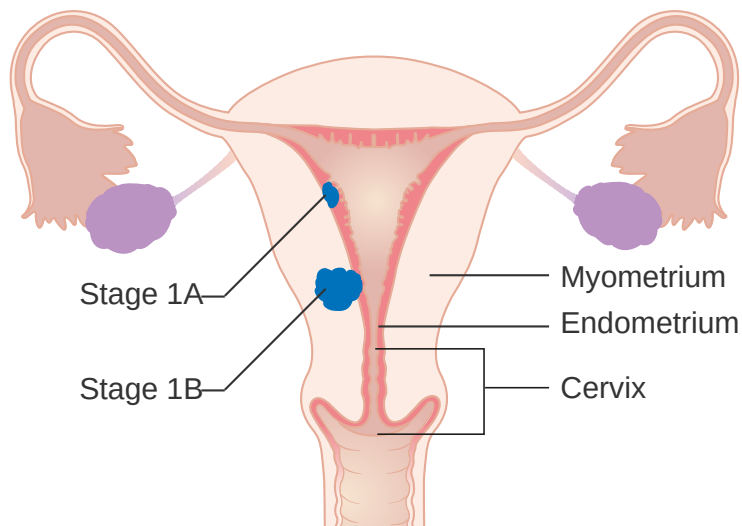


# Classification of Endometrial cancer using Machine Learning

Robert Róžański

# Endometrial cancer



Type I (less aggressive):

- Endometrioid adenocarcinoma

Type II (more aggressive):

- Serous cystadenocarcinoma
- Papillary serous adenocarcinoma
- Clear cell adenocarcinoma

# Data

## QIAseq Targeted DNA Panels

 Print



Digital DNA sequencing to confidently detect low-frequency variants

- **Digital sequencing enabled by molecular barcodes to remove PCR duplicates**
- **Complete Sample to Insight solution streamlines the workflow**
- **Compatibility with low-quality DNA enables efficient sequencing of FFPE and cfDNA samples**
- **Minimal DNA input to preserve precious samples**
- **Optimized buffers and conditions to achieve high coverage of GC-rich regions**

The QIAseq targeted DNA Panels have been developed as a complete Sample to Insight solution to enable digital DNA sequencing by utilizing molecular barcodes. Digital DNA sequencing is a unique approach to detect low-frequency variants with high confidence by overcoming the issues of PCR duplicates, false positives and library bias.

# Data

## Adenoma

Tubular Adenoma: BRAF, FBXW7, KRAS.

Tubulovillous Adenoma: BRAF, CTNNB1, KRAS, NRAS, PIK3CA (p110-alpha).

Villous Adenoma: BRAF, KRAS.

Other Adenoma-Related Genes: APC, DMD, SMAD4 (MADH4), STK11 (LKB1), TCF7L2.

## Carcinoid-Endocrine Tumor

APC, CTNNB1, TP53 (p53).

## Carcinoma

Adenocarcinoma: ACVR1B, AKT1, APC, ATM, ATP6V0D2, AXIN2, BAX, BLM, BMPR1A (ALK3), BRAF, BRCA1, BRCA2, BUB1B, CASP8 (FLICE), CDC27, CDH1 (E-Cadherin), CDK4, CDKN2A (p16INK4a), CHEK2 (RAD53), CTNNA1, CTNNB1, DCC, DMD, EGFR (ERBB1), ENG (EVI-1), EP300, EPCAM, ERBB2 (HER-2, NEU), FBXW7, FGFR3, FLCN, FZD3, GALNT12, GPC6, GREM1, KIT (CD117), KRAS, MAP2K4 (MKK4, JNKK1), MAP7, MET, MIER3, MLH1, MLH3, MSH2, MSH3, MSH6, MUTYH, MYO1B, NRAS, PALB2, PIK3CA (p110-alpha), PIK3R1 (p85-ALPHA), PMS1, PMS2, POLD1, POLE, PTEN, PTPN12, RET, RPS20, SLC9A9, SMAD2 (MADH2), SMAD4 (MADH4), SRC, STK11 (LKB1), TCERG1, TCF7L2, TGFBR2, TP53 (p53), WBSCR17.

Neuroendocrine: BRAF, KRAS.

Serrated Carcinoma: BRAF, PIK3CA (p110-alpha).

Squamous Cell Carcinoma: KRAS, TP53 (p53).

Other Carcinomas: BRCA2, CDKN2A (p16INK4a).

# Training Data

NIH

NATIONAL CANCER INSTITUTE

GDC Data Portal

[Home](#)[Projects](#)[Exploration](#)[Analysis](#)[Repository](#)[Quick Search](#)[Manage Sets](#)[Login](#)[Cart 0](#)[GDC Apps](#)

Cases

Genes

Mutations

[Add a Case Filter](#)

Case

input set

Upload Case Set

Primary Site

☐ Corpus uteri 538

☐ Uterus, NOS 182

Program

☐ TCGA 538

☐ CPTAC 101

☐ FM 81

Project

☐ TCGA-UCEC 538

Clear

Case

IN

input set

View Files in Repository

Cases (720)

Genes (22,160)

Mutations (913,229)

OncoGrid

Primary Site

Project

Disease Type

Gender

Vital Status

Showing 1 - 20 of 720 cases

Biospecimen

Clinical

JSON

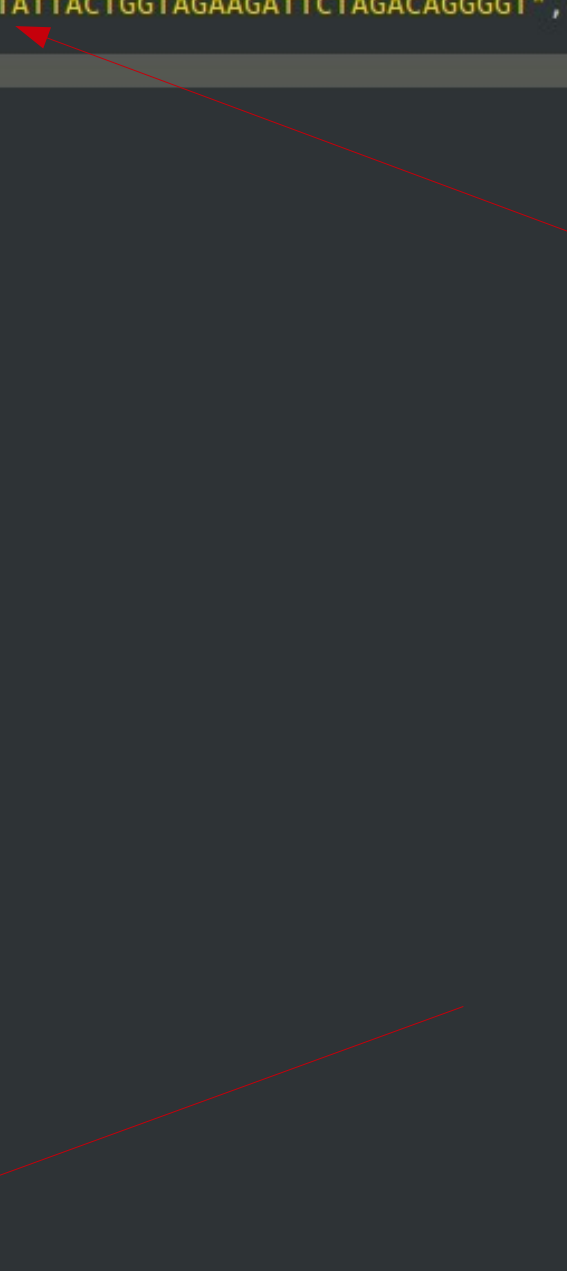
TSV

Save/Edit Case Set

| Case ID   | Project                   | Primary Site | Gender | Files              | Available Files per Data Category |                   |                    |                   |                   |                    |                    | # Mutations            | # Genes                | Slides            |
|---|---------------------------|--------------|--------|--------------------|-----------------------------------|-------------------|--------------------|-------------------|-------------------|--------------------|--------------------|------------------------|------------------------|-------------------|
|   |                           |              |        |                    | Seq                               | Exp               | SNV                | CNV               | Meth              | Clinical           | Bio                |                        |                        |                   |
| <input type="checkbox"/> <a href="#">TCGA-A5-A0G2</a> | <a href="#">TCGA-UCEC</a> | Corpus uteri | Female | <a href="#">58</a> | <a href="#">4</a>                 | <a href="#">5</a> | <a href="#">16</a> | <a href="#">5</a> | <a href="#">1</a> | <a href="#">10</a> | <a href="#">17</a> | <a href="#">42,051</a> | <a href="#">14,357</a> | <a href="#">3</a> |
| <input type="checkbox"/> <a href="#">TCGA-EO-A22U</a> | <a href="#">TCGA-UCEC</a> | Corpus uteri | Female | <a href="#">57</a> | <a href="#">4</a>                 | <a href="#">5</a> | <a href="#">16</a> | <a href="#">5</a> | <a href="#">1</a> | <a href="#">10</a> | <a href="#">16</a> | <a href="#">26,998</a> | <a href="#">12,629</a> | <a href="#">2</a> |
| <input type="checkbox"/> <a href="#">TCGA-FI-A2D5</a> | <a href="#">TCGA-UCEC</a> | Corpus uteri | Female | <a href="#">58</a> | <a href="#">4</a>                 | <a href="#">5</a> | <a href="#">16</a> | <a href="#">5</a> | <a href="#">1</a> | <a href="#">11</a> | <a href="#">16</a> | <a href="#">26,139</a> | <a href="#">12,482</a> | <a href="#">2</a> |

https://portal.gdc.cancer.gov/exploration?filters=%7B%22content%22%3A%5B%7B%22content%22%3A%7B%22field%22%3A%22cases.case\_id%22%2C%22value%22%3A%5B%22set\_id%3AAWxYARM-1FJwHPndQLUN%22%5D%7D%2C%22op%22%3A%22IN%22%7D%5D%2C%22op%22%3A%22AND%22%7D

```
97 }
98 }
99 ],
100 "genomic_dna_change": "chr12:g.71567009_71567010insTATATAAGTATTACTGGTAGAAGATTCTAGACAGGGGT",
101 "ssm_id": "1f850024-4e34-512b-880d-734e133bd006",
102 "mutation_subtype": "Small insertion"
103 },{
104 "consequence": [
105 {
106 "transcript": {
107 "is_canonical": false,
108 "gene": {
109 "symbol": "PDE12",
110 "gene_id": "ENSG00000174840"
111 },
112 "aa_change": null,
113 "annotation": {
114 "vep_impact": "MODIFIER",
115 "polyphen_impact": "",
116 "sift_impact": ""
117 },
118 "consequence_type": "3_prime_UTR_variant"
119 }
120 },
121 {
122 "transcript": {
123 "is_canonical": true,
124 "gene": {
125 "symbol": "PDE12",
126 "gene_id": "ENSG00000174840"
127 },
128 "aa_change": null,
129 "annotation": {
130 "vep_impact": "MODIFIER",
131 "polyphen_impact": "",
132 "sift_impact": ""
133 },
134 "consequence_type": "3_prime_UTR_variant"
135 }
136 }
137 ],
138 "genomic_dna_change": "chr3:g.57560088A>C",
139 "ssm_id": "c088a181-5490-5028-8d0f-01912288c38b",
140 "mutation_subtype": "Single base substitution"
141 },{
142 "consequence": [
143 {
```

Two red arrows are present in the image. The first arrow originates from the right side of the image and points to the string "chr12:g.71567009\_71567010insTATATAAGTATTACTGGTAGAAGATTCTAGACAGGGGT" on line 100. The second arrow also originates from the right side and points to the string "chr3:g.57560088A>C" on line 138.



```
97 }
98 }
99 ],
100 "genomic_dna_change": "chr12:g.71567009_71567010insTATATAAGTATTACTGGTAGAAGATTCTAGACAGGGGT",
101 "ssm_id": "1f850024-4e34-512b-880d-734e133bd006",
102 "mutation_subtype": "Small insertion"
103 },{
104 "consequence": [
105 {
106 "transcript": {
107 "is_canonical": false,
108 "gene": {
109 "symbol": "PDE12",
110 "gene_id": "ENSG00000174840"
111 },
112 "aa_change": null,
113 "annotation": {
114 "vep_impact": "MODIFIER",
115 "polyphen_impact": "",
116 "sift_impact": ""
117 },
118 "consequence_type": "3_prime_UTR_variant"
119 }
120 },
121 {
122 "transcript": {
123 "is_canonical": true,
124 "gene": {
125 "symbol": "PDE12",
126 "gene_id": "ENSG00000174840"
127 },
128 "aa_change": null,
129 "annotation": {
130 "vep_impact": "MODIFIER",
131 "polyphen_impact": "",
132 "sift_impact": ""
133 },
134 "consequence_type": "3_prime_UTR_variant"
135 }
136 }
137 ],
138 "genomic_dna_change": "chr3:g.57560088A>C",
139 "ssm_id": "c088a181-5490-5028-8d0f-01912288c38b",
140 "mutation_subtype": "Single base substitution"
141 },{
142 "consequence": [
143 {
```

```
97 }
98 }
99 ],
100 "genomic_dna_change": "chr12:g.71567009_71567010insTATATAAGTATTACTGGTAGAAGATTCTAGACAGGGGT",
101 "ssm_id": "1f850024-4e34-512b-880d-734e133bd006",
102 "mutation_subtype": "Small insertion"
103 },{
104 "consequence": [
105 {
106 "transcript": {
107 "is_canonical": false,
108 "gene": {
109 "symbol": "PDE12",
110 "gene_id": "ENSG00000174840"
111 },
112 "aa_change": null,
113 "annotation": {
114 "vep_impact": "MODIFIER",
115 "polyphen_impact": "",
116 "sift_impact": ""
117 },
118 "consequence_type": "3_prime_UTR_variant"
119 }
120 },
121 {
122 "transcript": {
123 "is_canonical": true,
124 "gene": {
125 "symbol": "PDE12",
126 "gene_id": "ENSG00000174840"
127 },
128 "aa_change": null,
129 "annotation": {
130 "vep_impact": "MODIFIER",
131 "polyphen_impact": "",
132 "sift_impact": ""
133 },
134 "consequence_type": "3_prime_UTR_variant"
135 }
136 }
137 ],
138 "genomic_dna_change": "chr3:g.57560088A>C",
139 "ssm_id": "c088a181-5490-5028-8d0f-01912288c38b",
140 "mutation_subtype": "Single base substitution"
141 },{
142 "consequence": [
143 {
```



```
97 }
98 }
99 ],
100 "genomic_dna_change": "chr12:g.71567009_71567010insTATATAAGTATTACTGGTAGAAGATTCTAGACAGGGGT",
101 "ssm_id": "1f850024-4e34-512b-880d-734e133bd006",
102 "mutation_subtype": "Small insertion"
103 },{
104 "consequence": [
105 {
106 "transcript": {
107 "is_canonical": false,
108 "gene": {
109 "symbol": "PDE12",
110 "gene_id": "ENSG00000174840"
111 },
112 "aa_change": null,
113 "annotation": {
114 "vep_impact": "MODIFIER",
115 "polyphen_impact": "",
116 "sift_impact": ""
117 },
118 "consequence_type": "3_prime_UTR_variant"
119 }
120 },
121 {
122 "transcript": {
123 "is_canonical": true,
124 "gene": {
125 "symbol": "PDE12",
126 "gene_id": "ENSG00000174840"
127 },
128 "aa_change": null,
129 "annotation": {
130 "vep_impact": "MODIFIER",
131 "polyphen_impact": "",
132 "sift_impact": ""
133 },
134 "consequence_type": "3_prime_UTR_variant"
135 }
136 }
137 ],
138 "genomic_dna_change": "chr3:g.57560088A>C",
139 "ssm_id": "c088a181-5490-5028-8d0f-01912288c38b",
140 "mutation_subtype": "Single base substitution"
141 },{
142 "consequence": [
143 {
```

# Feature Engineering

- gene
- VEP
- GO
- Signor

# Feature engineering: genes and VEPs

## genes

| case_id    | type | diagnosis                   | ACVR1B | AKT1 | ALK3 | APC | ATM | ATP6V0D2 |
|------------|------|-----------------------------|--------|------|------|-----|-----|----------|
| 8a0a58a0-5 | 2    | Serous cystadenocarcinoma   | 0      | 0    | 0    | 0   | 0   | 0        |
| 4db38349-2 | 1    | Endometrioid adenocarcinoma | 1      | 0    | 0    | 0   | 1   | 0        |
| 119b1761-a | 1    | Endometrioid adenocarcinoma | 0      | 0    | 1    | 0   | 1   | 1        |
| 435b57e7-5 | 1    | Endometrioid adenocarcinoma | 1      | 0    | 1    | 0   | 1   | 1        |

## VEP 1

| case_id    | type | diagnosis                   | ('ACVR1B', 'HIGH') | ('ACVR1B', 'LOW') | ('ACVR1B', 'MODERATE') | ('ACVR1B', 'MODIFIER') | ('AKT1', 'HIGH') | ('AKT1', 'LOW') |
|------------|------|-----------------------------|--------------------|-------------------|------------------------|------------------------|------------------|-----------------|
| 8a0a58a0-5 | 2    | Serous cystadenocarcinoma   | 0                  | 0                 | 0                      | 0                      | 0                | 0               |
| 4db38349-2 | 1    | Endometrioid adenocarcinoma | 0                  | 0                 | 1                      | 1                      | 0                | 0               |
| 119b1761-a | 1    | Endometrioid adenocarcinoma | 0                  | 0                 | 0                      | 0                      | 0                | 0               |
| 435b57e7-5 | 1    | Endometrioid adenocarcinoma | 0                  | 0                 | 0                      | 1                      | 0                | 0               |

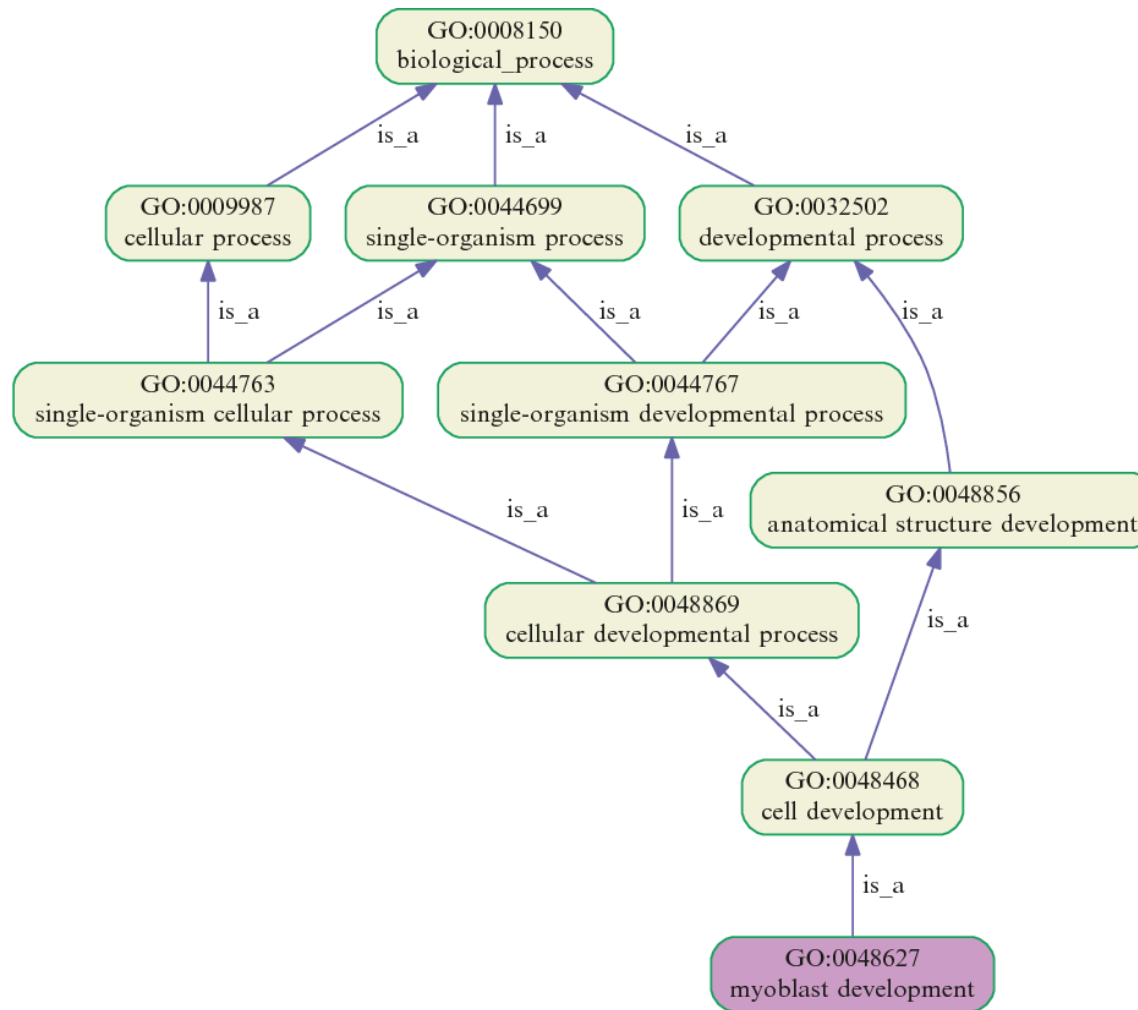
## VEP 2

| case_id    | type | diagnosis                   | ('ACVR1B', 'MODIFIER') | ('ACVR1B', 'VEP_SCORE') | ('AKT1', 'MODIFIER') | ('AKT1', 'VEP_SCORE') | ('ALK3', 'MODIFIER') | ('ALK3', 'VEP_SCORE') |
|------------|------|-----------------------------|------------------------|-------------------------|----------------------|-----------------------|----------------------|-----------------------|
| 8a0a58a0-5 | 2    | Serous cystadenocarcinoma   | 0                      | 0                       | 0                    | 0                     | 0                    | 0                     |
| 4db38349-2 | 1    | Endometrioid adenocarcinoma | 1                      | 2                       | 0                    | 0                     | 0                    | 0                     |
| 119b1761-a | 1    | Endometrioid adenocarcinoma | 0                      | 0                       | 1                    | 2                     | 0                    | 0                     |
| 435b57e7-5 | 1    | Endometrioid adenocarcinoma | 1                      | 0                       | 1                    | 2                     | 0                    | 0                     |

## VEP 3

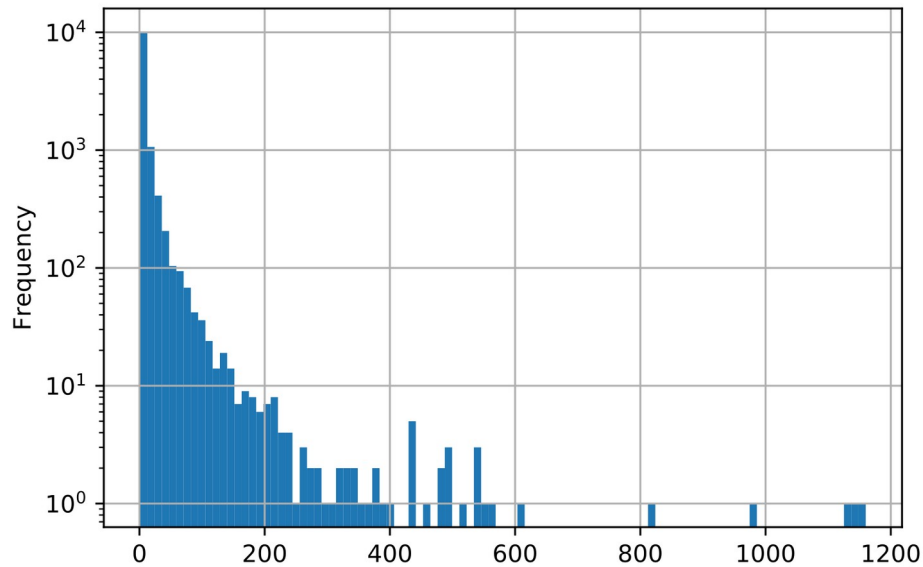
| case_id    | type | diagnosis                   | ACVR1B | AKT1 | ALK3 | APC | ATM | ATP6V0D2 |
|------------|------|-----------------------------|--------|------|------|-----|-----|----------|
| 8a0a58a0-5 | 2    | Serous cystadenocarcinoma   | 0      | 0    | 0    | 0   | 0   | 0        |
| 4db38349-2 | 1    | Endometrioid adenocarcinoma | 2      | 0    | 0    | 0   | 3   | 0        |
| 119b1761-a | 1    | Endometrioid adenocarcinoma | 0      | 0    | 2    | 0   | 2   | 3        |
| 435b57e7-5 | 1    | Endometrioid adenocarcinoma | 0      | 0    | 2    | 0   | 2   | 2        |

# Feature engineering: Gene Ontology

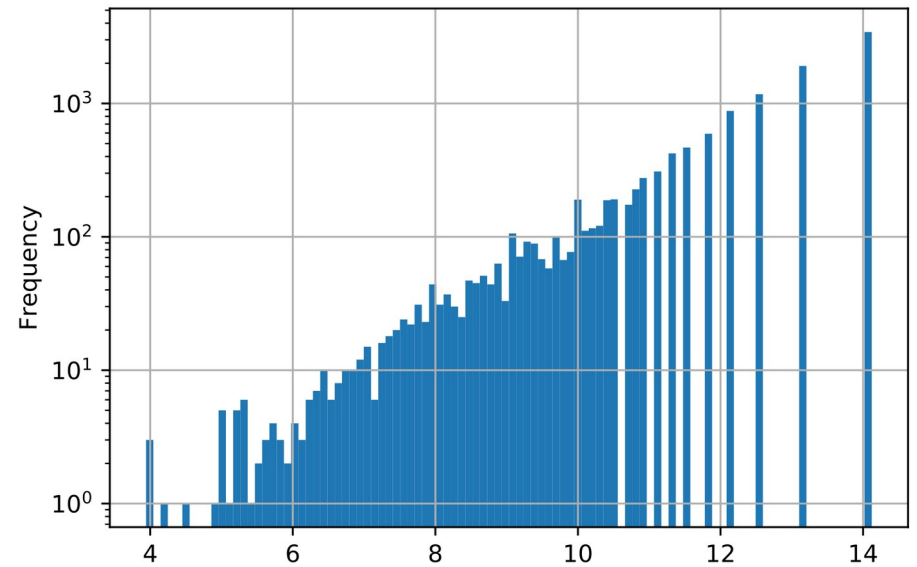


# Feature engineering: Gene Ontology

nb. of genes / GO term



information gain



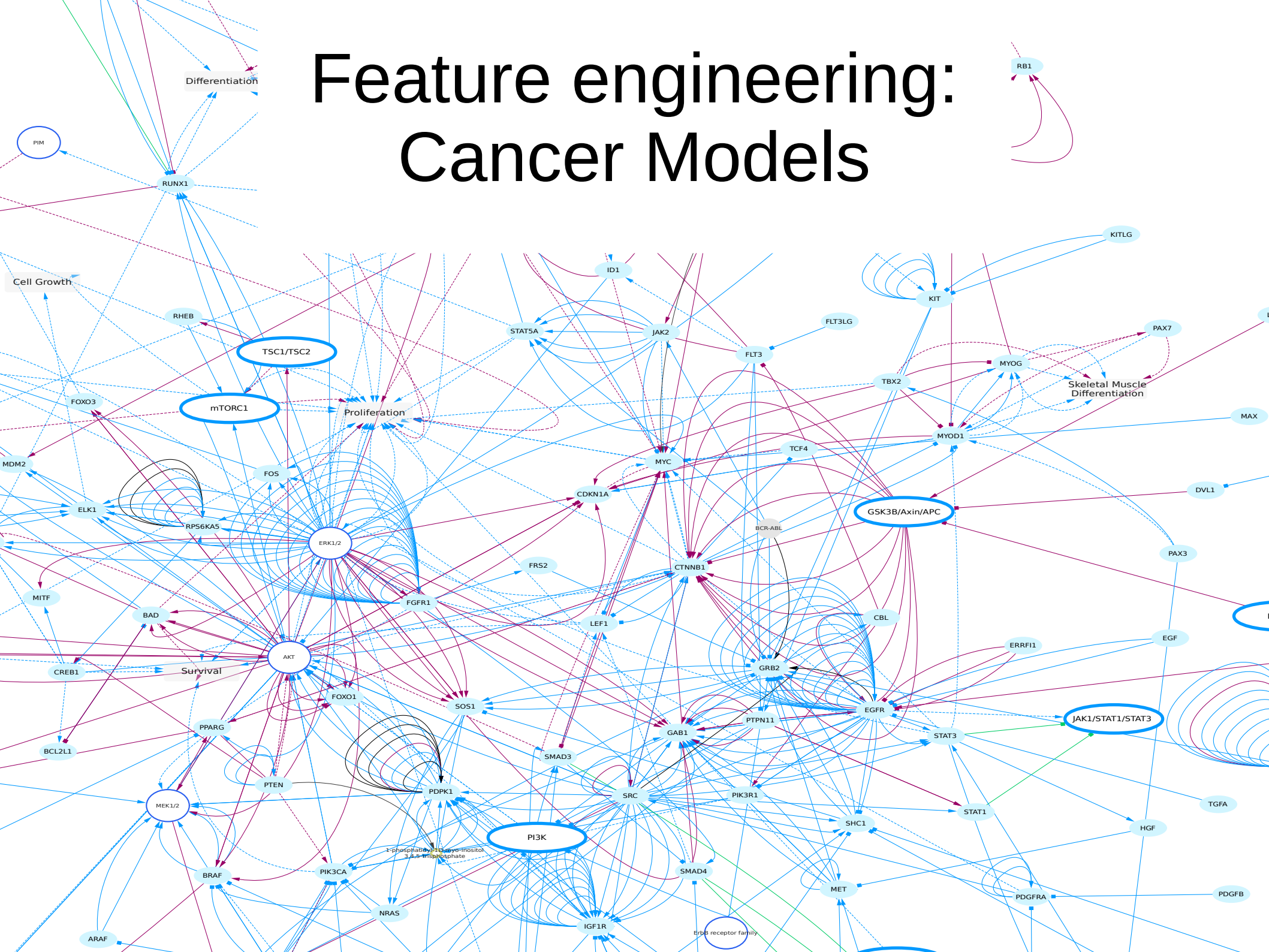
Alterovitz et al., “GO PaD: The Gene Ontology partition database”

# Feature engineering: Gene Ontology

| case_id    | GO:0006355 | GO:0006468 | GO:0006508 | GO:0006511 | GO:0006886 | GO:0006915 | GO:0006954 |
|------------|------------|------------|------------|------------|------------|------------|------------|
| 8a0a58a0-5 | 0          | 2          | 4          | 0          | 0          | 0          | 0          |
| 4db38349-2 | 0          | 10         | 16         | 3          | 0          | 0          | 14         |
| 119b1761-a | 0          | 4          | 14         | 2          | 0          | 0          | 11         |
| 435b57e7-5 | 0          | 6          | 17         | 3          | 3          | 0          | 20         |

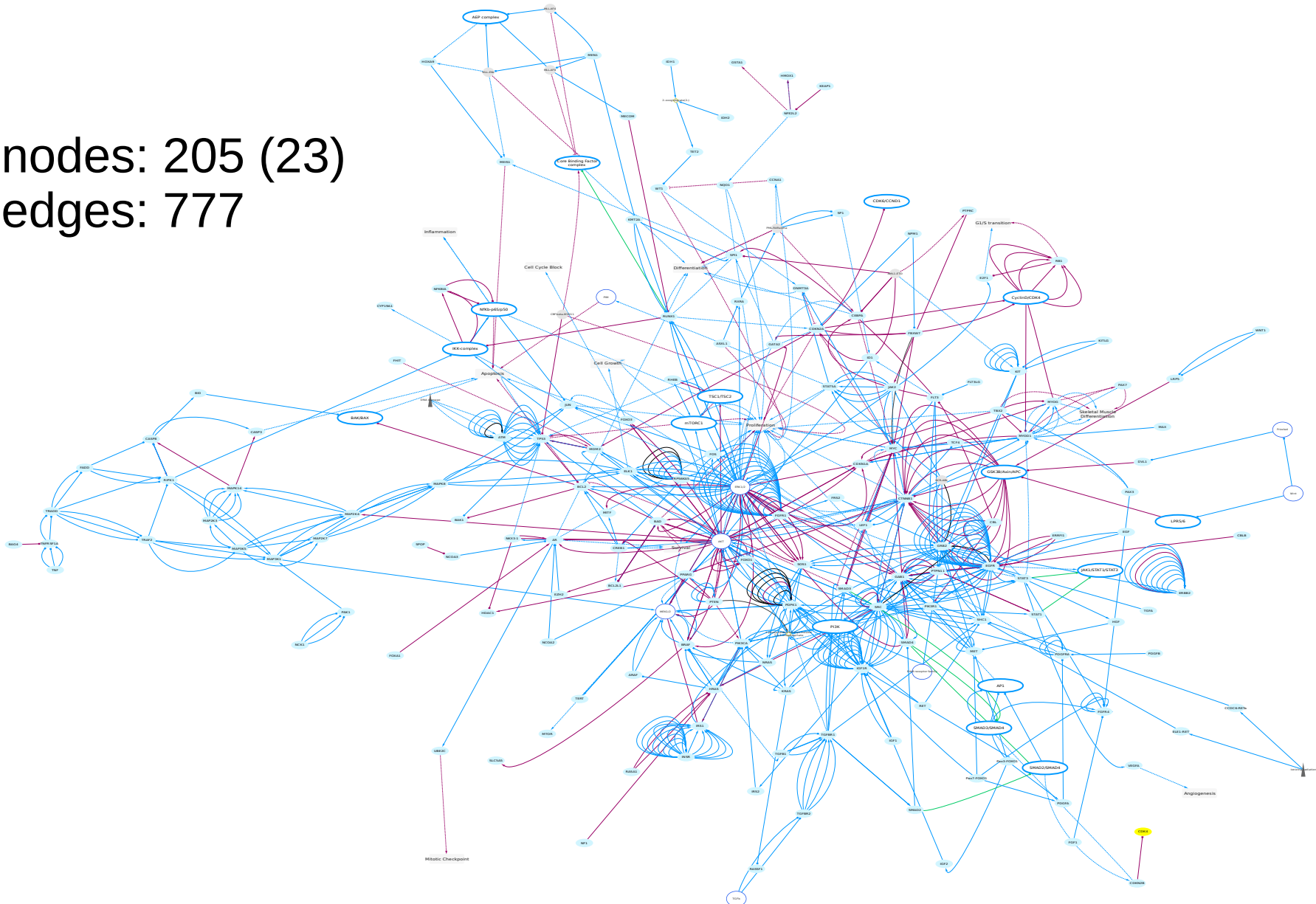


Diagram illustrating the RB1 protein structure, showing a blue oval labeled 'RB1' with a dashed line and two solid lines with arrows pointing towards it.



# Feature engineering: Cancer Models

nodes: 205 (23)  
edges: 777



# Feature engineering: Cancer Models

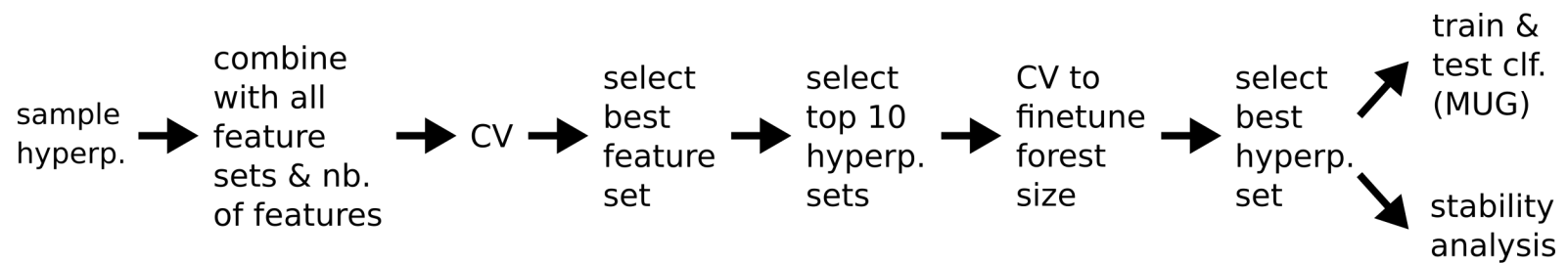
| case_id    | type | diagnosis    | Core Binding Factor comp | 1-phosphatidyl-1D-myo-ino | 2-oxoglutarate(2-) | AEP complex | AKT | AP1 |
|------------|------|--------------|--------------------------|---------------------------|--------------------|-------------|-----|-----|
| 8a0a58a0-5 | 2    | Serous cysta | 0                        | 2                         | 0                  | 0           | 2   | 0   |
| 4db38349-2 | 1    | Endometrioic | 0                        | 5                         | 0                  | 0           | 4   | 0   |
| 119b1761-a | 1    | Endometrioic | 0                        | 7                         | 0                  | 0           | 5   | 0   |
| 435b57e7-5 | 1    | Endometrioic | 0                        | 5                         | 0                  | 0           | 5   | 0   |

# Feature engineering: summary

**Supplementary Table 2.** Feature sets developed in this study.

| Name of the feature set |          | Number of features in each set | Set description   |
|-------------------------|----------|--------------------------------|---|
| <b>Gene</b>             |          | 71                             | Binary: presence or absence of mutations in genes   |
| <b>VEP 1</b>            |          | 281                            | one-hot encoded VEP scores  |
| <b>VEP 2</b>            |          | 142                            | one-hot encoded VEP scores without MODIFIER   |
| <b>VEP 3</b>            |          | 71                             | ordinal VEP encoding (0-3)  |
| <b>VEP 4</b>            |          | 71                             | binary VEP encoding, only MODERATE and HIGH   |
| <b>VEP 5</b>            |          | 71                             | binary VEP encoding, only HIGH  |
| <b>GO (3-13)</b>        | Gene     | 3-281                          | Mutations or VEP scores assigned to Gene Ontology biological process terms. First, the terms were stratified using information gain producing 10 sets of terms (GO 3-13), each set containing terms of similar level of generality. Then, for each GO term set, values from feature sets <i>gene</i> and <i>vep_3-5</i> were mapped to corresponding GO terms and their values summed. For brevity, the feature sets based on Gene Ontology were not described separately. Because the number of GO terms differs between the levels, the number of features also differs (roughly, there are fewer general terms than the more specific ones). |
|                         | VEP 3    | 3-281                          |   |
|                         | VEP 4    | 3-281                          |   |
|                         | VEP 5    | 3-281                          |   |
| <b>SIGNOR</b>           | Gene     | 65                             | Mutations and VEP scores assigned to objects in a tumour signal transduction model. First, a number of tumour-related signaling models from SIGNOR was merged forming a network with 205 nodes and 777 edges representing causal relationships important in oncogenesis. Next, for each node, a set of all predecessor nodes was found. Next, each predecessor was assigned a value copied from feature sets <i>gene</i> and <i>vep_3-5</i> . Finally, all nodes were assigned a value by summing up values of all predecessors. Features that had zero values across all cases were dropped.   |
|                         | VEP 3    | 65                             |   |
|                         | VEP 4    | 65                             |   |
|                         | VEP 5    | 65                             |   |
| <b>SIGNOR enriched</b>  | Gene/VEP | 113                            | As SIGNOR, but genes that were missing from the network were added to the feature sets, with values identical as in the original feature sets.  |

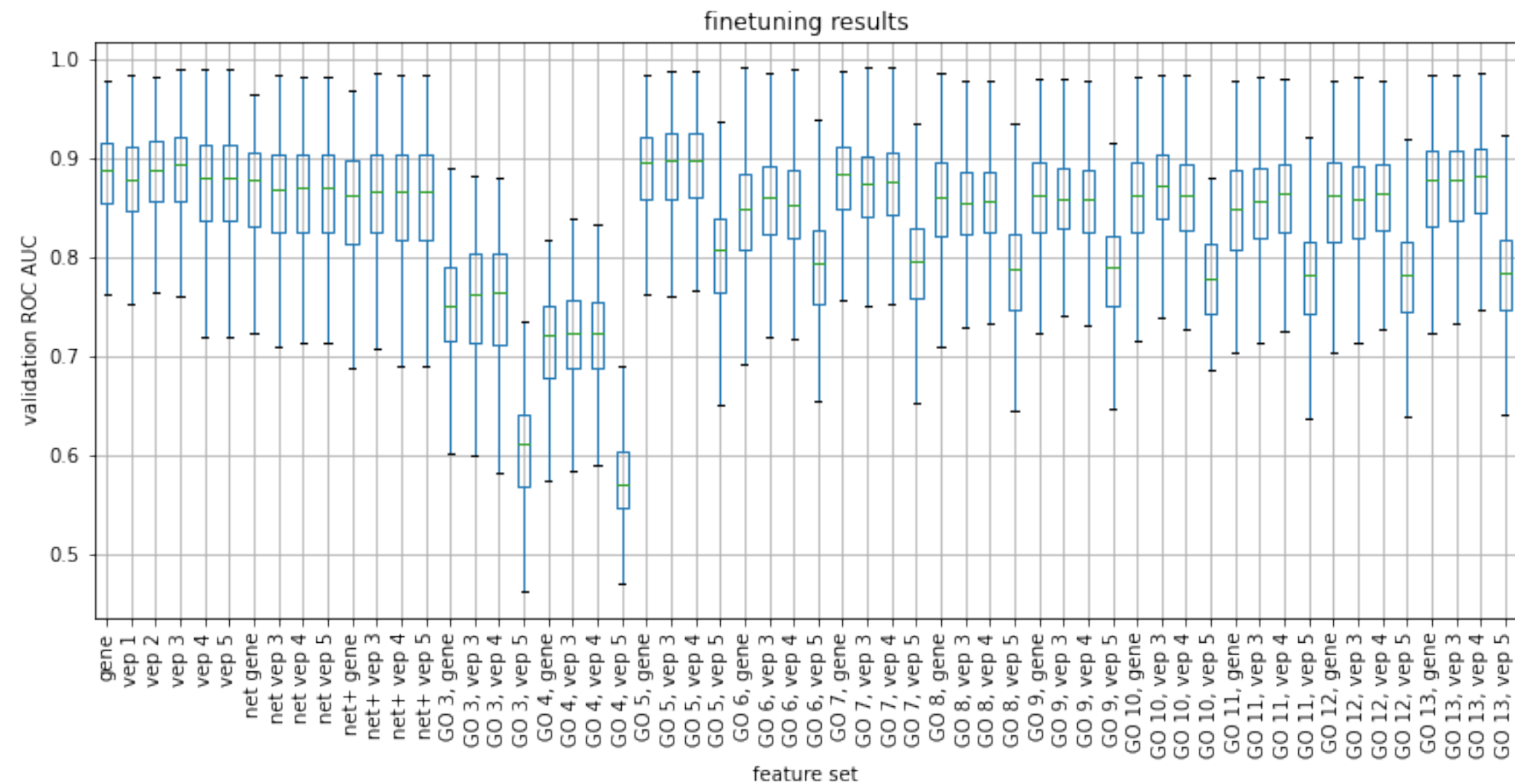
# Finetuning



feature set  
nb. of features  
forest max\_leaf\_nodes  
max\_depth  
min\_samples\_split  
max\_features  
size

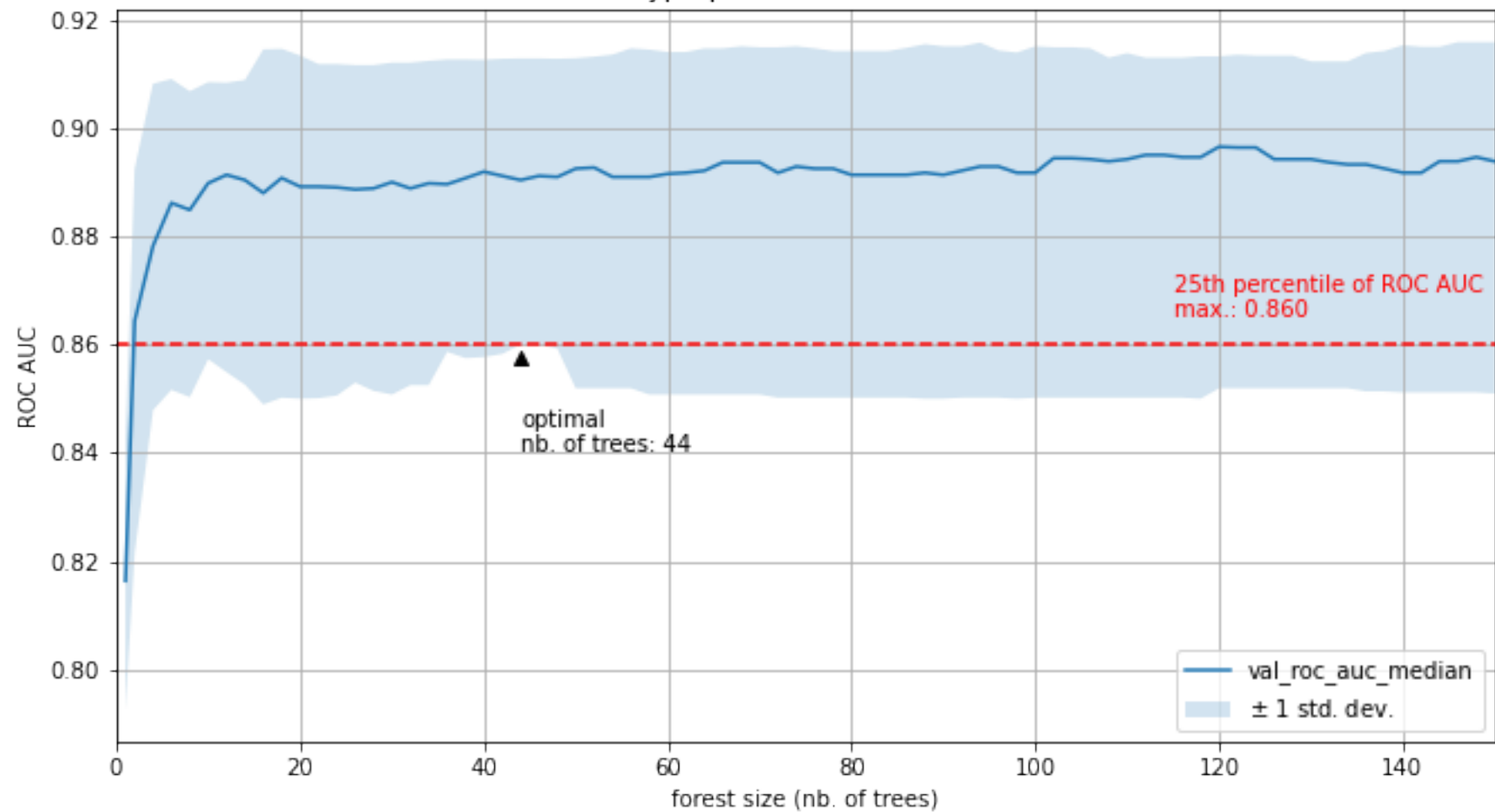


# Best Feature Set

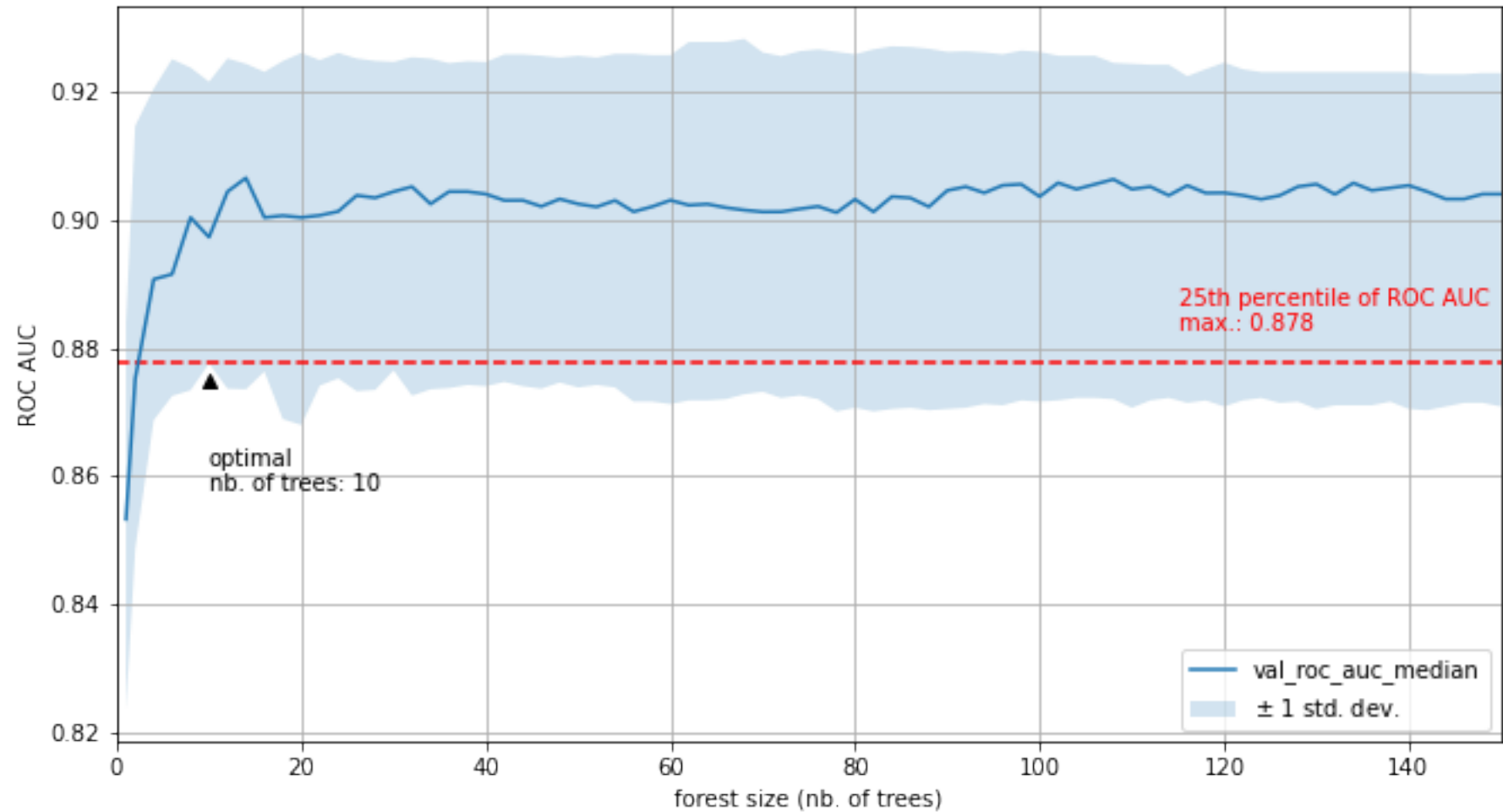




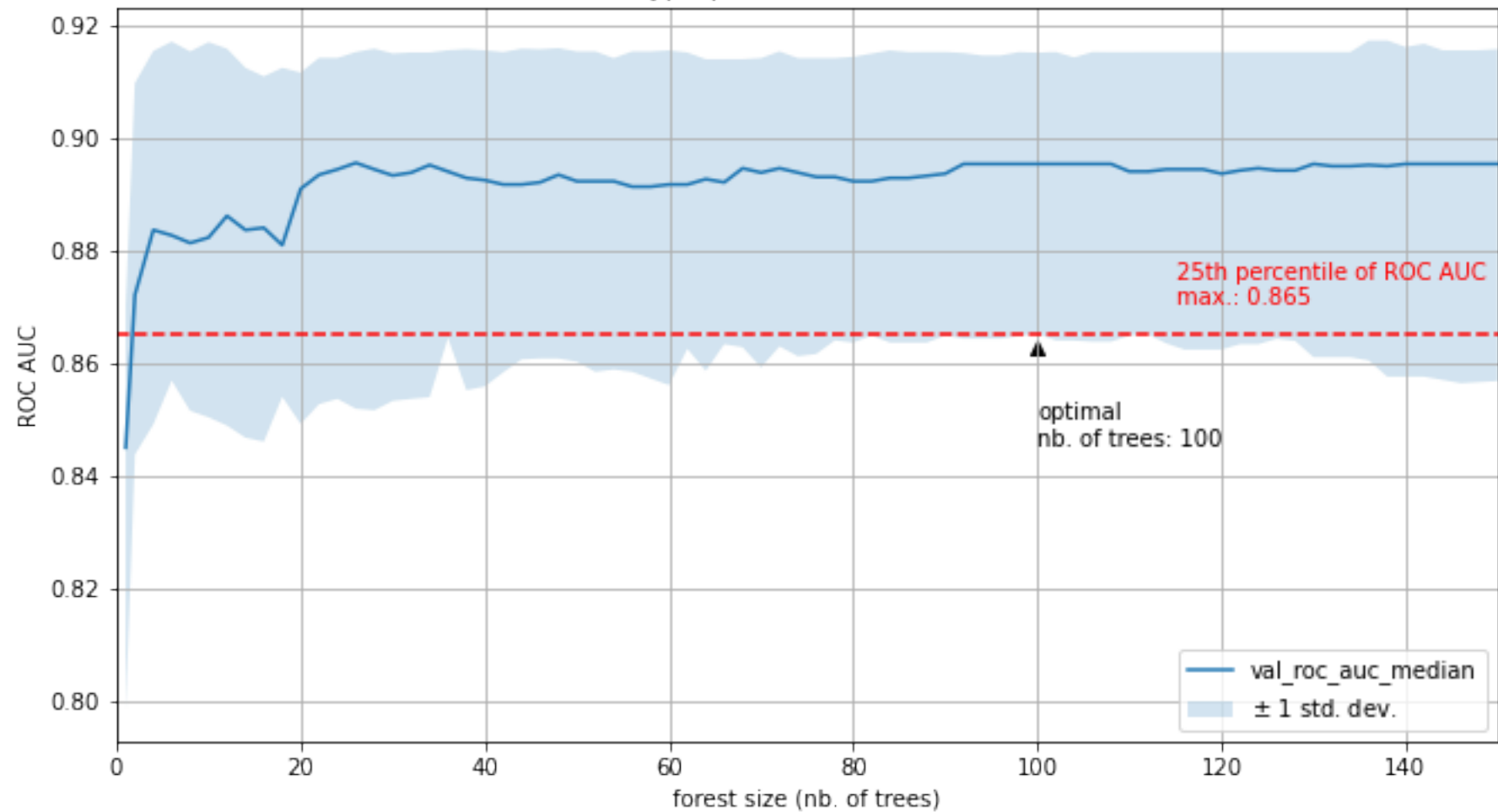
hyperparameter set nb. 8

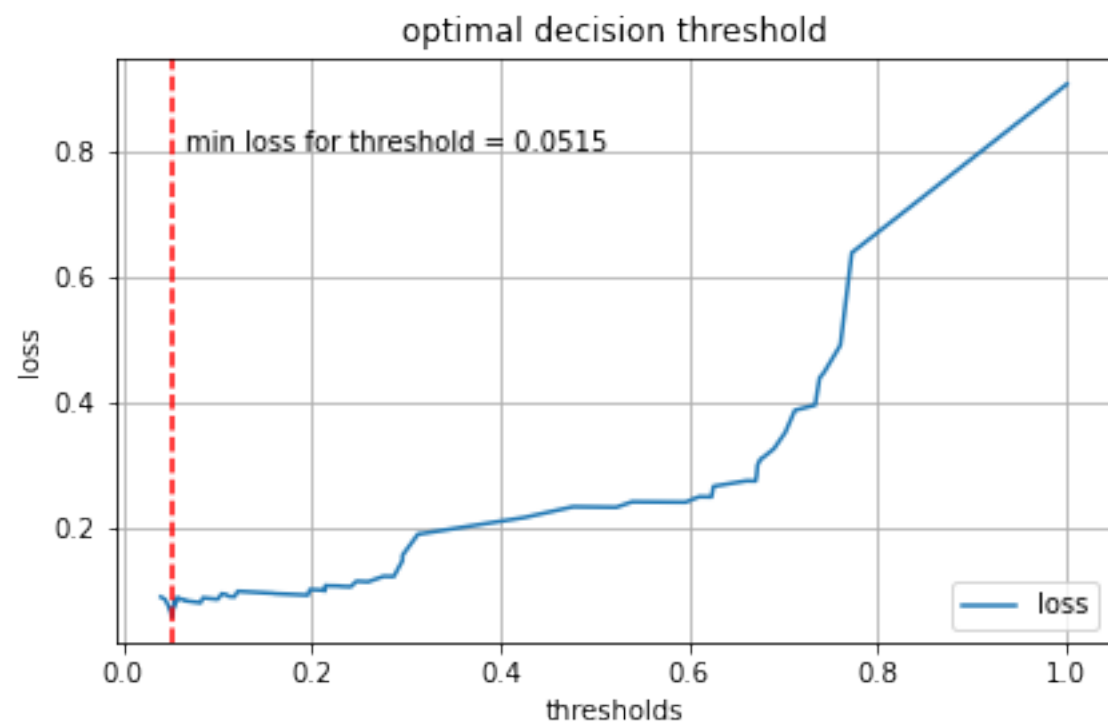


hyperparameter set nb. 0

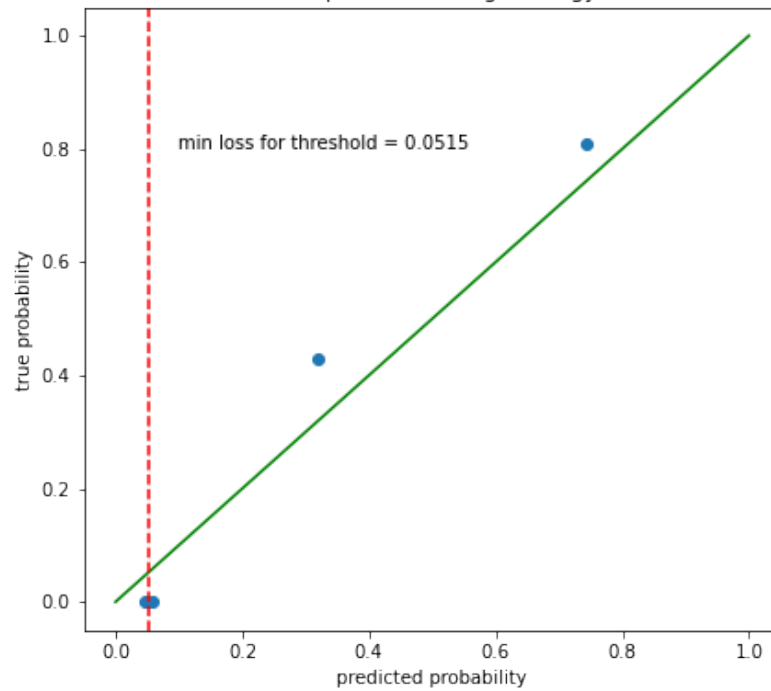


hyperparameter set nb. 2

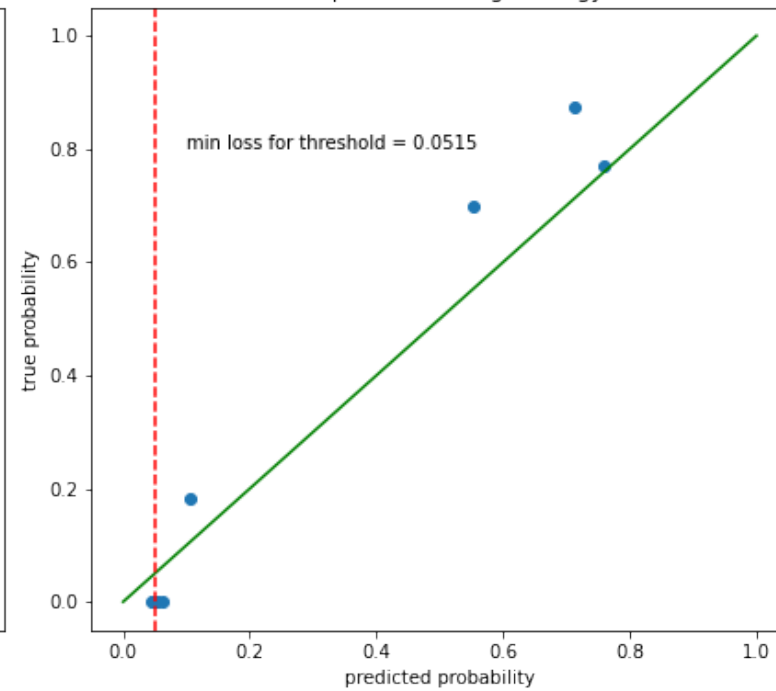




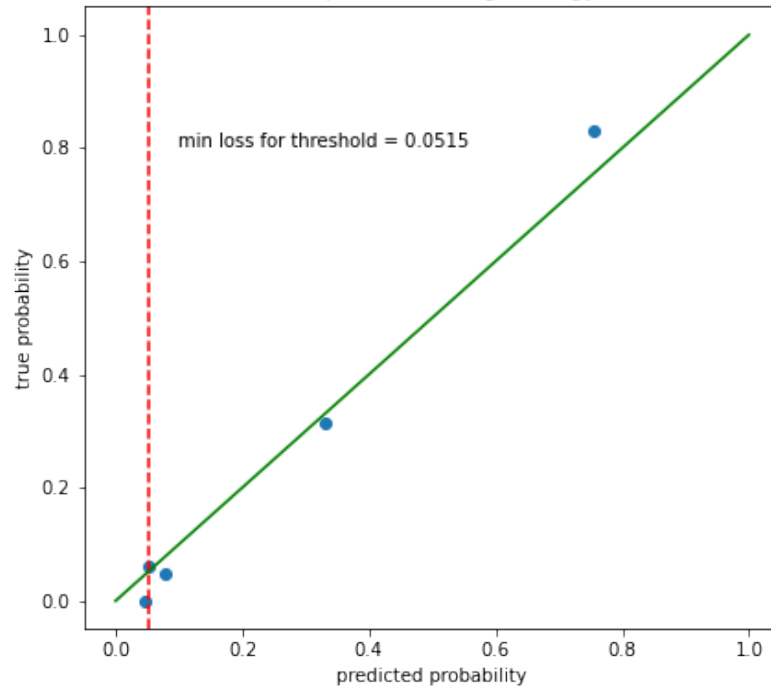
calibration quality on test data  
(n=104) with 5 bins  
and quantile binning strategy



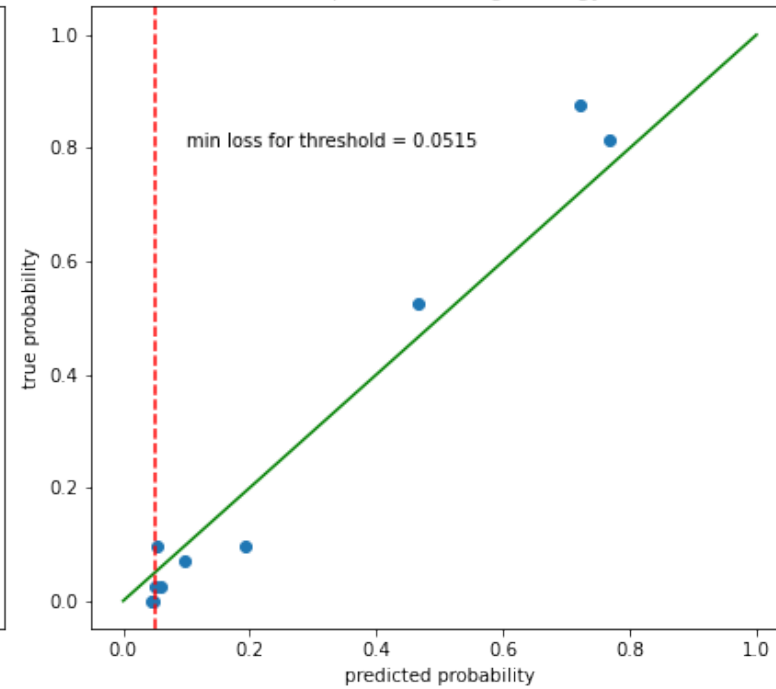
calibration quality on test data  
(n=104) with 10 bins  
and quantile binning strategy

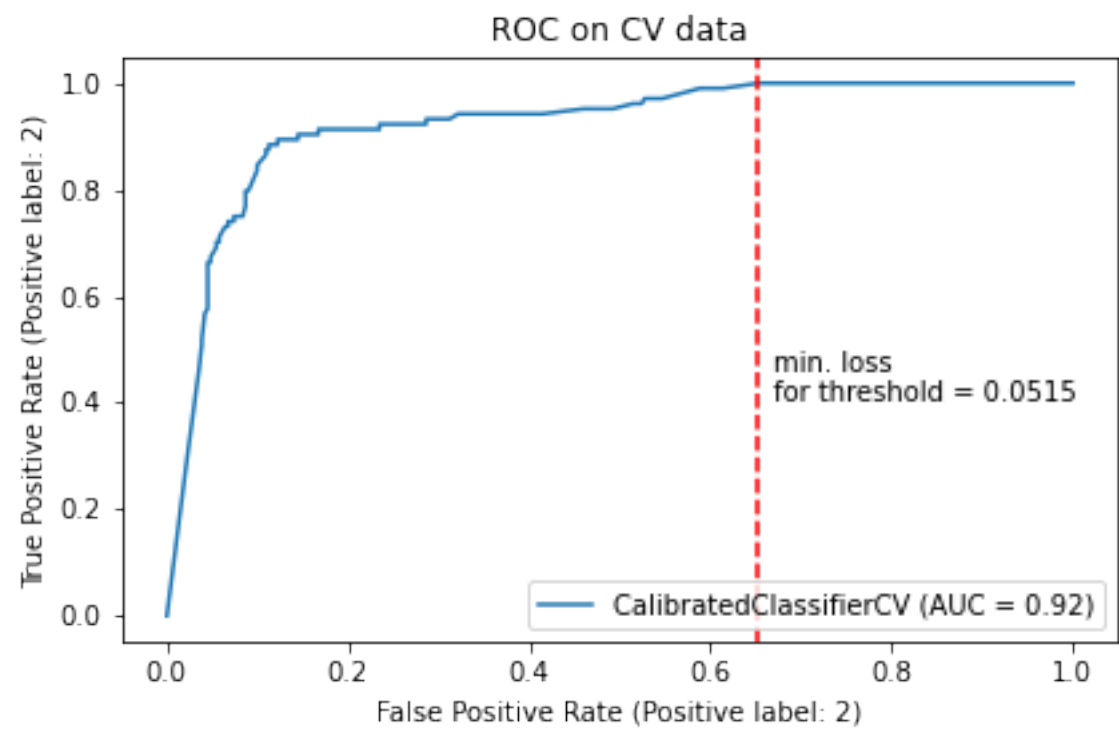


calibration quality on CV data  
(n=415) with 5 bins  
and quantile binning strategy

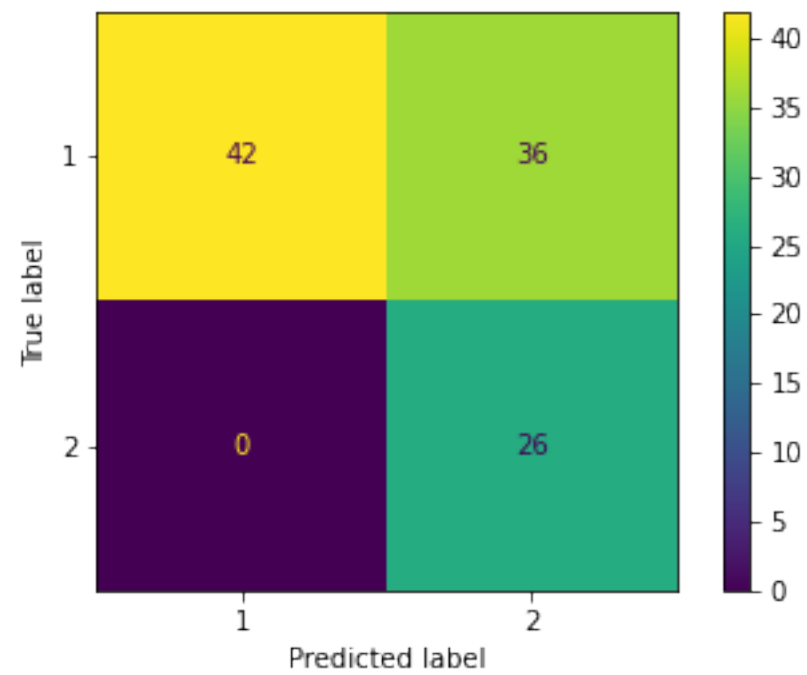


calibration quality on CV data  
(n=415) with 10 bins  
and quantile binning strategy

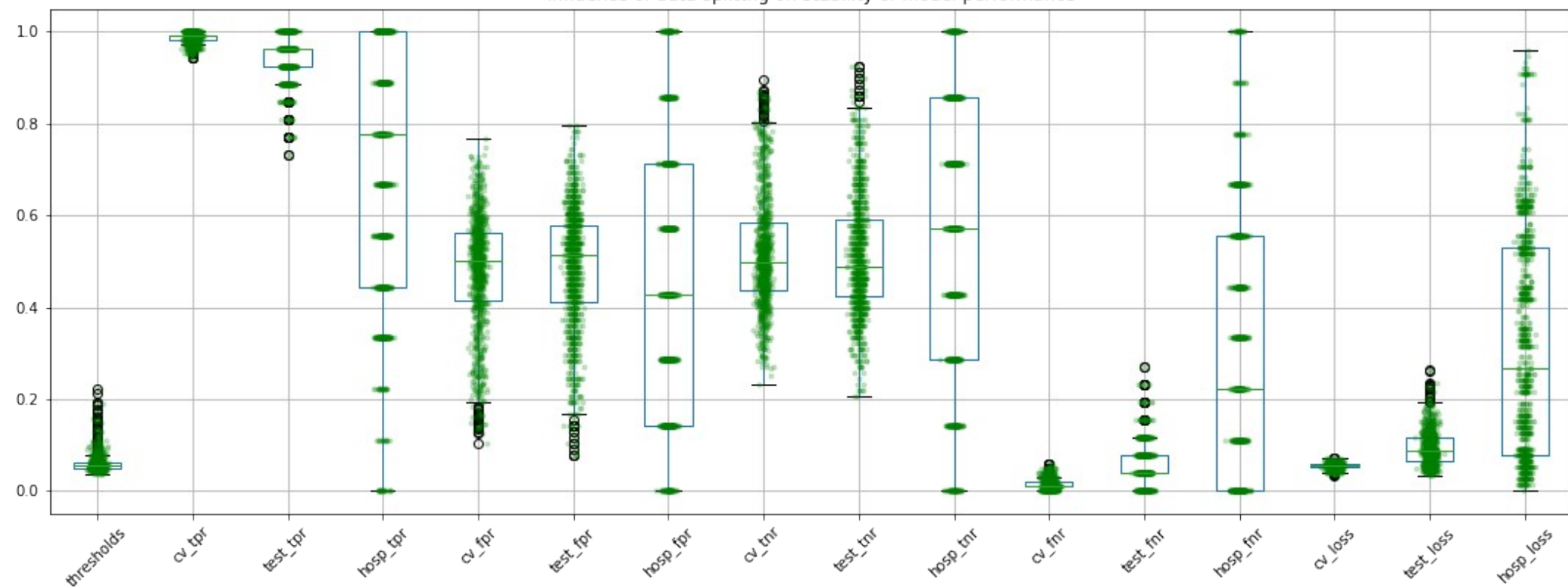




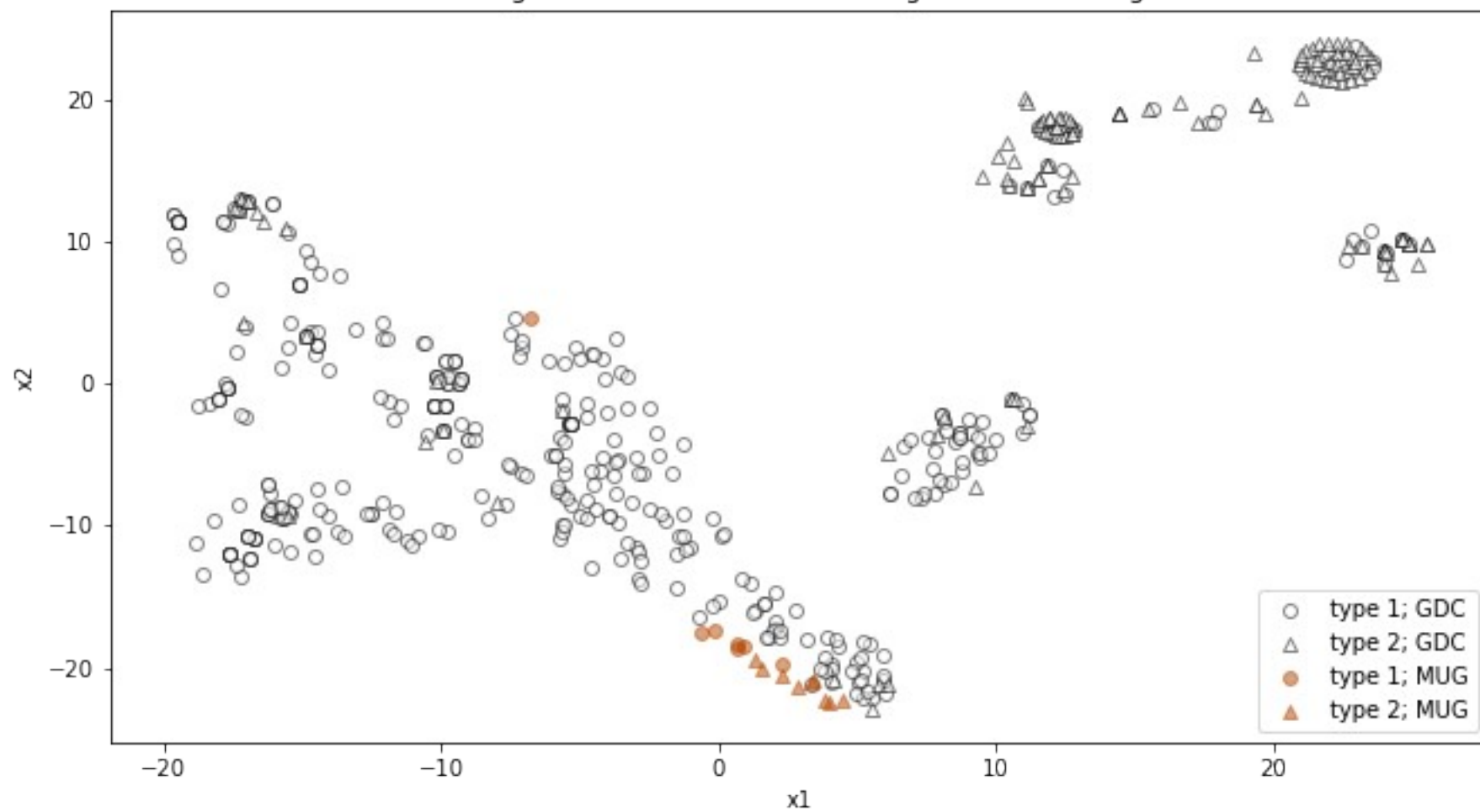




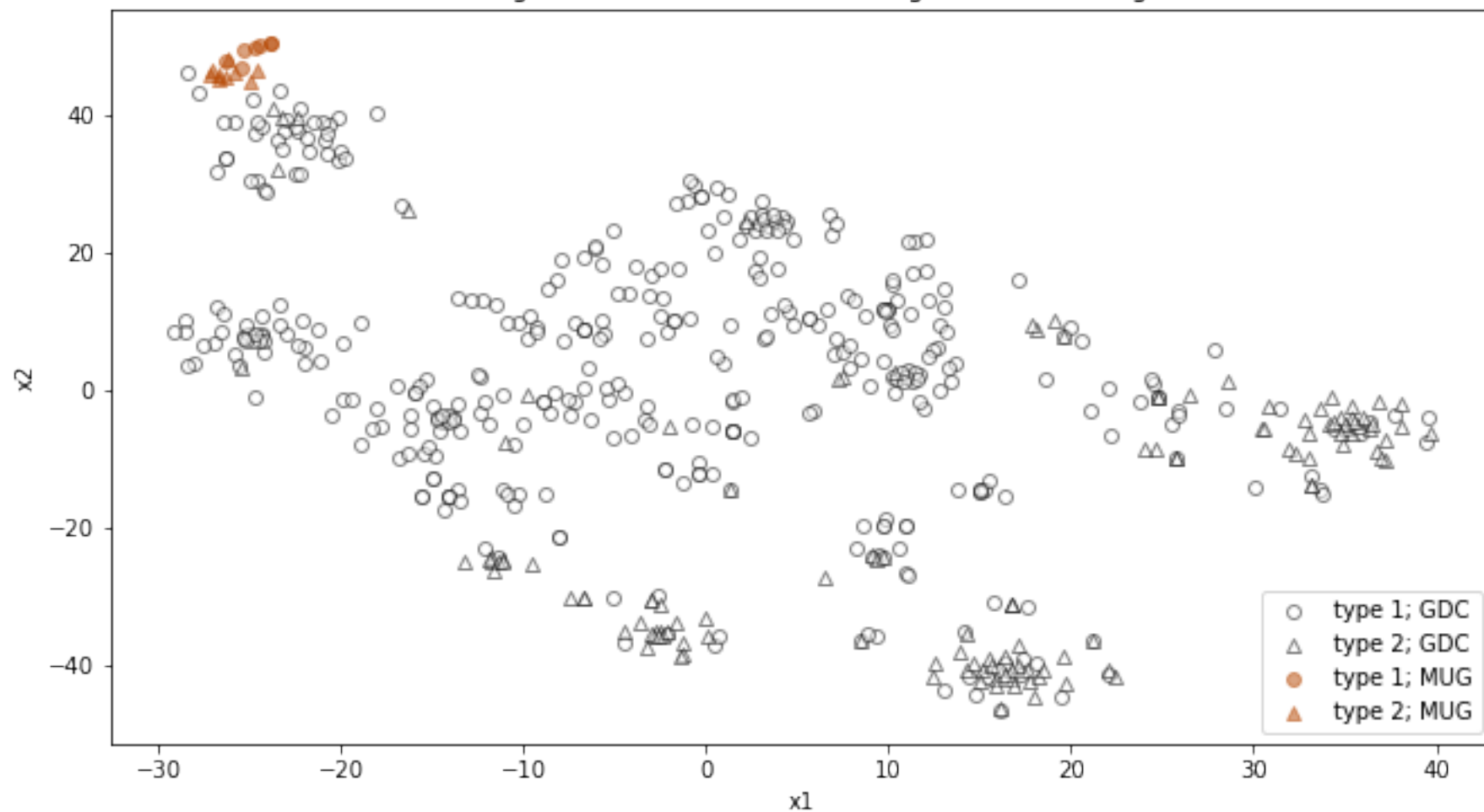
influence of data splitting on stability of model performance



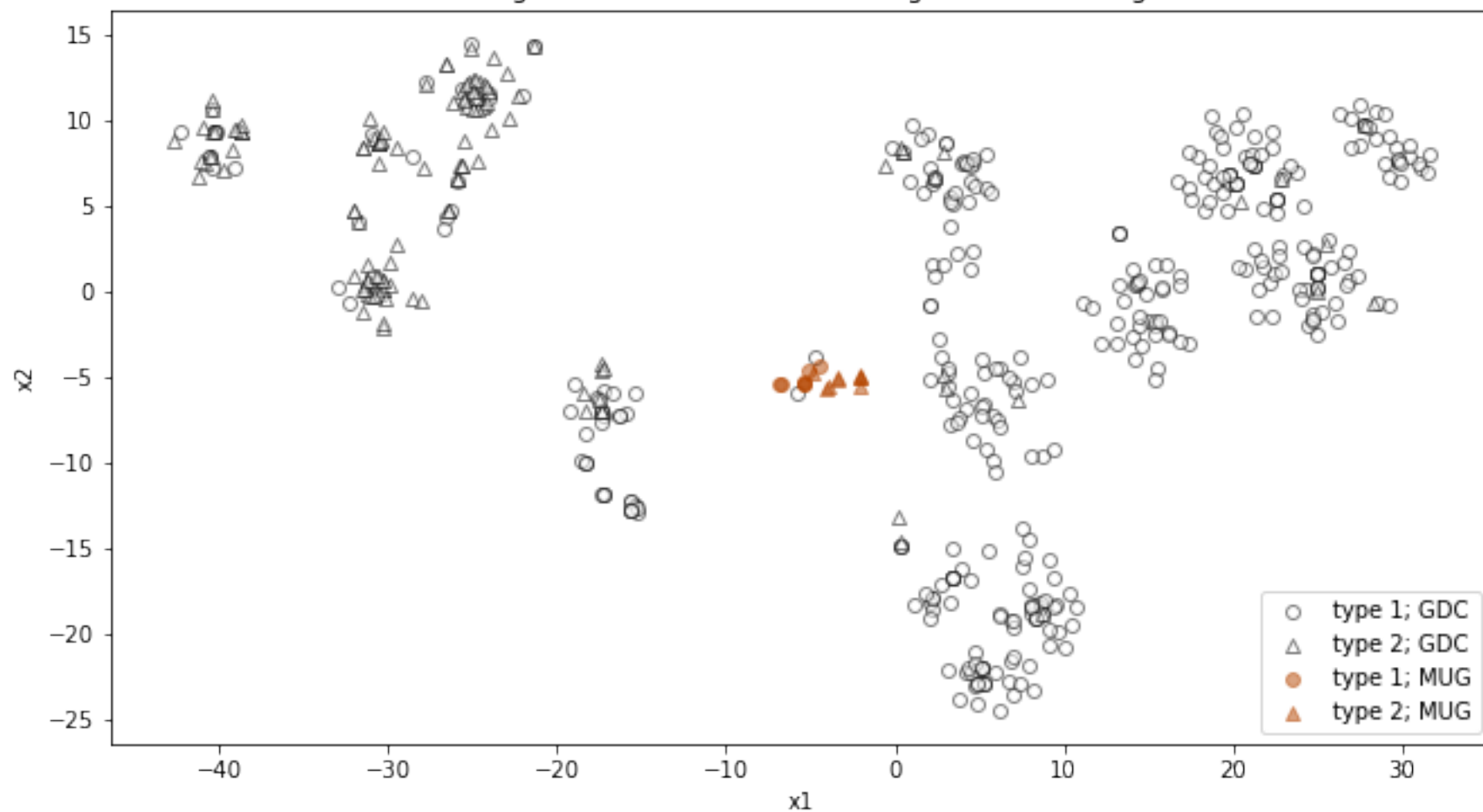
distribution of GDC and MUG datapoints projected on 2d space  
using t-distributed Stochastic Neighbor Embedding



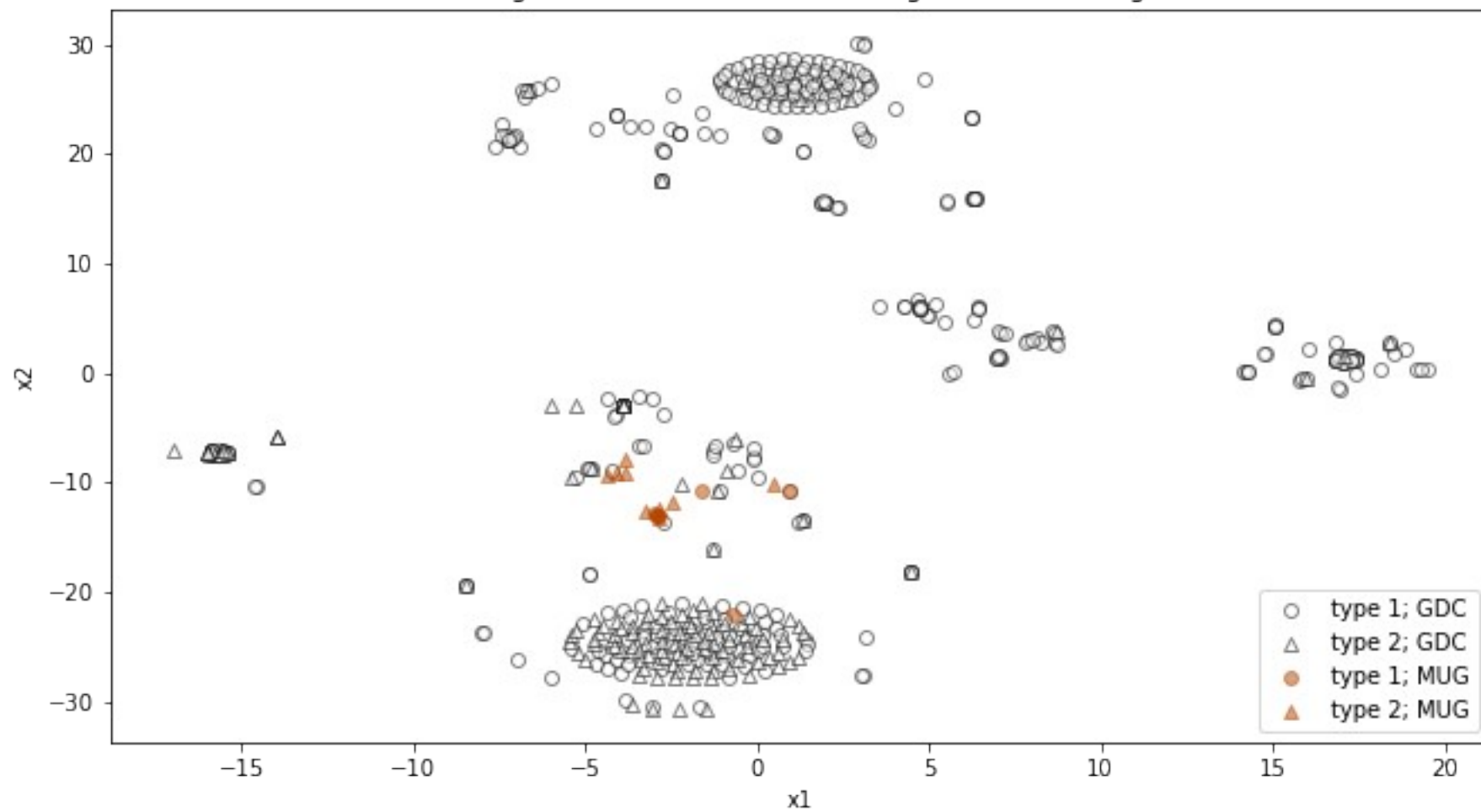
distribution of GDC and MUG datapoints projected on 2d space  
using t-distributed Stochastic Neighbor Embedding



distribution of GDC and MUG datapoints projected on 2d space  
using t-distributed Stochastic Neighbor Embedding



distribution of GDC and MUG datapoints projected on 2d space  
using t-distributed Stochastic Neighbor Embedding





distribution of GDC and MUG datapoints projected on 2d space  
using t-distributed Stochastic Neighbor Embedding

