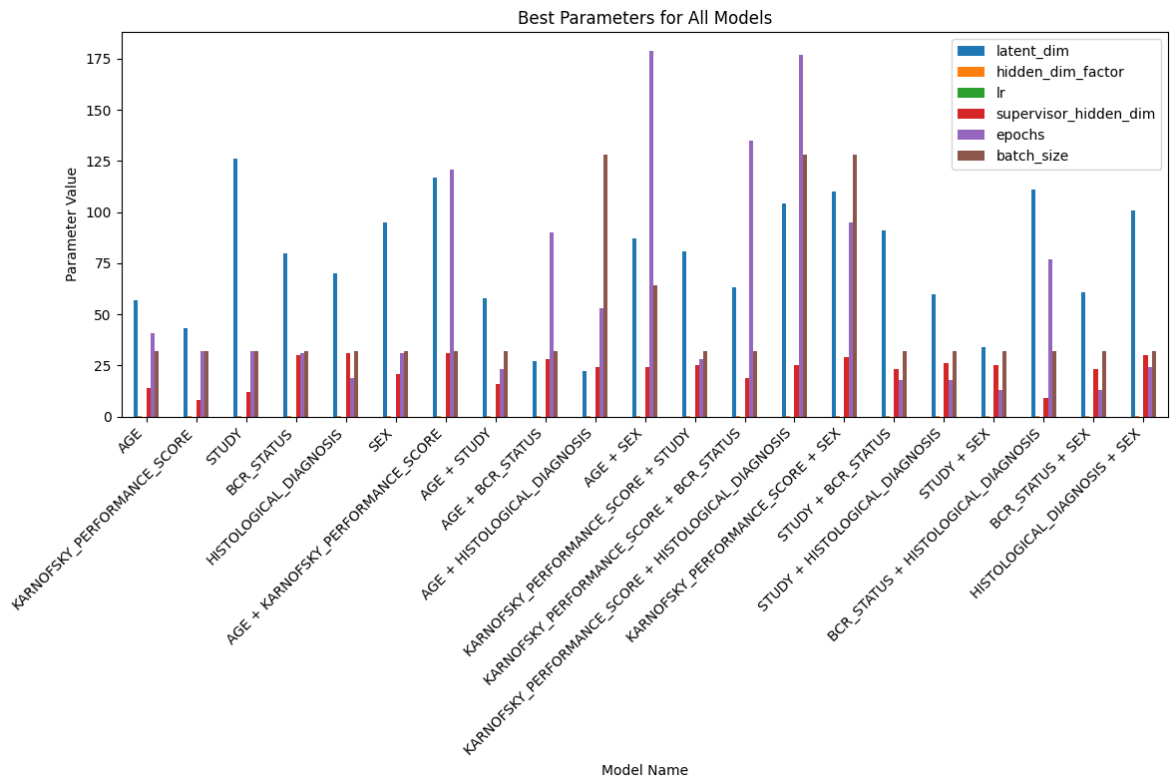         for idx, row in results_df.iterrows():
```

```
        model_name = row["Model Name"]
        best_params = row["Best Parameters"]
        # Create a new row dictionary with the model name and its best parame
        row_dict = {"Model Name": model_name}
        for param, value in best_params.items():
            row_dict[param] = value
        param_rows.append(row_dict)

params_df = pd.DataFrame(param_rows)
params_df = params_df.set_index("Model Name")
display(params_df)

# Visualize the best parameters for each model as a grouped bar chart.
# This assumes that all models share the same set of hyperparameter keys.
plt.figure(figsize=(12, 8))
params_df.plot(kind="bar", figsize=(12, 8))
plt.title("Best Parameters for All Models")
plt.xlabel("Model Name")
plt.ylabel("Parameter Value")
plt.xticks(rotation=45, ha="right")
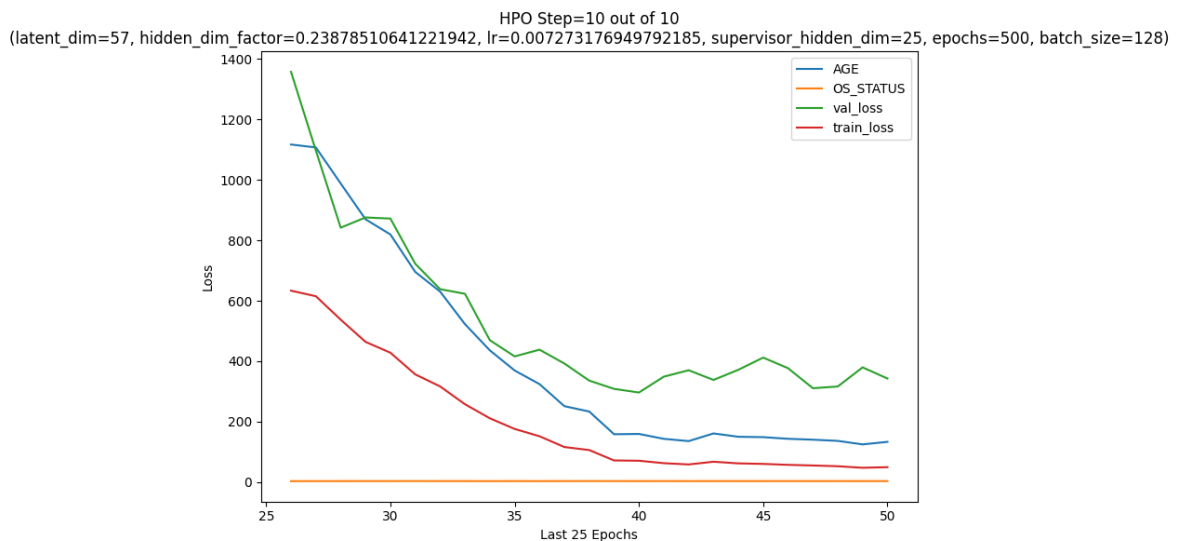plt.tight_layout()
plt.show()
```

| Model Name | latent_dim | hidden_dim_factor | lr | supervisor |
|---|---|---|---|---|
| AGE | 57 | 0.297327 | 0.004124 | |
| KARNOFSKY_PERFORMANCE_SCORE | 43 | 0.415413 | 0.009521 | |
| STUDY | 126 | 0.315875 | 0.000129 | |
| BCR_STATUS | 80 | 0.359499 | 0.000431 | |
| HISTOLOGICAL_DIAGNOSIS | 70 | 0.487421 | 0.000749 | |
| SEX | 95 | 0.466202 | 0.000101 | |
| AGE + KARNOFSKY_PERFORMANCE_SCORE | 117 | 0.357150 | 0.000788 | |
| AGE + STUDY | 58 | 0.492938 | 0.008728 | |
| AGE + BCR_STATUS | 27 | 0.215583 | 0.000928 | |
| AGE + HISTOLOGICAL_DIAGNOSIS | 22 | 0.360345 | 0.009413 | |
| AGE + SEX | 87 | 0.336906 | 0.000910 | |
| KARNOFSKY_PERFORMANCE_SCORE + STUDY | 81 | 0.401936 | 0.005484 | |
| KARNOFSKY_PERFORMANCE_SCORE + BCR_STATUS | 63 | 0.446120 | 0.001121 | |
| KARNOFSKY_PERFORMANCE_SCORE + HISTOLOGICAL_DIAGNOSIS | 104 | 0.467163 | 0.002526 | |
| KARNOFSKY_PERFORMANCE_SCORE + SEX | 110 | 0.327746 | 0.005415 | |
| STUDY + BCR_STATUS | 91 | 0.237357 | 0.000764 | |
| STUDY + HISTOLOGICAL_DIAGNOSIS | 60 | 0.475891 | 0.003096 | |
| STUDY + SEX | 34 | 0.313385 | 0.002520 | |
| BCR_STATUS + HISTOLOGICAL_DIAGNOSIS | 111 | 0.449693 | 0.000209 | |
| BCR_STATUS + SEX | 61 | 0.445594 | 0.002500 | |
| HISTOLOGICAL_DIAGNOSIS + SEX | 101 | 0.358604 | 0.009337 | |

```
<Figure size 1200x800 with 0 Axes>
```

Best Parameters for All Models



```
In [91]:  tuner = flexynesis.HyperparameterTuning(train_dataset,
                                     model_class = flexynesis.DirectPr
                                     config_name = "DirectPred",
                                     surv_event_var="OS_STATUS",
                                     surv_time_var="OS_MONTHS",
                                     target_variables = ["AGE"],
                                     n_iter=HPO_ITER,
                                     plot_losses=True,
                                     early_stop_patience=10)
          model_age, best_params_age = tuner.perform_tuning()
```

HPO Step=10 out of 10
(latent_dim=57, hidden_dim_factor=0.23878510641221942, lr=0.007273176949792185, supervisor_hidden_dim=25, epochs=500, batch_size=128)



Validation: |
…

| Validate metric | DataLoader 0 |
|---|---|
| AGE | 339.4956359863281 |
| OS_STATUS | 3.1707317073170738 |
| val_loss | 342.66636769364516 |

```
Tuning Progress: 100%|
                                                                                    | 10/10 [14:13<00:00, 85.38s/it, Iteration=10,
Best Loss=142]
[INFO] current best val loss: 141.8916905439046; best params: {'latent_di
m': np.int64(90), 'hidden_dim_factor': 0.40477655101031473, 'lr': 0.008136
437121084042, 'supervisor_hidden_dim': np.int64(17), 'epochs': 500, 'batch
_size': np.int64(32)} since 5 hpo iterations
```

In [92]: `best_params_age`

Out[92]: 
```
{'latent_dim': np.int64(90),
 'hidden_dim_factor': 0.40477655101031473,
 'lr': 0.008136437121084042,
 'supervisor_hidden_dim': np.int64(17),
 'epochs': 36,
 'batch_size': np.int64(32)}
```

In [93]: 
```python
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model_age.
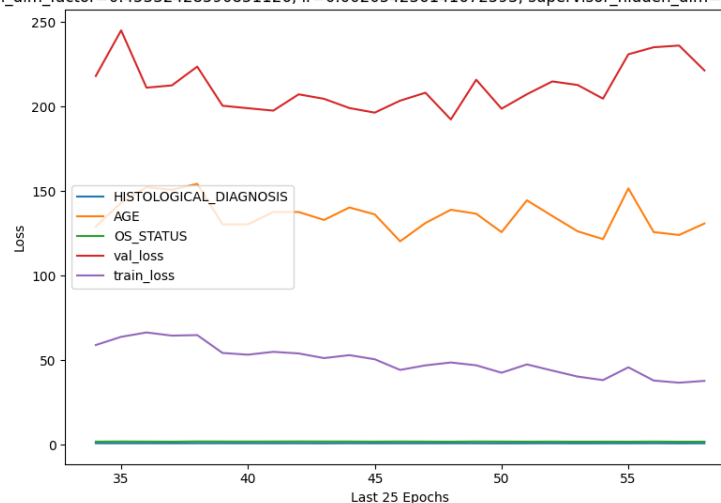                            surv_event_var=model_age.surv_event_var, surv
```

Out[93]:

|   | method | var | variable_type | metric | value |
|---|--------|-----|---------------|--------|-------|
| 0 | DirectPred | AGE | numerical | mse | 189.384308 |
| 1 | DirectPred | AGE | numerical | r2 | 0.280345 |
| 2 | DirectPred | AGE | numerical | pearson_corr | 0.529476 |
| 3 | DirectPred | OS_STATUS | numerical | cindex | 0.735526 |

In [94]: 
```python
tuner = flexynesis.HyperparameterTuning(train_dataset,
                                        model_class = flexynesis.DirectPr
                                        config_name = "DirectPred",
                                        surv_event_var="OS_STATUS",
                                        surv_time_var="OS_MONTHS",
                                        target_variables = ["HISTOLOGICAL
                                        n_iter=HPO_ITER,
                                        plot_losses=True,
                                        early_stop_patience=10)
model_histo_age, best_params_histo_age = tuner.perform_tuning()
```

HPO Step=10 out of 10
(latent_dim=35, hidden_dim_factor=0.45332428590851126, lr=0.0020542361416672593, supervisor_hidden_dim=27, epochs=500, batch_size=32)



```
Validation: |
…
```

| Validate metric | DataLoader 0 |
|---|---|
| AGE | 218.32479858398438 |
| HISTOLOGICAL_DIAGNOSIS | 0.8972963094711304 |
| OS_STATUS | 2.207113363363363 |
| val_loss | 221.429224146021 |

```
Tuning Progress: 100%|
████████████████████████████████████████████████████████████
███████████████████████| 10/10 [12:58<00:00, 77.87s/it, Iteration=10,
Best Loss=195]
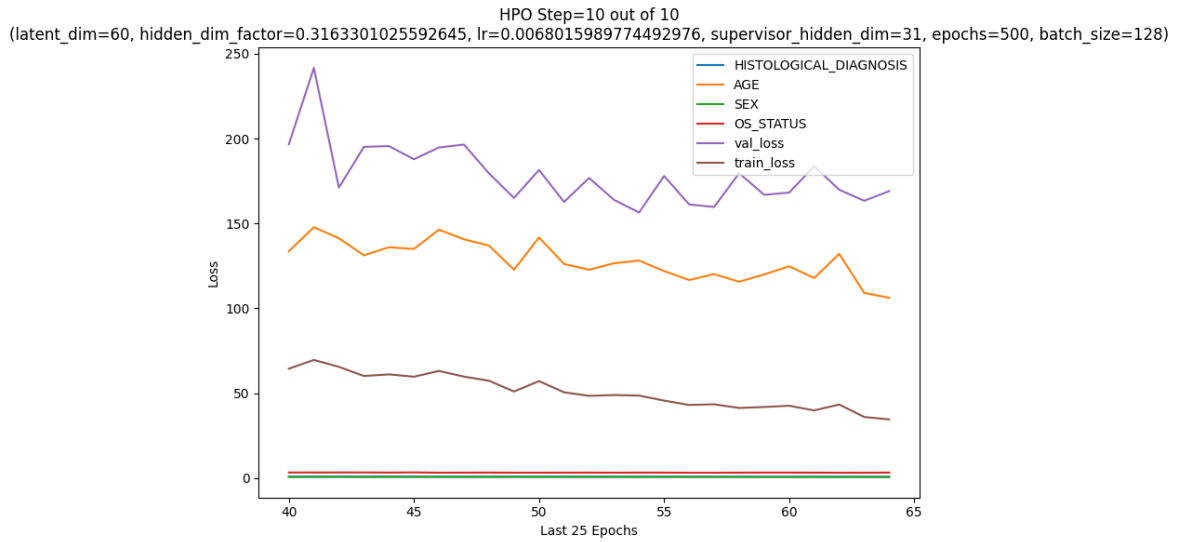[INFO] current best val loss: 194.86954287967168; best params: {'latent_di
m': np.int64(100), 'hidden_dim_factor': 0.45326239118533657, 'lr': 0.00756
9812865018211, 'supervisor_hidden_dim': np.int64(29), 'epochs': 500, 'batc
h_size': np.int64(64)} since 6 hpo iterations
```

In [95]:
```python
best_params_histo_age
```

Out[95]:
```
{'latent_dim': np.int64(100),
 'hidden_dim_factor': 0.45326239118533657,
 'lr': 0.007569812865018211,
 'supervisor_hidden_dim': np.int64(29),
 'epochs': 35,
 'batch_size': np.int64(64)}
```

In [96]:
```python
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model_hist
                            surv_event_var=model_histo_age.surv_event_var
```

Out[96]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | balanced_acc | 0.562536 |
| 1 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | f1_score | 0.594229 |
| 2 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | kappa | 0.490923 |
| 3 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | average_auroc | NaN |
| 4 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | average_aupr | NaN |
| 5 | DirectPred | AGE | numerical | mse | 212.005554 |
| 6 | DirectPred | AGE | numerical | r2 | 0.251678 |
| 7 | DirectPred | AGE | numerical | pearson_corr | 0.501675 |
| 8 | DirectPred | OS_STATUS | numerical | cindex | 0.759876 |

In [97]:
```python
tuner = flexynesis.HyperparameterTuning(train_dataset,
                                        model_class = flexynesis.DirectPr
                                        config_name = "DirectPred",
                                        surv_event_var="OS_STATUS",
                                        surv_time_var="OS_MONTHS",
                                        target_variables = ["HISTOLOGICAL
                                        n_iter=HPO_ITER,
                                        plot_losses=True,
                                        early_stop_patience=10)
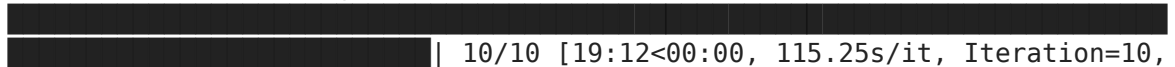model_histo_age_sex, best_params_histo_age_sex = tuner.perform_tuning()
```

HPO Step=10 out of 10
(latent_dim=60, hidden_dim_factor=0.3163301025592645, lr=0.0068015989774492976, supervisor_hidden_dim=31, epochs=500, batch_size=128)

Validation: |
…

| Validate metric | DataLoader 0 |
|---|---|
| AGE | 164.3182830810547 |
| HISTOLOGICAL_DIAGNOSIS | 0.7894374132156372 |
| OS_STATUS | 3.1578947368421053 |
| SEX | 0.852466881275177 |
| val_loss | 169.11808455617805 |

```
Tuning Progress: 100%|
███████████████████████████████████████████████████
███████████████████████| 10/10 [19:12<00:00, 115.25s/it, Iteration=10,
Best Loss=169]
[INFO] current best val loss: 169.11808455617805; best params: {'latent_di
m': np.int64(60), 'hidden_dim_factor': 0.3163301025592645, 'lr': 0.0068015
989774492976, 'supervisor_hidden_dim': np.int64(31), 'epochs': 500, 'batch
_size': np.int64(128)} since 0 hpo iterations
```

In [98]: `best_params_histo_age_sex`

Out[98]: 
```
{'latent_dim': np.int64(60),
 'hidden_dim_factor': 0.3163301025592645,
 'lr': 0.0068015989774492976,
 'supervisor_hidden_dim': np.int64(31),
 'epochs': 64,
 'batch_size': np.int64(128)}
```

In [99]: 
```
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model_hist
                            surv_event_var=model_histo_age_sex.surv_event
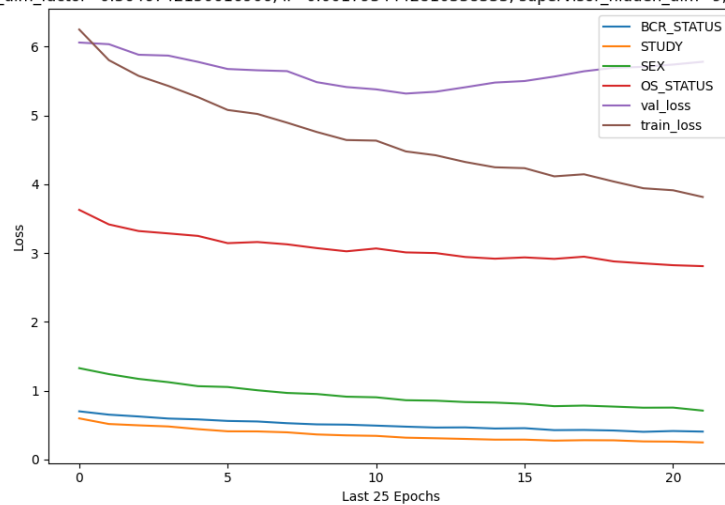```

Out[99]:

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | balanced_acc | 0.513490 |
| 1 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | f1_score | 0.559774 |
| 2 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | kappa | 0.415131 |
| 3 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | average_auroc | NaN |
| 4 | DirectPred | HISTOLOGICAL_DIAGNOSIS | categorical | average_aupr | NaN |
| 5 | DirectPred | AGE | numerical | mse | 206.294189 |
| 6 | DirectPred | AGE | numerical | r2 | 0.266506 |
| 7 | DirectPred | AGE | numerical | pearson_corr | 0.516242 |
| 8 | DirectPred | SEX | categorical | balanced_acc | 0.527388 |
| 9 | DirectPred | SEX | categorical | f1_score | 0.518168 |
| 10 | DirectPred | SEX | categorical | kappa | 0.059810 |
| 11 | DirectPred | SEX | categorical | average_auroc | NaN |
| 12 | DirectPred | SEX | categorical | average_aupr | NaN |
| 13 | DirectPred | OS_STATUS | numerical | cindex | 0.743298 |

In [120…
```python
tuner = flexynesis.HyperparameterTuning(train_dataset,
                                        model_class = flexynesis.DirectPr
                                        config_name = "DirectPred",
                                        surv_event_var="OS_STATUS",
                                        surv_time_var="OS_MONTHS",
                                        target_variables = ["BCR_STATUS",
                                        n_iter=HPO_ITER,
                                        plot_losses=True,
                                        early_stop_patience=10)
model_bcr_study_sex, best_params_bcr_study_sex = tuner.perform_tuning()
```

HPO Step=10 out of 10
(latent_dim=92, hidden_dim_factor=0.3646742150616906, lr=0.0017934442810338333, supervisor_hidden_dim=9, epochs=500, batch_size=128)



Validation: |
…

| Validate metric | DataLoader 0 |
|:---:|:---:|
| BCR_STATUS | 0.41967713832855225 |
| OS_STATUS | 3.426470588235294 |
| SEX | 1.568159818649292 |
| STUDY | 0.3639817535877228 |
| val_loss | 5.778289149789249 |

```
 Progress:  90%|
████████████████████████████████████████████████████████████████
████████████████████         | 9/10 [03:10<00:17, 17.11s/it, Iteration=10,
Best Loss=3.3]

Tuning Progress: 100%|
████████████████████████████████████████████████████████████████
██████████████████████| 10/10 [03:10<00:00, 19.01s/it, Iteration=10,
Best Loss=3.3]
[INFO] current best val loss: 3.2993653112592036; best params: {'latent_di
m': np.int64(24), 'hidden_dim_factor': 0.23022943258822356, 'lr': 0.006074
57995389591, 'supervisor_hidden_dim': np.int64(28), 'epochs': 500, 'batch_
size': np.int64(32)} since 7 hpo iterations
```

In [121…   `best_params_bcr_study_sex`

Out[121…
```
{'latent_dim': np.int64(24),
 'hidden_dim_factor': 0.23022943258822356,
 'lr': 0.00607457995389591,
 'supervisor_hidden_dim': np.int64(28),
 'epochs': 19,
 'batch_size': np.int64(32)}
```

In [122…
```
flexynesis.evaluate_wrapper(method = 'DirectPred', y_pred_dict=model_bcr_
                            surv_event_var=model_bcr_study_sex.surv_event
```

Out[122…

| | method | var | variable_type | metric | value |
|---|---|---|---|---|---|
| 0 | DirectPred | BCR_STATUS | categorical | balanced_acc | 0.504757 |
| 1 | DirectPred | BCR_STATUS | categorical | f1_score | 0.624258 |
| 2 | DirectPred | BCR_STATUS | categorical | kappa | 0.012116 |
| 3 | DirectPred | BCR_STATUS | categorical | average_auroc | 0.631254 |
| 4 | DirectPred | BCR_STATUS | categorical | average_aupr | 0.356507 |
| 5 | DirectPred | STUDY | categorical | balanced_acc | 0.739869 |
| 6 | DirectPred | STUDY | categorical | f1_score | 0.769609 |
| 7 | DirectPred | STUDY | categorical | kappa | 0.492617 |
| 8 | DirectPred | STUDY | categorical | average_auroc | 0.842753 |
| 9 | DirectPred | STUDY | categorical | average_aupr | 0.669224 |
| 10 | DirectPred | SEX | categorical | balanced_acc | 0.506929 |
| 11 | DirectPred | SEX | categorical | f1_score | 0.513457 |
| 12 | DirectPred | SEX | categorical | kappa | 0.014689 |
| 13 | DirectPred | SEX | categorical | average_auroc | NaN |
| 14 | DirectPred | SEX | categorical | average_aupr | NaN |
| 15 | DirectPred | OS_STATUS | numerical | cindex | 0.712861 |

## 3.2 Survival-risk subtypes

Use the best model from the above exercises to inspect sample embeddings categorized by survival risk scores.

Let's group the samples by predicted survival risk scores into 2 groups and visualize the sample embeddings colored by risk subtypes.

**Notice**: You can use the code-below to get survival risk groups, however, notice that you must have built a model with "OS_STATUS" already.

In [123…
```python
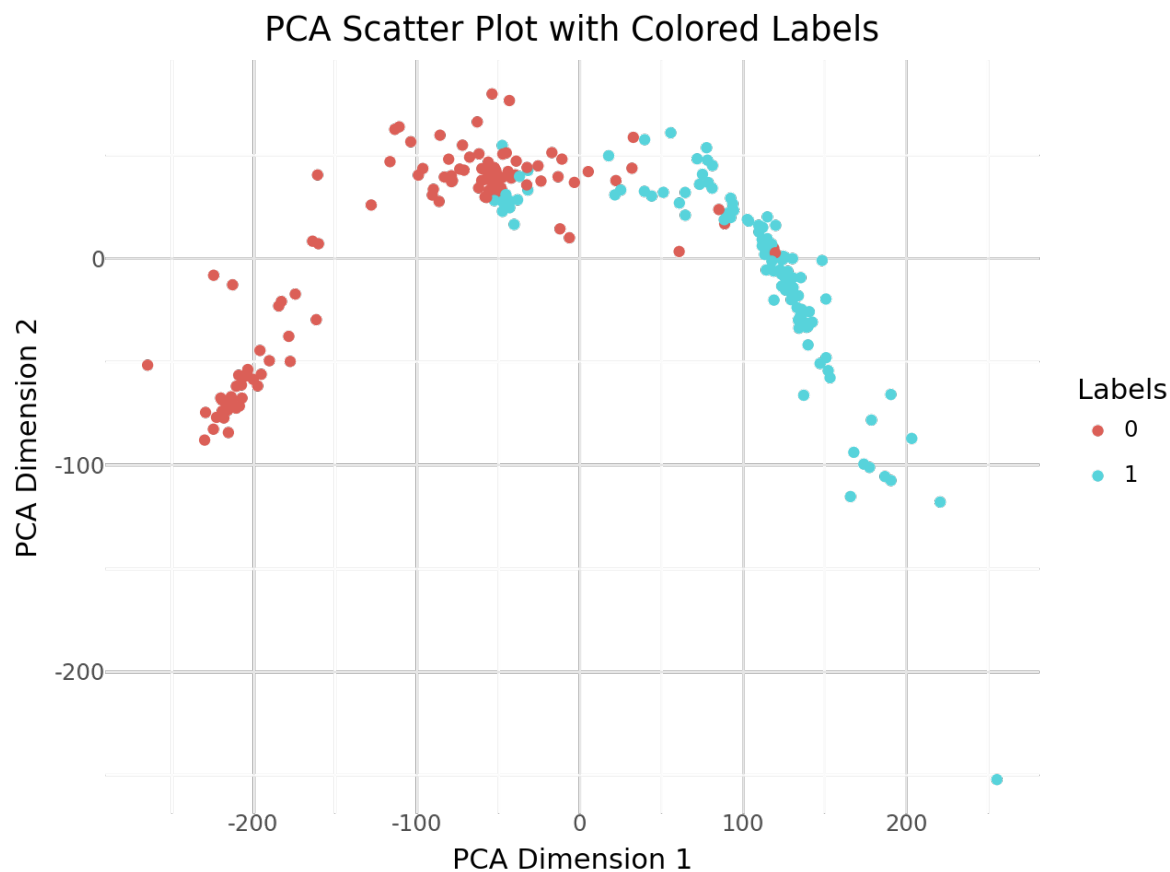# get model outputs for survival variable
model = model_histo_age
outputs = model.predict(test_dataset)['OS_STATUS'].flatten()
risk_scores = np.exp(outputs)
# Define quantile thresholds
quantiles = np.quantile(risk_scores, [0.5])
# Assign groups based on quantiles
groups = np.digitize(risk_scores, quantiles)
```

In [124…
```python
# Extract sample embeddings
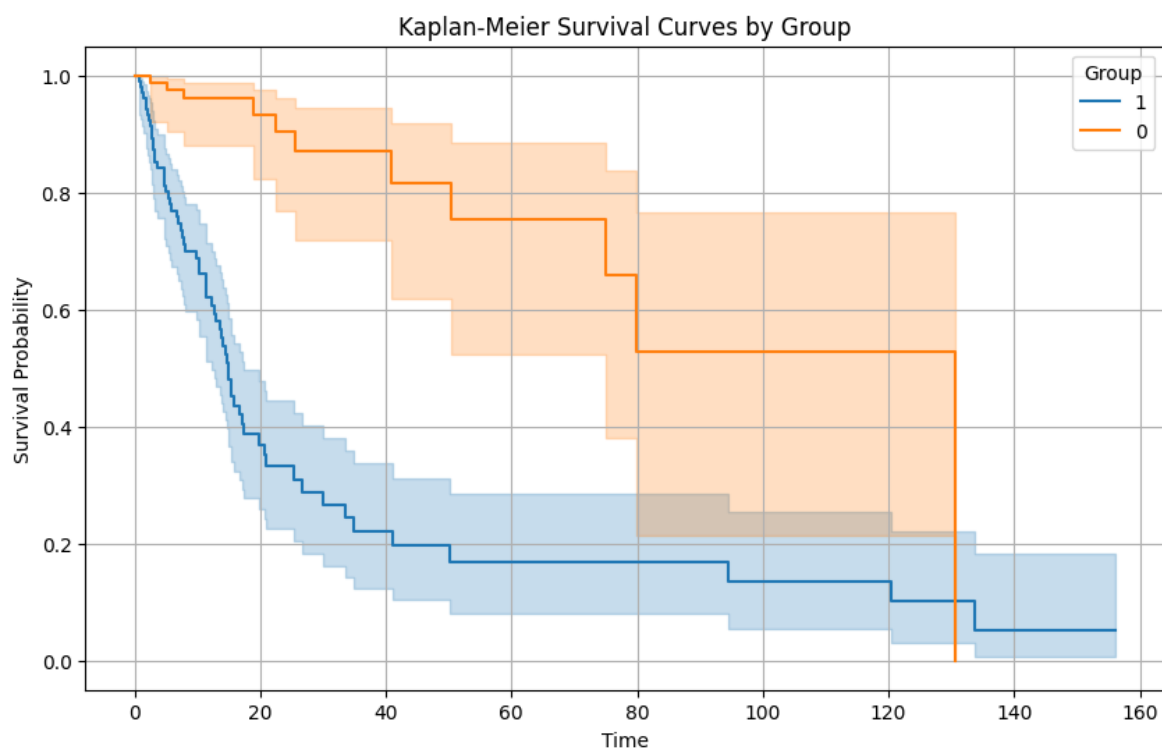E = model.transform(test_dataset)
```

In [125…
```python
flexynesis.plot_dim_reduced(E, groups)
```

## PCA Scatter Plot with Colored Labels



Let's also see the Kaplan Meier Curves of the risk subtypes

```
In [126…   # remove samples with NA values first
           durations = test_dataset.ann['OS_MONTHS']
           events = test_dataset.ann['OS_STATUS']
           valid_indices = ~torch.isnan(durations) & ~torch.isnan(events)
```

```
In [127…   flexynesis.plot_kaplan_meier_curves(durations[valid_indices], events[vali
```

## Finding survival-associated markers

We can also compute feature importance scores for prediction of overall survival.

```
In [128…  model.compute_feature_importance(train_dataset, 'OS_STATUS')
```

```
In [129…  # get top 10 features
          flexynesis.get_important_features(model, var = 'OS_STATUS', top=10)
```

Out[129…

|   | target_variable | target_class | target_class_label | layer | name | importance |
|---|---|---|---|---|---|---|
| 0 | OS_STATUS | 0 | | mut | IDH1 | 0.579877 |
| 1 | OS_STATUS | 0 | | mut | ATRX | 0.302072 |
| 2 | OS_STATUS | 0 | | mut | CIC | 0.042751 |
| 3 | OS_STATUS | 0 | | mut | IDH2 | 0.027195 |
| 4 | OS_STATUS | 0 | | mut | COL6A3 | 0.026273 |
| 5 | OS_STATUS | 0 | | mut | TEKT4 | 0.026143 |
| 6 | OS_STATUS | 0 | | mut | PIK3CA | 0.021109 |
| 7 | OS_STATUS | 0 | | mut | RELN | 0.020368 |
| 8 | OS_STATUS | 0 | | mut | TP53 | 0.020182 |
| 9 | OS_STATUS | 0 | | mut | BRD3 | 0.019807 |

## Comparing top markers with clinical covariates

Let's build a linear Cox-PH model including the top 5 markers and other clinical variables such as histological diagnosis, disease type (STUDY), age, and sex.

```
In [130…  # define a data.frame with clinical covariates and top markers along with
          vars = ['AGE', 'SEX', 'HISTOLOGICAL_DIAGNOSIS', 'STUDY', 'OS_MONTHS', 'OS_
          # read clinical variables
          df_clin = pd.concat(
              [pd.DataFrame({x: train_dataset.ann[x] for x in vars}, index=train_da
               pd.DataFrame({x: test_dataset.ann[x] for x in vars}, index=test_data
              axis = 0)
          # get top 5 survival markers and extract the input data for these markers
          imp = flexynesis.get_important_features(model, var = 'OS_STATUS', top=5)
          df_imp = pd.concat([train_dataset.get_feature_subset(imp), test_dataset.g

          # combine markers with clinical variables
          df = pd.concat([df_imp, df_clin], axis = 1)
          # remove samples without survival endpoints
          df = df[df['OS_STATUS'].notna()]
          df
```

| Out[130… | | mut_IDH1 | mut_ATRX | mut_CIC | mut_IDH2 | mut_COL6A3 | AGE | SEX | H |
|---|---|---|---|---|---|---|---|---|---|
| TCGA-S9-A6TV | | 0.982173 | 1.707482 | -0.344546 | -0.148522 | -0.177595 | 50.0 | 1.0 | |
| TCGA-HW-8322 | | 0.982173 | -0.585658 | -0.344546 | -0.148522 | -0.177595 | 39.0 | 1.0 | |
| TCGA-06-5415 | | -1.018150 | -0.585658 | -0.344546 | -0.148522 | -0.177595 | 60.0 | 1.0 | |
| TCGA-VM-A8CB | | 0.982173 | -0.585658 | 2.902366 | -0.148522 | -0.177595 | 33.0 | 1.0 | |
| TCGA-HT-7860 | | -1.018150 | -0.585658 | -0.344546 | -0.148522 | -0.177595 | 60.0 | 0.0 | |
| … | | … | … | … | … | … | … | … | |
| TCGA-HT-7855 | | 0.982173 | 1.707482 | -0.344546 | -0.148522 | -0.177595 | 39.0 | 1.0 | |
| TCGA-S9-A7IY | | 0.982173 | -0.585658 | 2.902366 | -0.148522 | -0.177595 | 39.0 | 1.0 | |
| TCGA-S9-A6TY | | -1.018150 | -0.585658 | -0.344546 | 6.733003 | -0.177595 | 50.0 | 1.0 | |
| TCGA-TM-A84G | | 0.982173 | -0.585658 | -0.344546 | -0.148522 | -0.177595 | 54.0 | 0.0 | |
| TCGA-FG-A87Q | | -1.018150 | -0.585658 | -0.344546 | -0.148522 | -0.177595 | 61.0 | 0.0 | |

729 rows × 11 columns

In [131… 
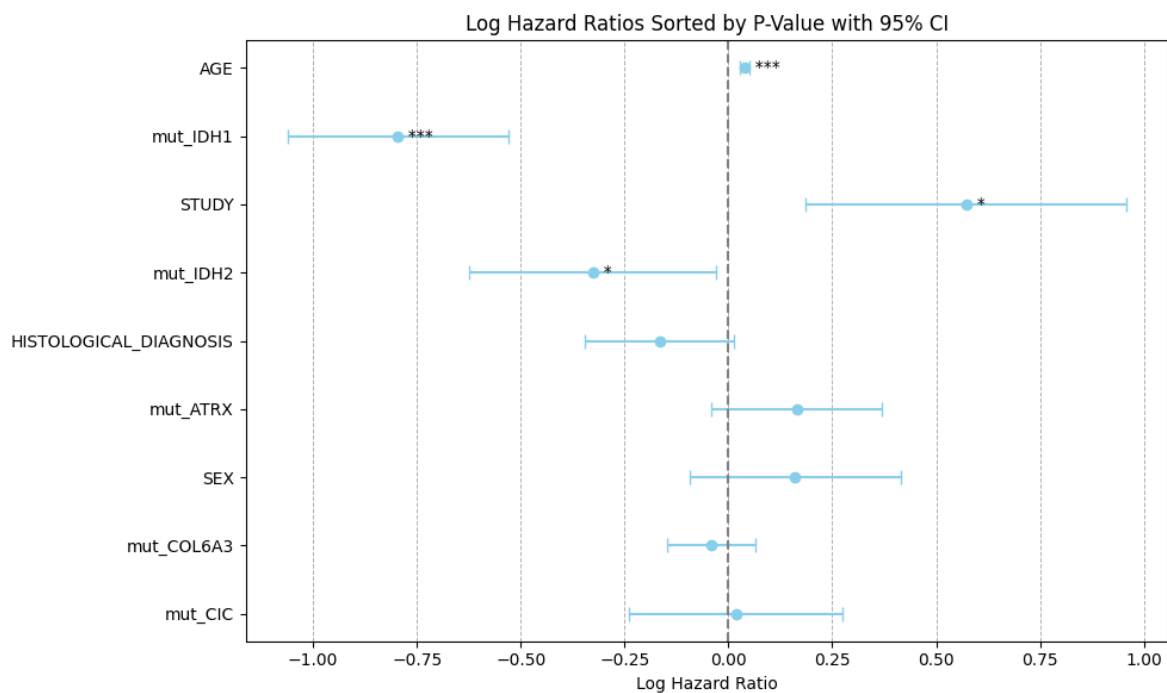```
# build a cox model
coxm = flexynesis.build_cox_model(df, 'OS_MONTHS', 'OS_STATUS')
```

No low variance features were removed based on event conditioning.

In [132… 
```
# visualize log-hazard ratios sorted by p-values
flexynesis.plot_hazard_ratios(coxm)
```

/home/thesamurai/micromamba/envs/flexynesisenv/lib/python3.11/site-package
s/flexynesis/utils.py:765: FutureWarning: Series.__getitem__ treating keys
as positions is deprecated. In a future version, integer keys will always
be treated as labels (consistent with DataFrame behavior). To access a val
ue by position, use `ser.iloc[pos]`

Log Hazard Ratios Sorted by P-Value with 95% CI

## 3.3 Final Exercise

- Inspect the top 10 markers from section 3.2 and see if they have been characterized in the literature as important markers for Glioma disease progression.

The following is a summary of the genes listed in the data table concerning their characterization as markers for glioma disease progression. Each gene's importance is explicitly supported by relevant references.

IDH1: Mutations in IDH1 are frequently examined as significant markers in gliomas, particularly in lower-grade gliomas, where they are associated with improved survival outcomes. Studies show that IDH1 mutations are linked to distinctive metabolic changes within tumor cells, affecting tumor behavior and patient prognosis Zhang et al. (2020)Wang et al., 2020; Ji et al., 2019). These mutations are crucial in distinguishing between different glioma subtypes and can guide therapeutic strategies (Koso et al., 2012).

ATRX: ATRX mutations are often found in conjunction with IDH1 mutations and are significant in glioma progression and patient outcomes. Loss of ATRX function is correlated with increased genomic instability and aggressive tumor characteristics. Research demonstrates that ATRX mutations can indicate better prognosis due to their association with favorable tumor characteristics, suggesting their dual role as both markers of progression and potential therapeutic targets (Ji et al., 2019; Wang et al., 2020).

CIC: Although not as extensively characterized as IDH1 or ATRX, CIC mutations have been implicated in gliomagenesis. They often occur in the context of IDH1 mutations and may contribute to aggressive tumor behavior. Evidence suggests that mutations in CIC could influence the tumoral environment and patient prognosis, indicating a need for further clinical studies (Ji et al., 2019; Choi et al., 2024).

IDH2: Similar to IDH1, mutations in IDH2 are emerging as significant markers in gliomas. While their direct prognostic value remains less established compared to IDH1, their role in the metabolic and epigenetic landscape of gliomas indicates a potential link to tumor progression and patient outcomes. Investigations are ongoing to fully elucidate their implications in glioma pathology (Li et al., 2013).

TP53: TP53 mutations are frequently encountered in gliomas and are associated with both tumor progression and adverse prognosis. This tumor suppressor gene plays a pivotal role in regulating the cell cycle and apoptosis, so its mutations can lead to unchecked progression of gliomas. The frequency of TP53 mutations in high-grade gliomas underscores its relevance as a prognostic marker and therapeutic target (Li et al., 2013; Luo et al., 2023).

PIK3CA: PIK3CA mutations are implicated in glioma tumorigenesis and progression. They contribute to alterations in cellular signaling pathways that promote tumor growth and survival. Studies indicate that PIK3CA is linked with aggressive glioma phenotypes, making it an essential marker in understanding glioma pathology (Chen et al., 2013; Choi et al., 2024).

COL6A3: COL6A3, while less frequently cited as a direct marker of glioma progression, has been recognized for its role in modifying tumor microenvironments. Its expression has been correlated with tumor invasiveness and may serve as a potential marker of disease progression warranting further exploration (Chen et al., 2013; Choi et al., 2024).

RELN: Research into RELN has revealed its involvement in signaling pathways that affect glioma cell proliferation and migration. As a neuromodulator, disruptions in RELN signaling may correlate with glioma aggressiveness, suggesting its utility in prognostic evaluations (Cai et al., 2020; Choi et al., 2024).

BRD3: Recent studies have linked BRD3 with the regulation of cell proliferation and survival pathways in gliomas. Its expression levels might provide insights into patient prognosis and therapeutic responses, highlighting its relevance as a biomarker (Chen et al., 2013; Li & Lan, 2021; Kang et al., 2021).

These references highlight the critical role that genetic mutations play in glioma progression and the potential for utilizing these markers in clinical diagnostics and therapeutic strategies.

References:

Zhang et al. (2020): Zhang et al. "Identification of hub genes related to prognosis in glioma" Bioscience reports (2020) doi:10.1042/bsr20193377

Wang et al. (2020): Wang et al. "Systematically Dissecting the Function of RNA-Binding Proteins During Glioma Progression" Frontiers in genetics (2020) doi:10.3389/fgene.2019.01394

Cai et al. (2020): Cai et al. "LncRNA LINC00998 inhibits the malignant glioma phenotype via the CBX3-mediated c-Met/Akt/mTOR axis" Cell death and disease (2020)

doi:10.1038/s41419-020-03247-6

Koso et al. (2012): Koso et al. "Transposon mutagenesis identifies genes that transform neural stem cells into glioma-initiating cells" Proceedings of the national academy of sciences (2012) doi:10.1073/pnas.1215899109

Ji et al. (2019): Ji et al. "A panel of synapse assembly genes as a biomarker for Gliomas" (2019) doi:10.1101/19011114

Li et al. (2013): Li et al. "Comprehensive analysis of the functional microRNA–mRNA regulatory network identifies miRNA signatures associated with glioma malignant progression" Nucleic acids research (2013) doi:10.1093/nar/gkt1054

Luo et al. (2023): Luo et al. "Hypermethylation of HIC2 is a potential prognostic biomarker and tumor suppressor of glioma based on bioinformatics analysis and experiments" CNS neuroscience & therapeutics (2023) doi:10.1111/cns.14093

Chen et al. (2013): Chen et al. "Knockdown of FRAT1 Expression by RNA Interference Inhibits Human Glioblastoma Cell Growth, Migration and Invasion" PLOS One (2013) doi:10.1371/journal.pone.0061206

Choi et al. (2024): Choi et al. "Significant Genes Associated with Mortality and Disease Progression in Grade II and III Glioma" Biomedicines (2024) doi:10.3390/biomedicines12040858

Li & Lan (2021): Li and Lan "Bioinformatics analysis reveals a stem cell-expressed circ-Serpine2-mediated miRNA-mRNA regulatory subnetwork in the malignant progression of glioma" Journal of translational medicine (2021) doi:10.1186/s12967-021-03118-4

Kang et al. (2021): Kang et al. "Genomic instability in lower-grade glioma: Prediction of prognosis based on lncRNA and immune infiltration" Molecular therapy — oncolytics (2021) doi:10.1016/j.omto.2021.07.011