

# Employing fingerprinting of medicinal plants by means of LC-MS and machine learning for species identification task

Pavel Kharyuk, Dmitry Nazarenko, Ivan Oseledets, Igor Rodin, Oleg Shpigun, Andrey Tsitsilin, Mikhail Lavrentyev

## Contents

<b>1</b>	<b>Confusion matrices</b>	<b>1</b>
<b>2</b>	<b>Top5 predictions</b>	<b>14</b>
<b>3</b>	<b>Hierarchical clustering analysis (HCA)</b>	<b>20</b>
<b>4</b>	<b>Sparse non-negative components</b>	<b>21</b>
<b>5</b>	<b>Autoencoder: structure, t-SNE plots and selection of last layer size</b>	<b>24</b>
<b>6</b>	<b>Dataset</b>	<b>27</b>
<b>7</b>	<b>Confusion matrices for prediction of plant parts</b>	<b>29</b>
<b>8</b>	<b>Github repository structure.</b>	<b>32</b>

## 1 Confusion matrices

In this section we provide confusion matrices measured on test1 part as mean and median of 5 times repeated 5-fold cross validation (25 runs in total). Columns: predicted labels; rows: true labels.

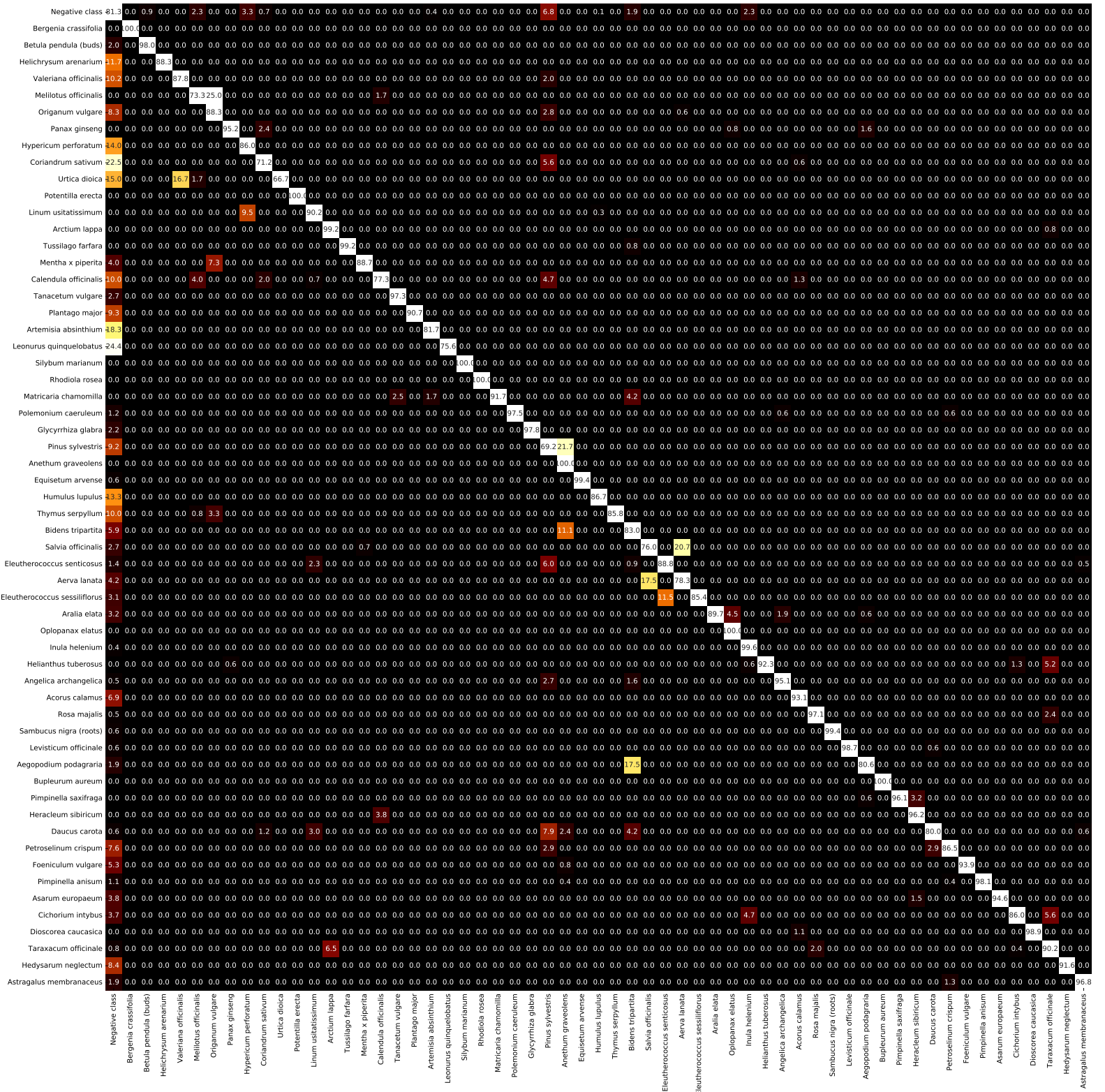


Figure S1.1(a): Mean confusion matrix for Large discrete Bayesian Network

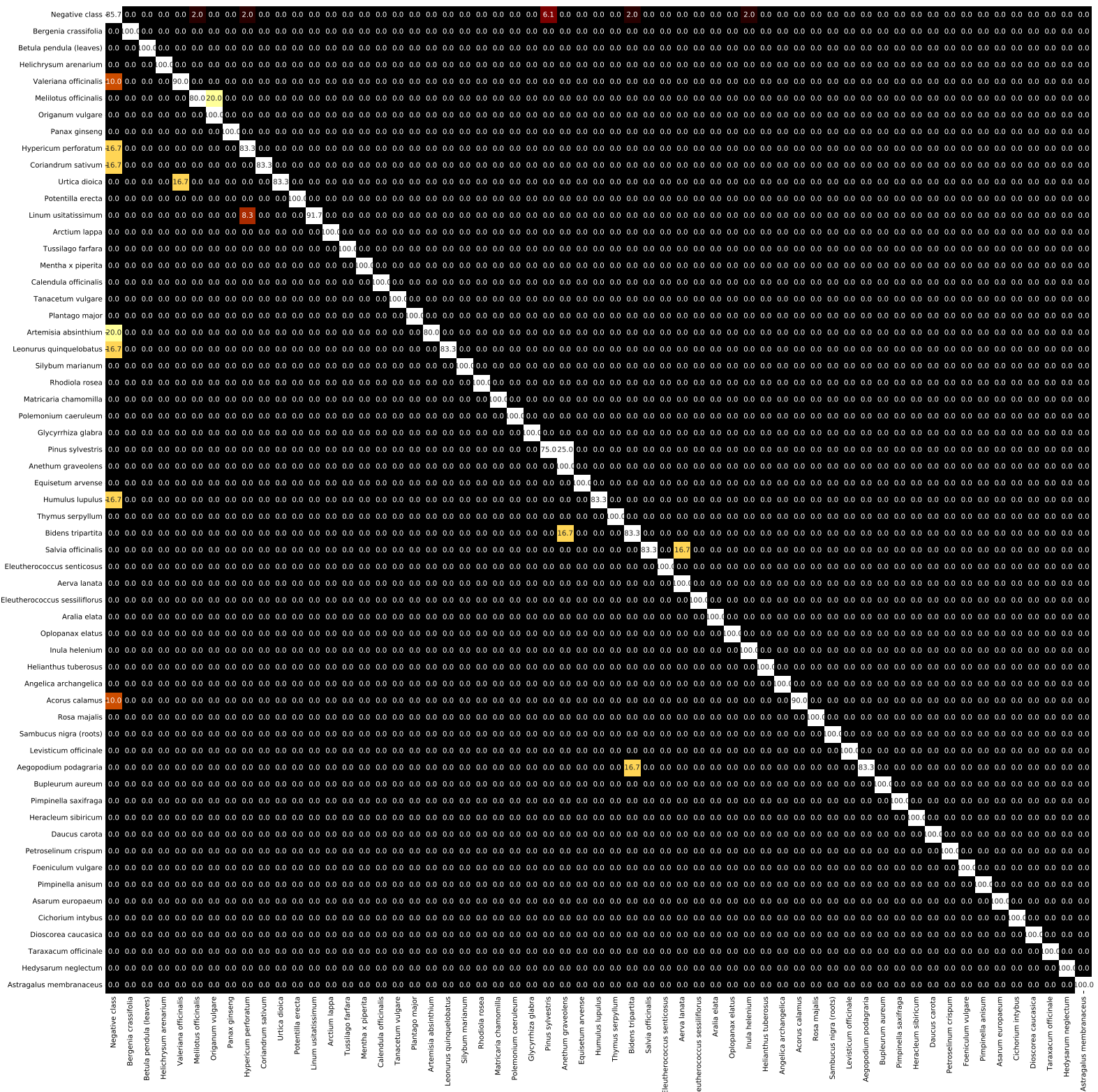


Figure S1.1(b): Median confusion matrix for Large discrete Bayesian Network

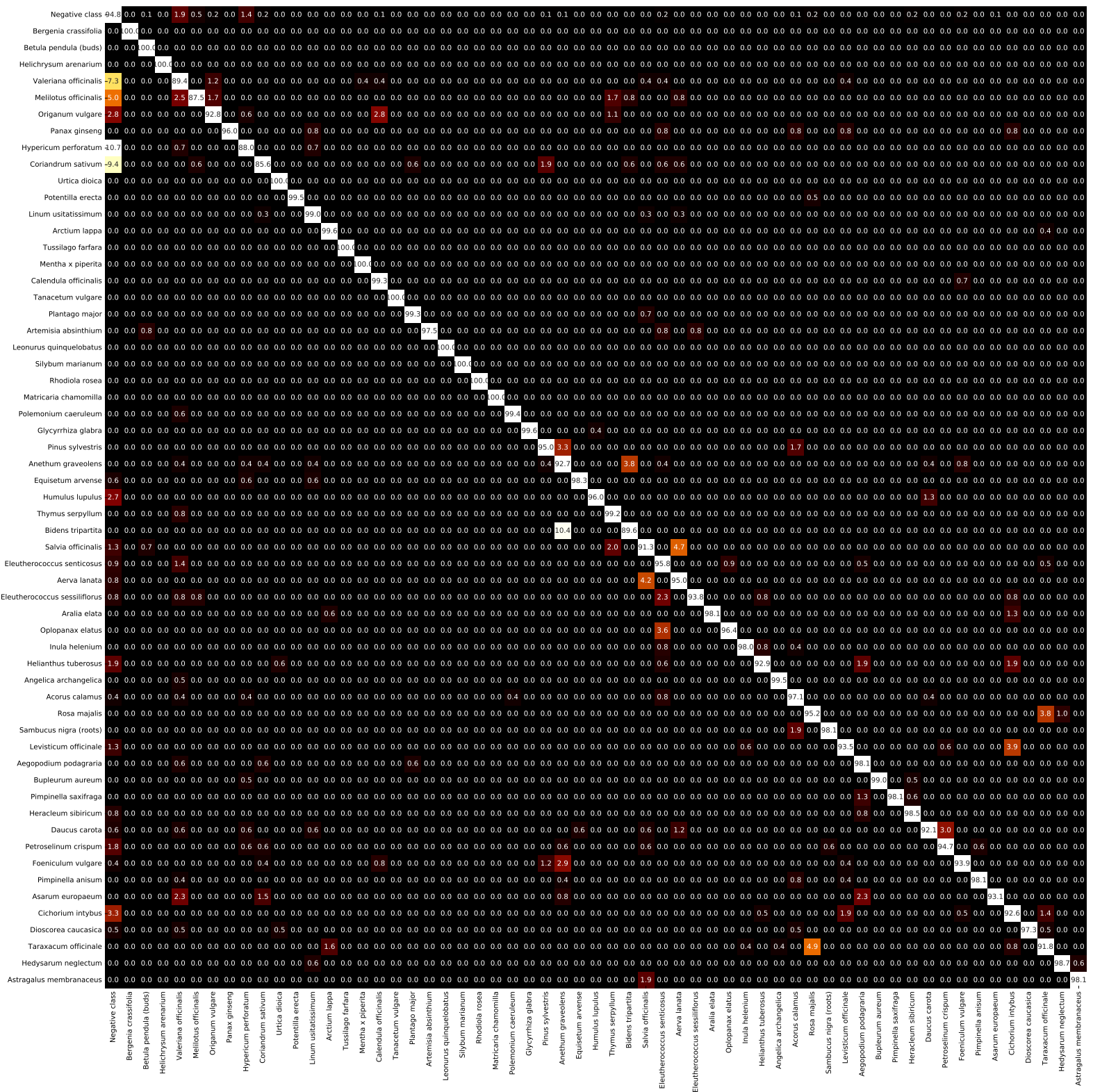


Figure S1.1(c): Mean confusion matrix for Logistic Regression trained on auto-encoded feature space

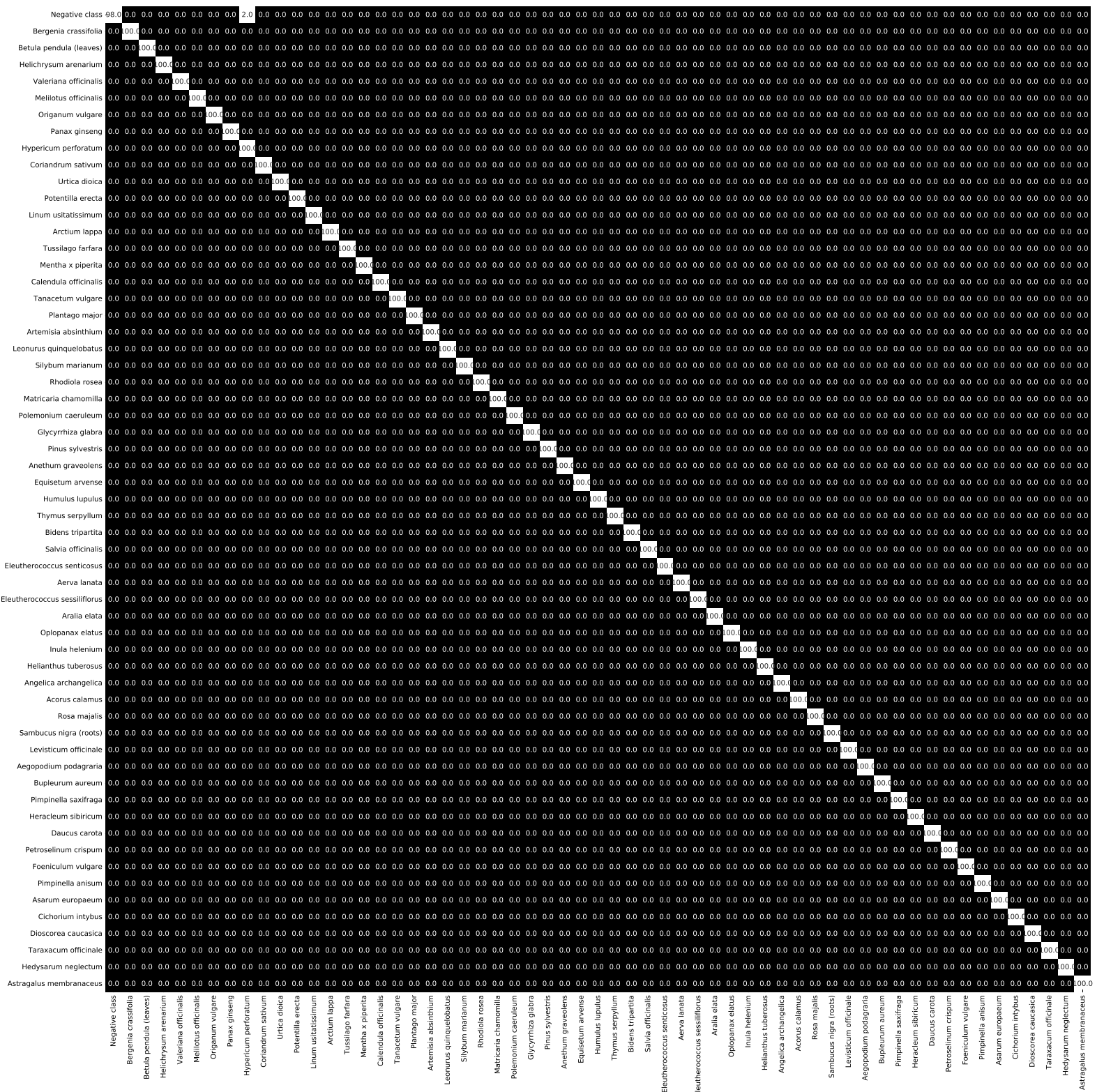


Figure S1.1(d): Median confusion matrix for Logistic Regression trained on autoencoded feature space

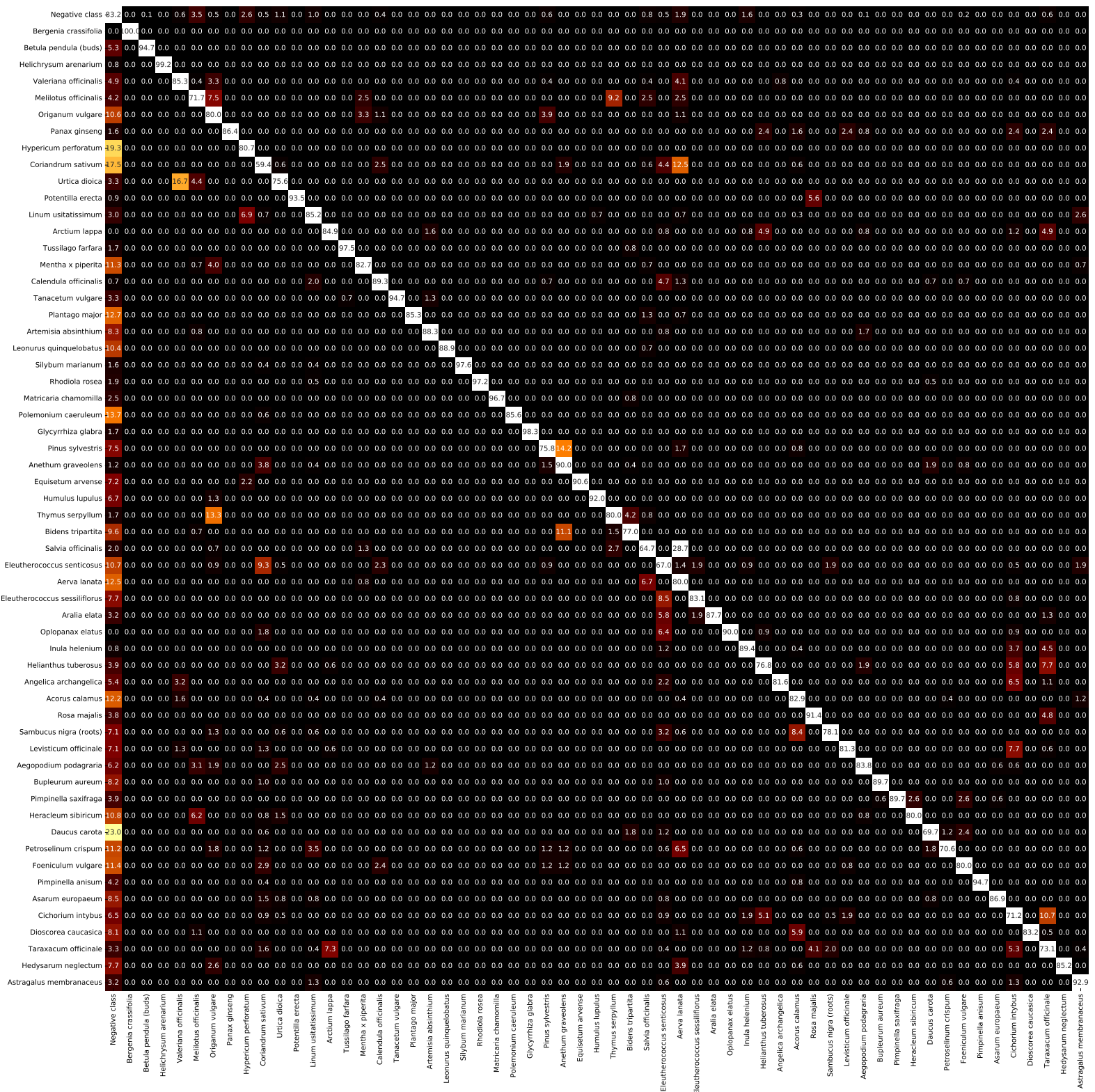


Figure S1.1(e): Mean confusion matrix for Naive Bayes classifier trained on autoencoded feature space





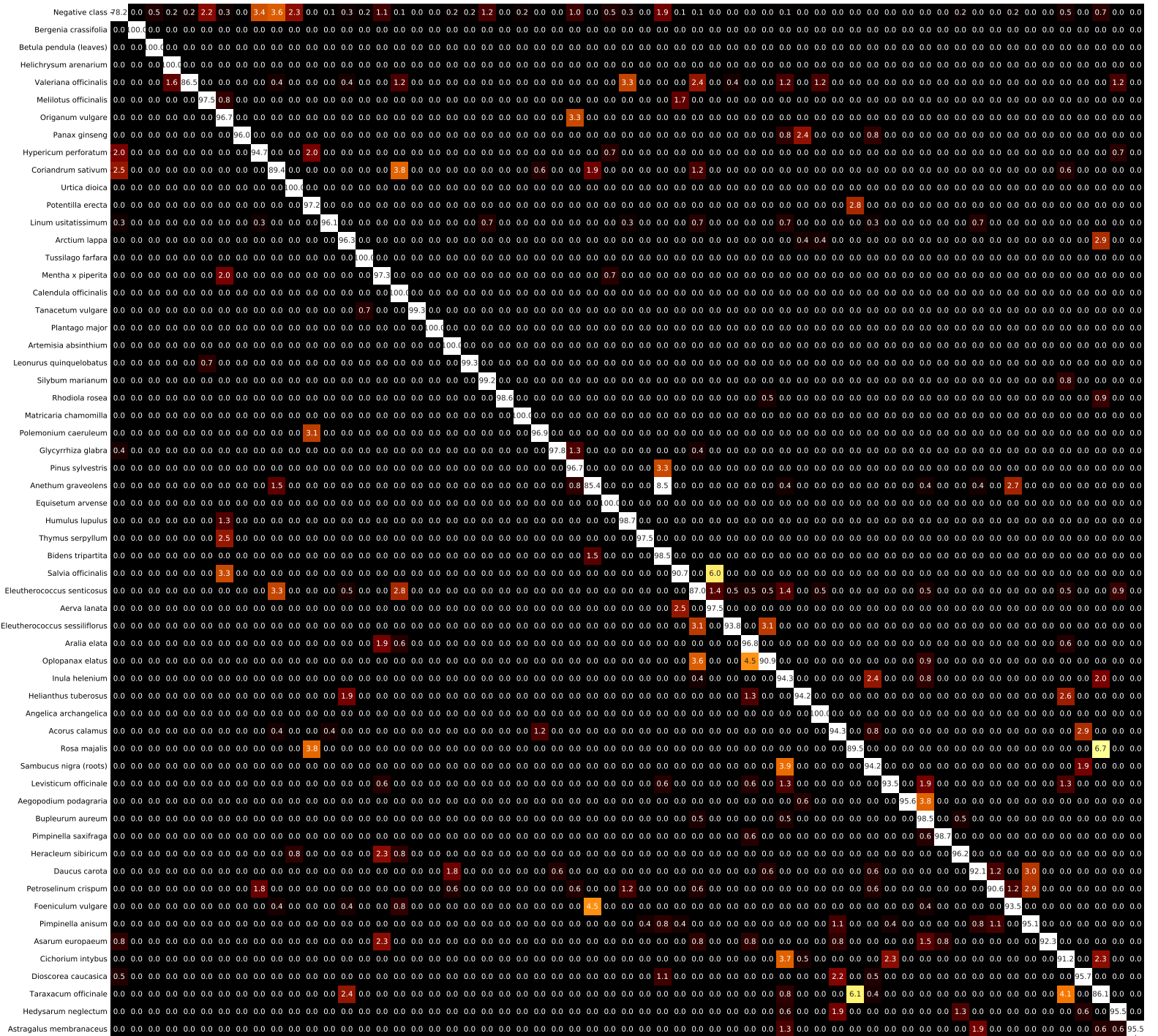


Figure S1.1(g): Mean confusion matrix for sparse non-negative Tucker decomposition based classifier (with principal angle)





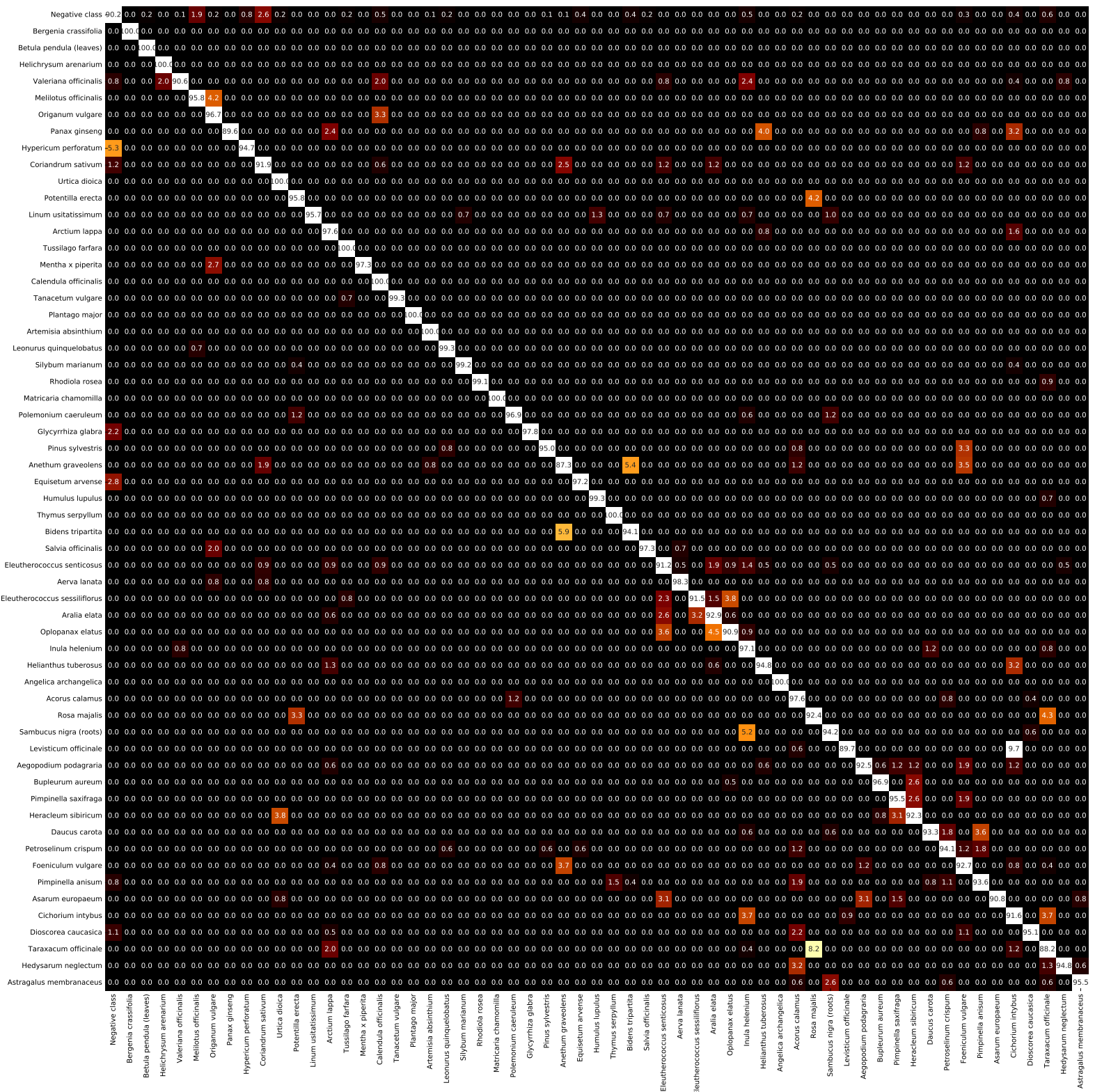


Figure S1.1(i): Mean confusion matrix for sparse non-negative matrix factorization based classifier (with principal angle)

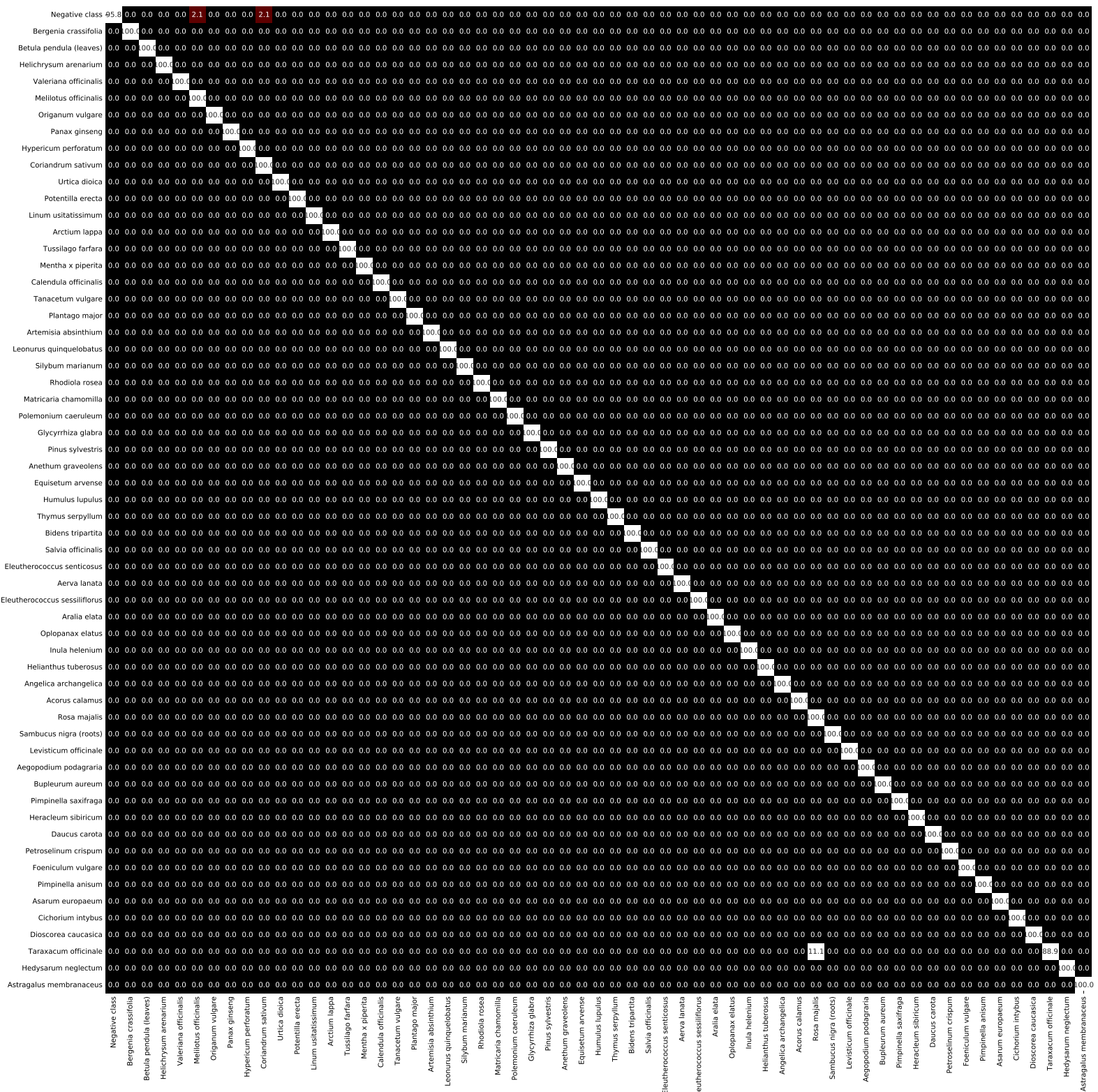


Figure S1.1(j): Median confusion matrix for sparse non-negative matrix factorization based classifier (with principal angle)

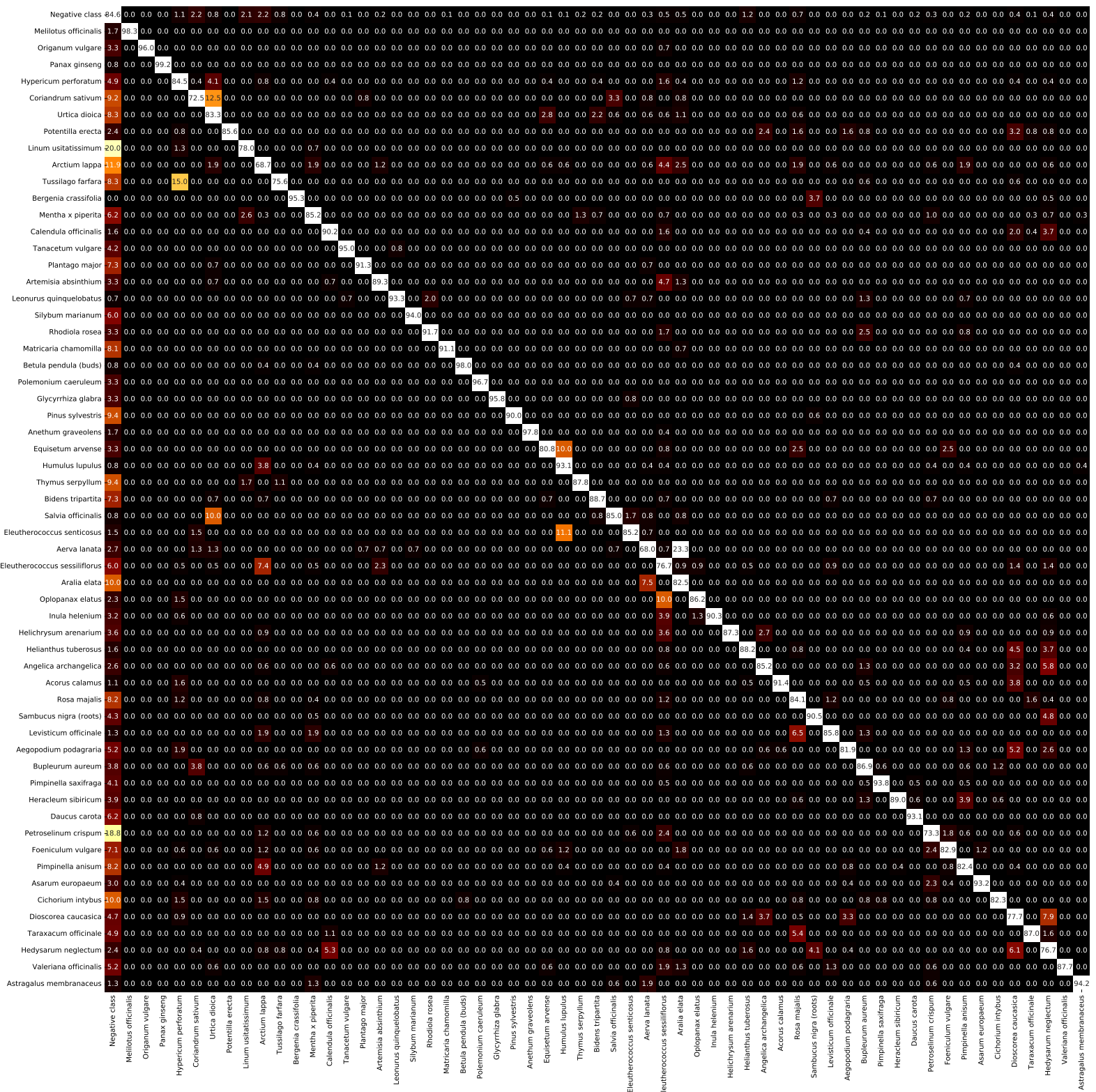


Figure S1.1(k): Mean confusion matrix for Hybrid Bayesian Network-based classifier trained on autoencoded feature space







## 2 Top5 predictions

Top5 prediction results for each sample of a class (inside each fold) were extracted and pooled together. Top5 frequent results from this list were selected and this top5s were then pooled from all 25 parts of CV (5 repetitions of 5 folds). Resulting top5 list was selected as “neighbors” for that class. Columns: true labels; rows: “neighbors”.

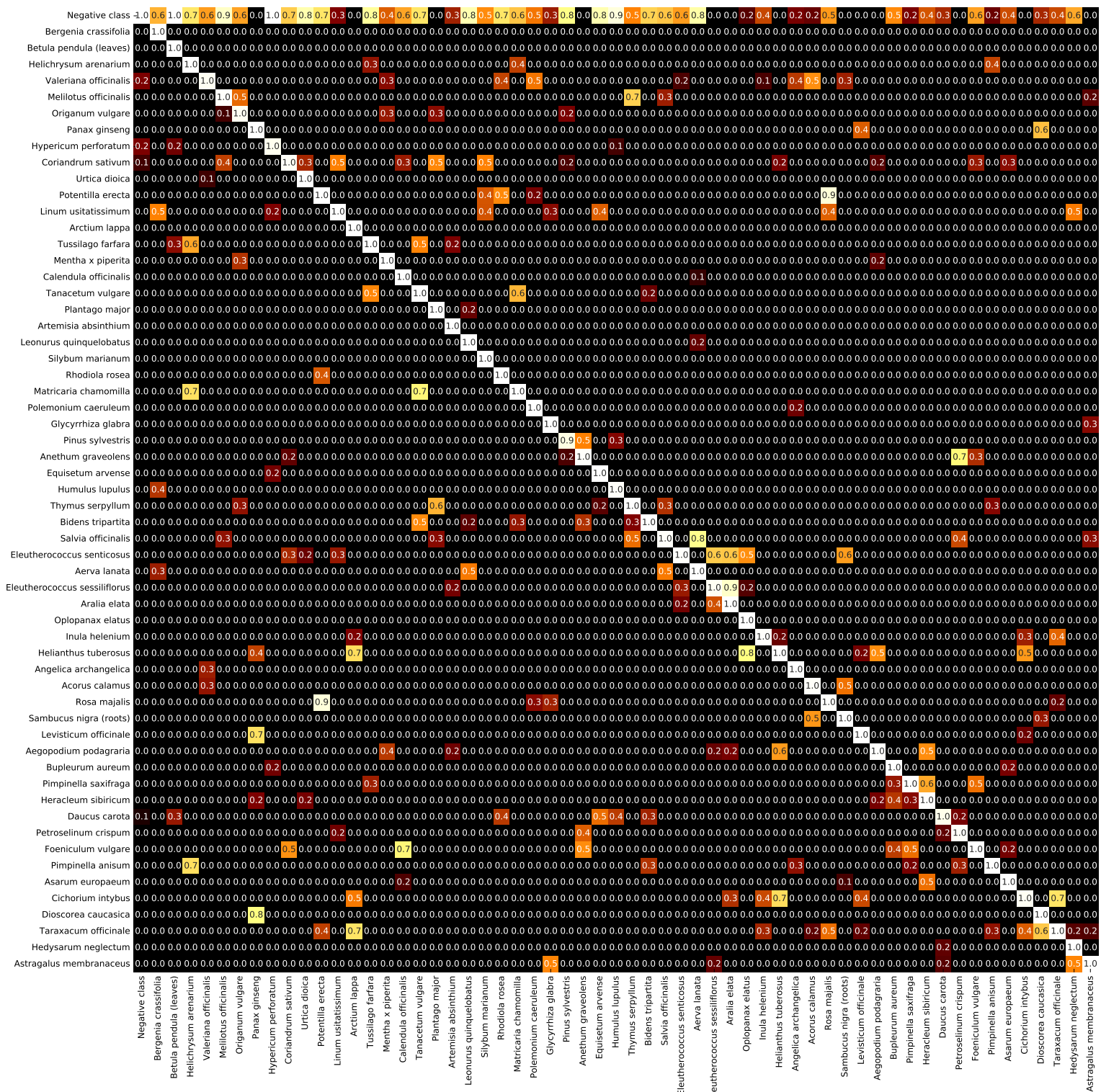


Figure S1.2(a): Logistic regression on features encoded by autoencoder



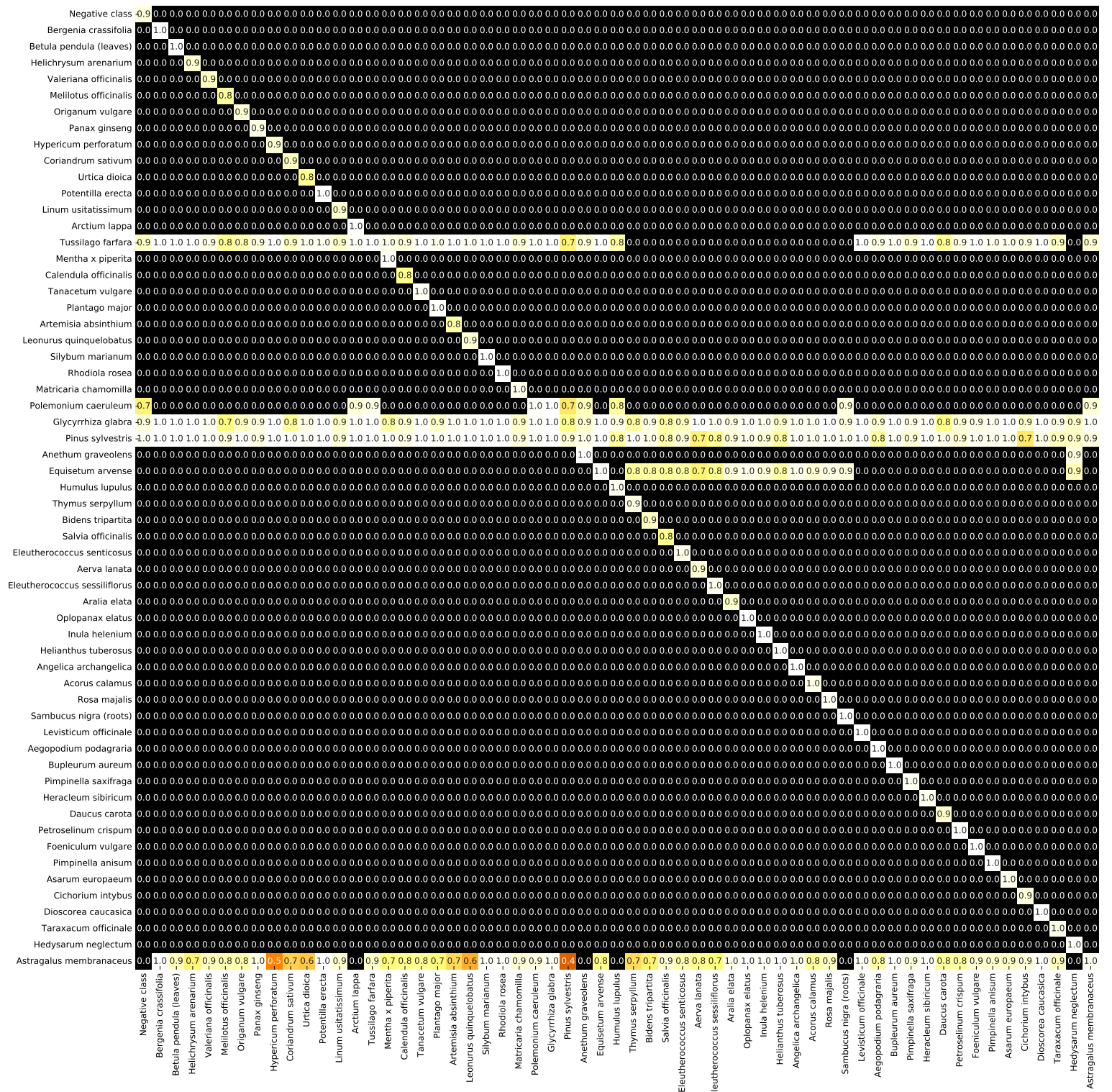


Figure S1.2(c): Large discrete Bayesian Network on binarized features



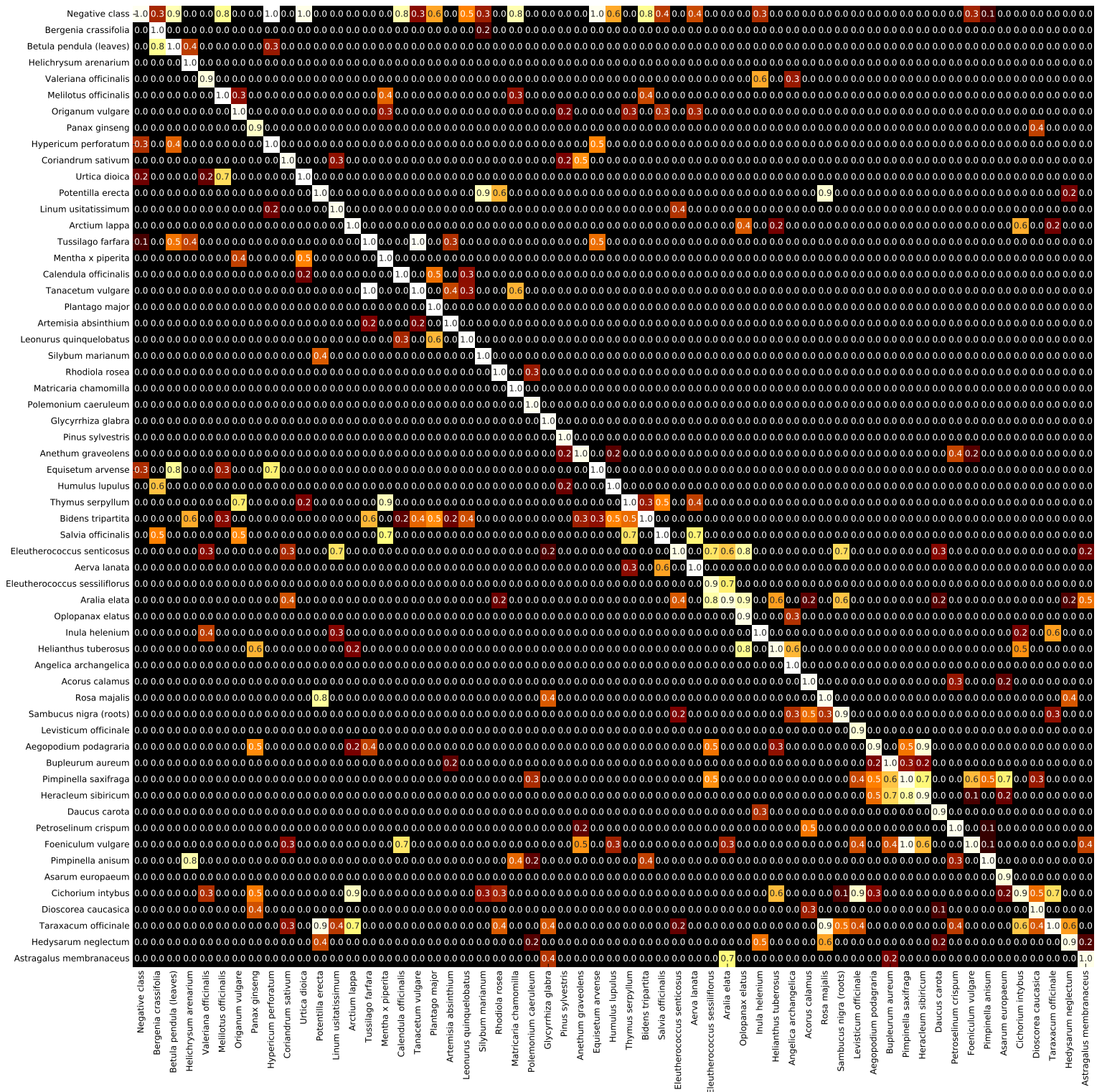


Figure S1.2(d): Sparse NMF classifier



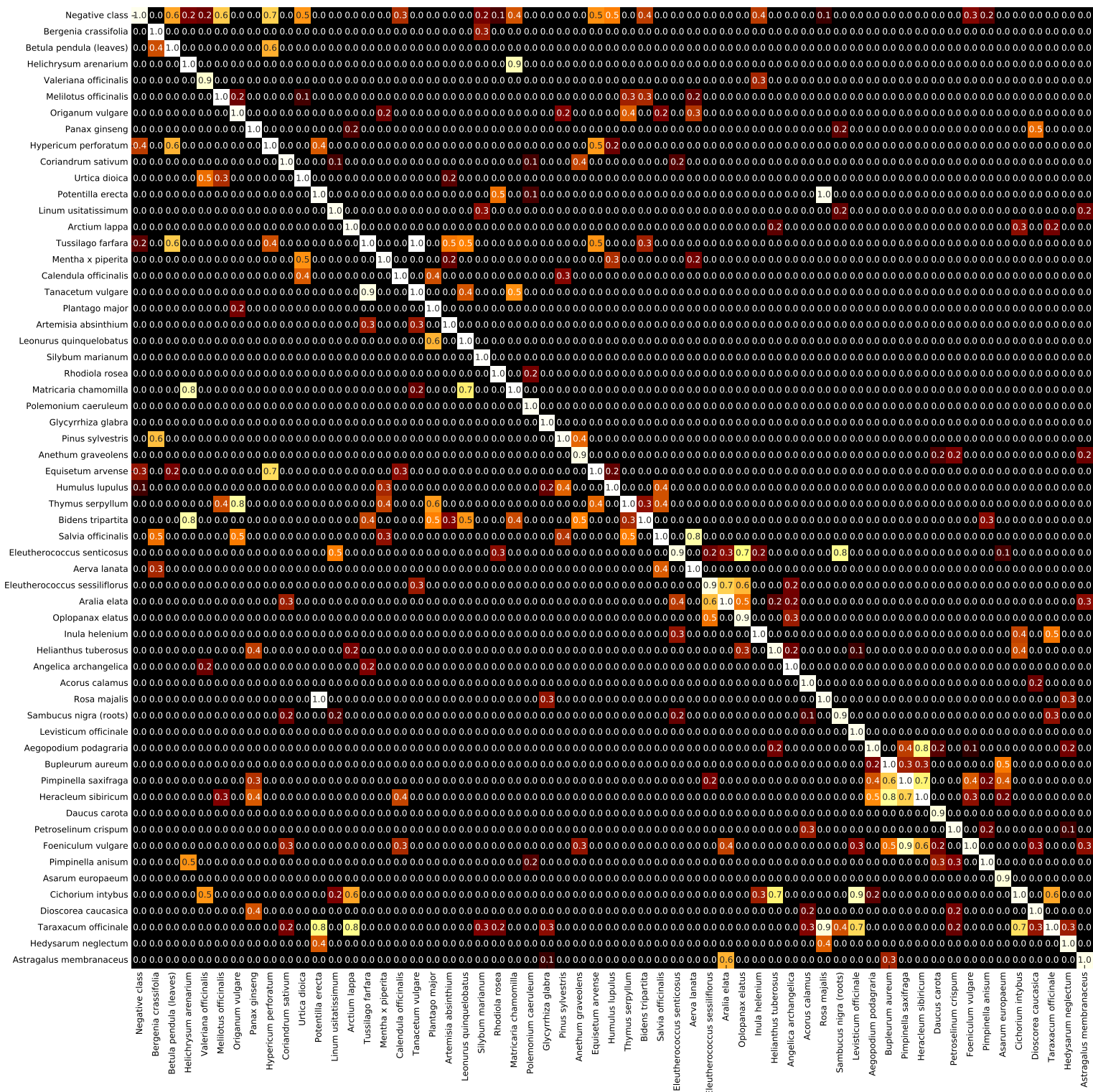


Figure S1.2(e): Sparse NTD classifier

### 3 Hierarchical clustering analysis (HCA)

HCA has been performed with hierarchical agglomerative clustering algorithm. Two variants of analysis have been performed: (1) clustering of mean samples for each of 76 classes (74 species); (2) clustering of all samples.

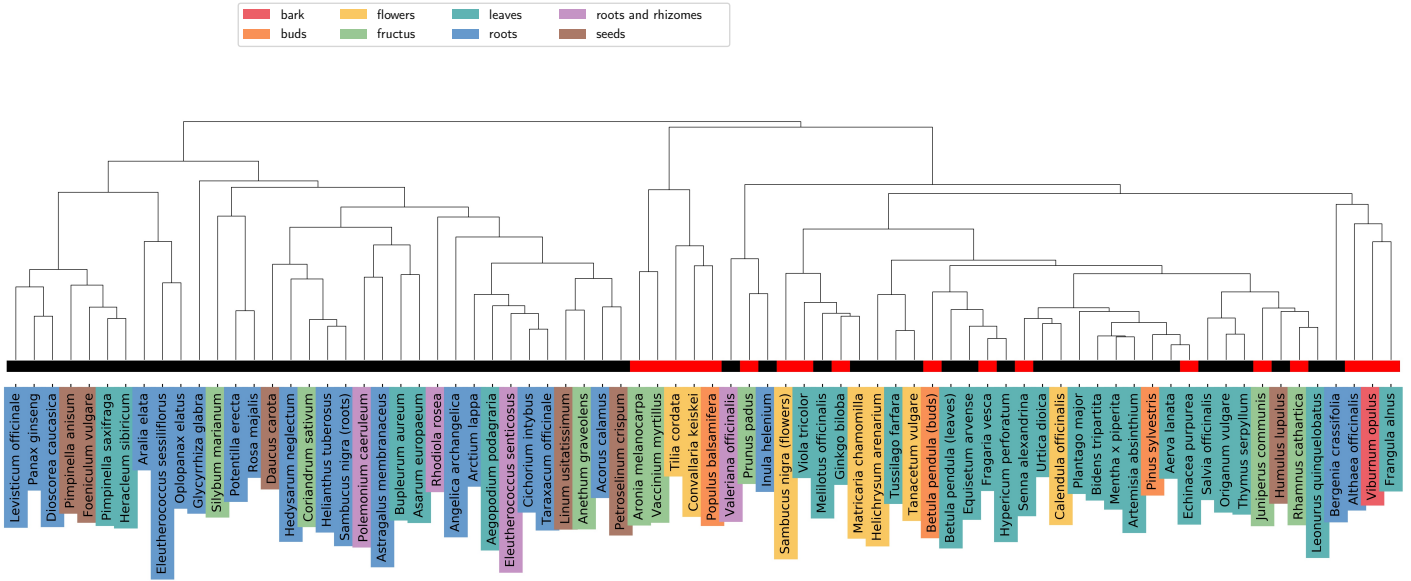


Figure S1.3(a): HCA of mean samples in the original feature space (linkage: complete, metric: correlation)

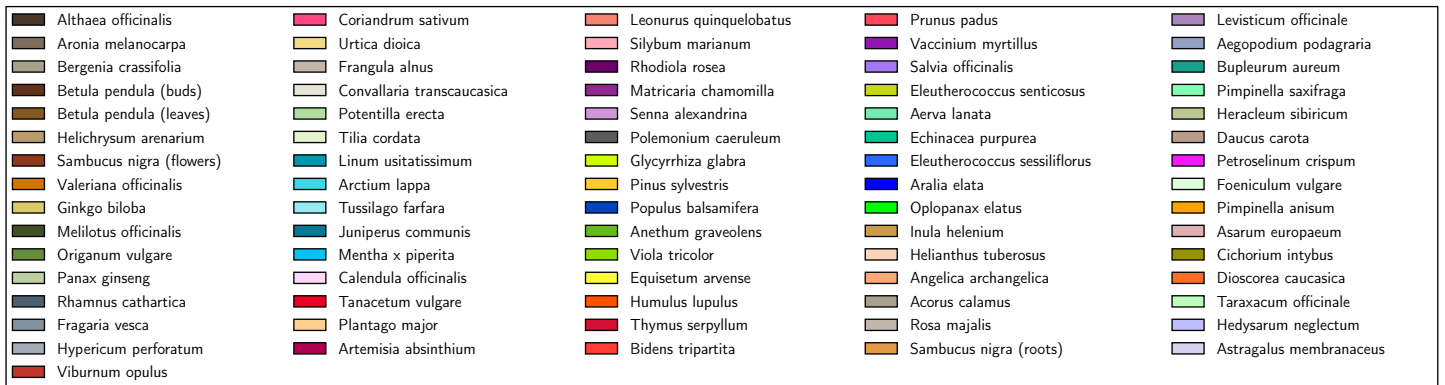


Figure S1.3(b): Colour codes for 76 classes

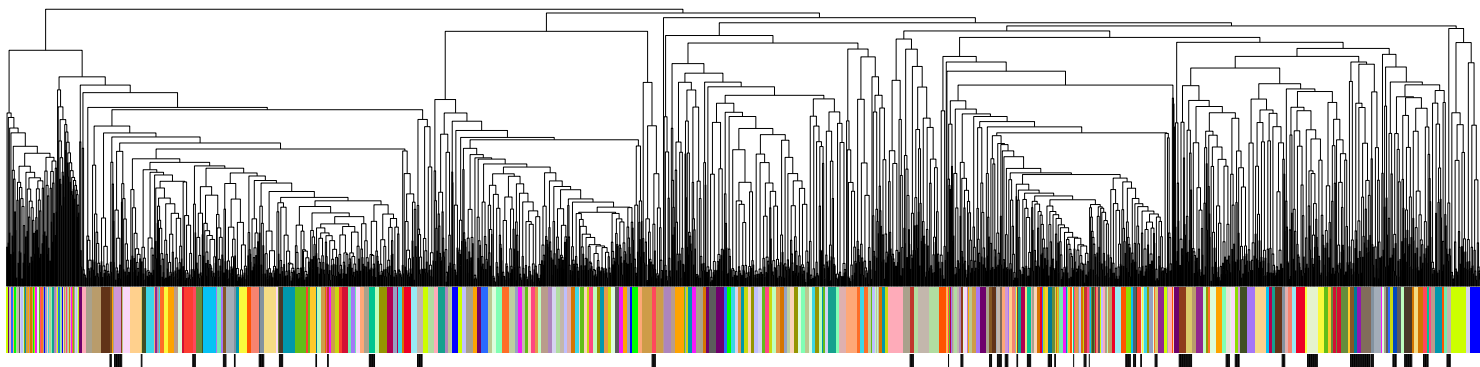


Figure S1.3(c): HCA of all samples in the original feature space (linkage: weighted, metric: euclidean); black vertical lines indicate samples from the negative class

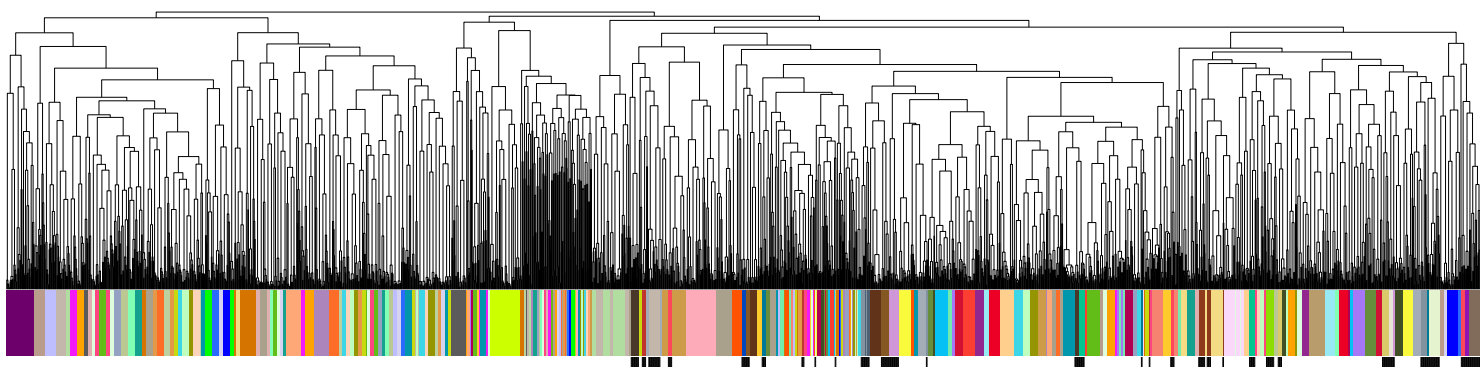


Figure S1.3(d): HCA of whole samples in the autoencoded feature space (linkage: weighted, metric: euclidean); black vertical lines indicate samples from the negative class

## 4 Sparse non-negative components

Components obtained after sparse non-negative matrix factorization (SNMF) and sparse non-negative Tucker decomposition (SNTD). On Figures S1.4(a), S1.4(b) first three components extracted by SNMF (and then splitted into positive and negative polarity parts) and SNTD algorithms are displayed. The following criteria was used to select three components:

$$\min_{i=1,r;i \notin \Omega} \left( \frac{1}{n_c - 1} \sum_{j \neq l} |B[j, i]| \right) - \log(|A[1, i]| + \varepsilon), \quad (1)$$

where  $\Omega$  is a set of already selected components,  $B \in \mathbb{R}^{n_c-1 \times r_l}$  - medians (taken by sample axis) of m/z-wise correlation matrix between components for current class  $l$  and samples drawn from other classes, and  $A \in \mathbb{R}^{1 \times r_l}$  is a vector of sample-wise medians of correlation matrix between components of current class  $l$  and samples drawn from current class,  $n_c$  - summarized number of classes.

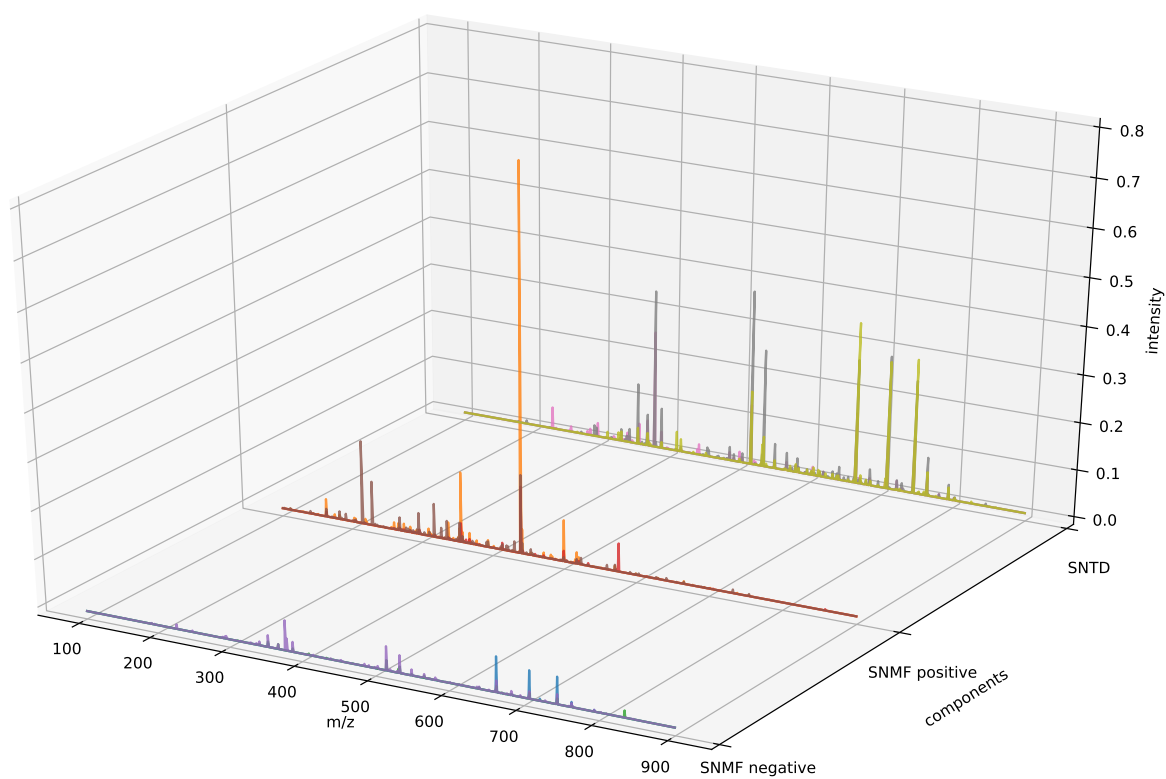


Figure S1.4(a): *Aralia elata*. Colours indicate separate components.

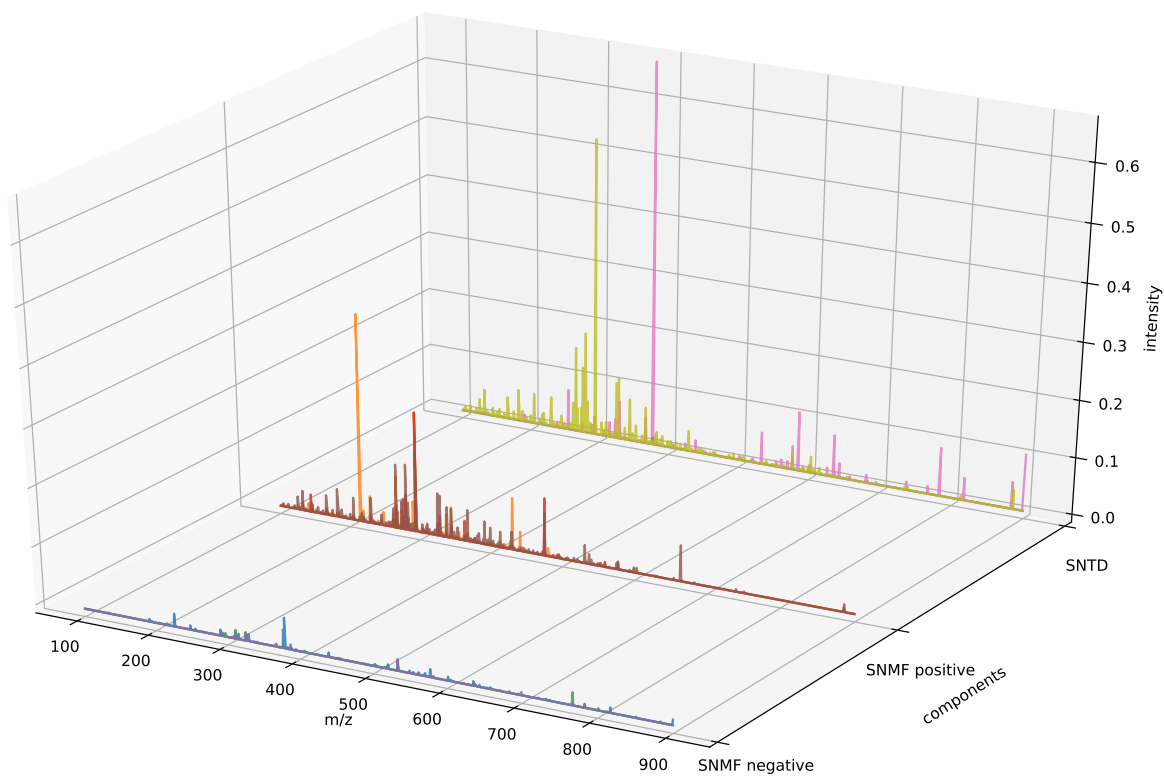


Figure S1.4(b): *Dioscorea caucasica*. Colours indicate separate components.



## 5 Autoencoder: structure, t-SNE plots and selection of last layer size

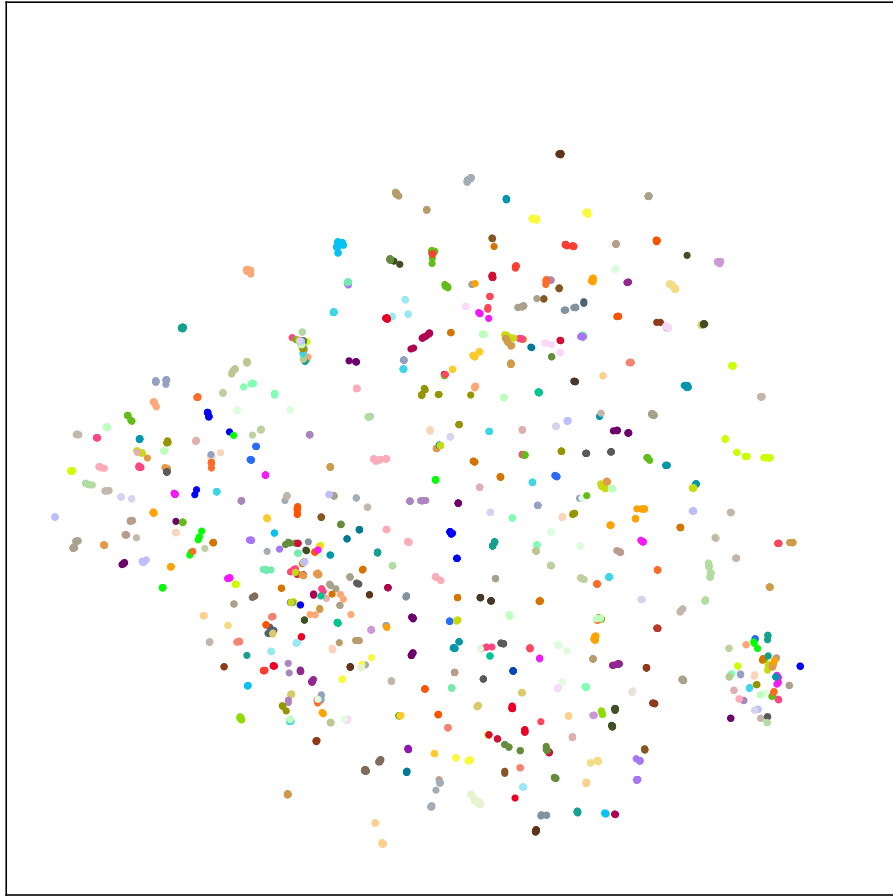


Figure S1.5(a): t-SNE plot for the main dataset with original feature space (1600 variables). Colour codes are the same as in Figure S1.3(b). Perplexity=10.

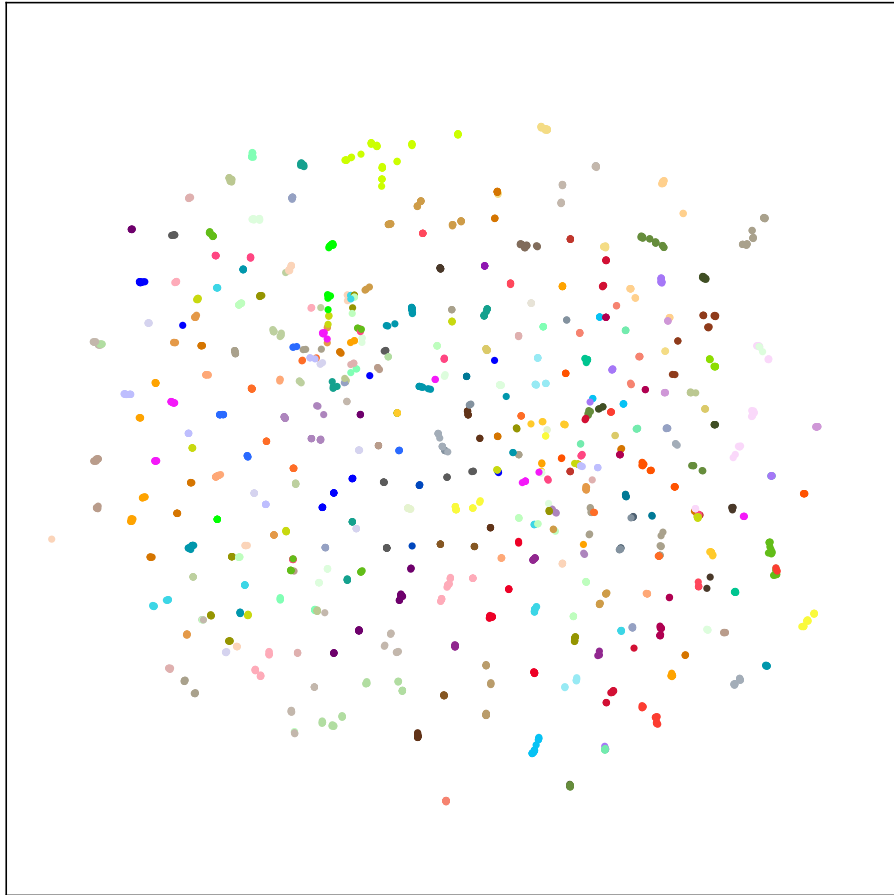


Figure S1.5(b): t-SNE plot for the main dataset with autoencoded feature space (25 variables). Colour codes are the same as in Figure S1.3(b). Perplexity=10.

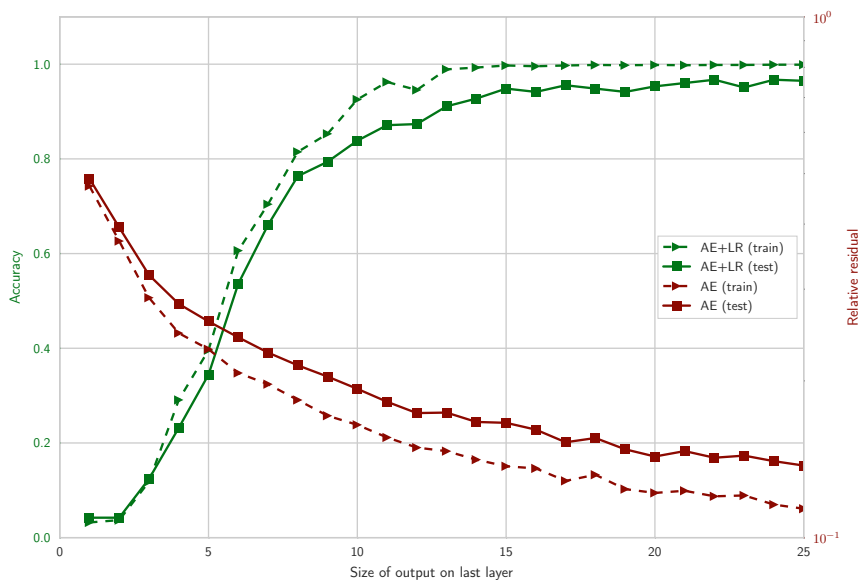


Figure S1.5(c): Accuracy on train/test1 parts of the dataset on 5-fold CV (green lines) and median of relative residual error among samples (red lines) depending on last layer size

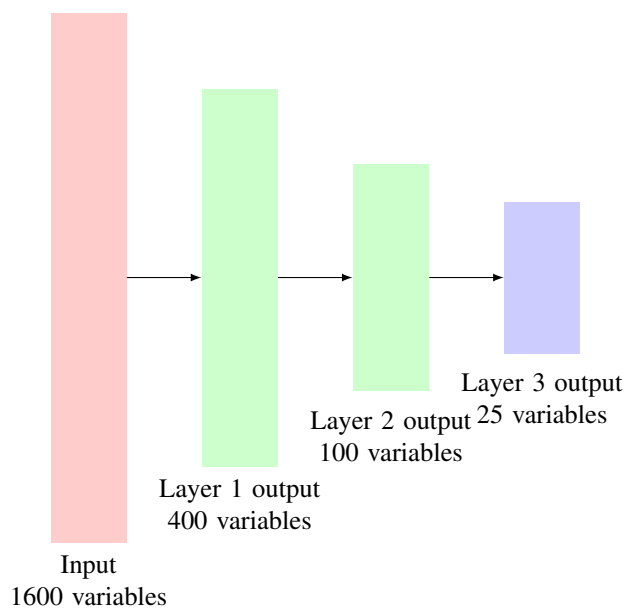


Figure S1.5(d): Final structure of autoencoder's encoding part

## 6 Dataset

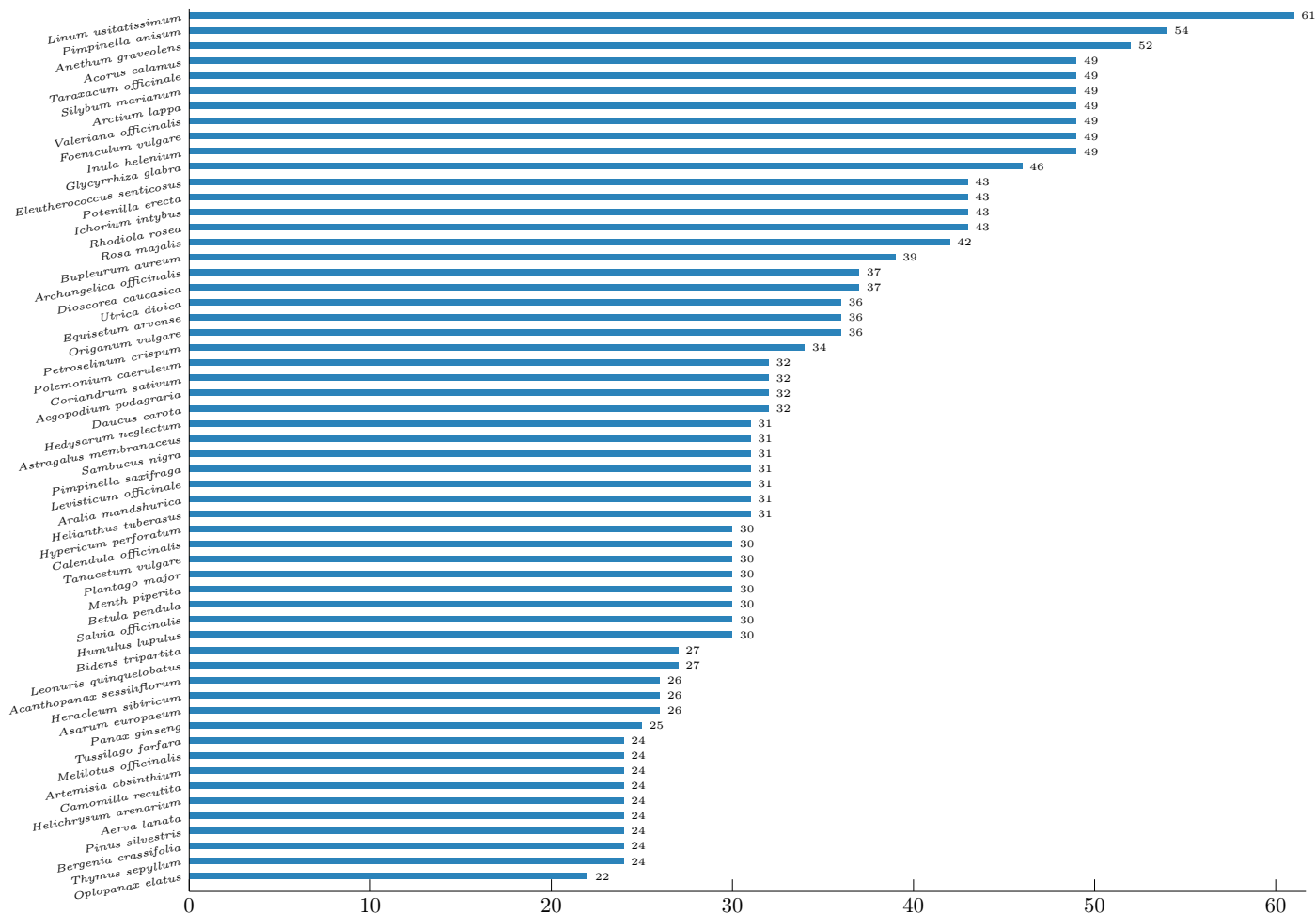


Figure S1.6(a): Dataset composition by plant species. Only 58 species with 20 or more available chromatograms are displayed.

Table S1.6(a): Plant species used in experiment and corresponding class labels. Corresponding file in computer readable format is available at Github repository <https://github.com/kharyuk/chemfin-plasp> as data/species.csv .

Class	Plant species	Organ	Class	Plant species	Organ
1	<i>Althaea officinalis</i>	roots	41	<i>Glycyrrhiza glabra</i>	roots
2	<i>Aronia melanocarpa</i>	fructus	42	<i>Pinus sylvestris</i>	buds
3	<i>Bergenia crassifolia</i>	roots	43	<i>Populus balsamifera</i>	buds
4	<i>Betula pendula</i>	leaves	44	<i>Anethum graveolens</i>	fructus
5	<i>Betula pendula</i>	buds	45	<i>Viola tricolor</i>	leaves
6	<i>Helichrysum arenarium</i>	flowers	46	<i>Equisetum arvense</i>	leaves
7	<i>Sambucus nigra</i>	flowers	47	<i>Humulus lupulus</i>	seeds
8	<i>Valeriana officinalis</i>	roots, rhizomes	48	<i>Thymus serpyllum</i>	leaves
9	<i>Ginkgo biloba</i>	leaves	49	<i>Bidens tripartita</i>	leaves
10	<i>Melilotus officinalis</i>	leaves	50	<i>Prunus padus</i>	fructus
11	<i>Origanum vulgare</i>	leaves	51	<i>Vaccinium myrtillus</i>	fructus
12	<i>Panax ginseng</i>	roots	52	<i>Salvia officinalis</i>	leaves
13	<i>Rhamnus cathartica</i>	fructus	53	<i>Eleutherococcus senticosus</i>	roots, rhizomes
14	<i>Fragaria vesca</i>	leaves	54	<i>Aerva lanata</i>	leaves
15	<i>Hypericum perforatum</i>	leaves	55	<i>Echinacea purpurea</i>	leaves
16	<i>Viburnum opulus</i>	bark	56	<i>Eleutherococcus sessiliflorus</i>	roots
17	<i>Coriandrum sativum</i>	fructus	57	<i>Aralia elata</i>	roots
18	<i>Urtica dioica</i>	leaves	58	<i>Oplopanax elatus</i>	roots
20	<i>Frangula alnus</i>	leaves	59	<i>Inula helenium</i>	roots
23	<i>Convallaria keiskei</i>	flowers	60	<i>Helianthus tuberosus</i>	roots
24	<i>Potentilla erecta</i>	roots	61	<i>Angelica archangelica</i>	roots
25	<i>Tilia cordata</i>	flowers	62	<i>Acorus calamus</i>	roots
26	<i>Linum usitatissimum</i>	seeds	63	<i>Rosa majalis</i>	roots
27	<i>Arctium lappa</i>	roots	64	<i>Sambucus nigra</i>	roots
28	<i>Tussilago farfara</i>	leaves	65	<i>Levisticum officinale</i>	roots
29	<i>Juniperus communis</i>	fructus	66	<i>Aegopodium podagraria</i>	leaves
30	<i>Mentha x piperita</i>	leaves	67	<i>Bupleurum aureum</i>	leaves
31	<i>Calendula officinalis</i>	flowers	68	<i>Pimpinella saxifraga</i>	leaves
32	<i>Tanacetum vulgare</i>	flowers	69	<i>Heracleum sibiricum</i>	leaves
33	<i>Plantago major</i>	leaves	70	<i>Daucus carota</i>	seeds
34	<i>Artemisia absinthium</i>	leaves	71	<i>Petroselinum crispum</i>	seeds
35	<i>Leonurus quinquelobatus</i>	leaves	72	<i>Foeniculum vulgare</i>	seeds
36	<i>Silybum marianum</i>	fructus	73	<i>Pimpinella anisum</i>	seeds
37	<i>Rhodiola rosea</i>	roots, rhizomes	75	<i>Asarum europaeum</i>	leaves
38	<i>Matricaria chamomilla</i>	flowers	76	<i>Cichorium intybus</i>	roots
39	<i>Senna alexandrina</i>	leaves	77	<i>Dioscorea caucasica</i>	roots
40	<i>Polemonium caeruleum</i>	roots, rhizomes	78	<i>Taraxacum officinale</i>	roots
			79	<i>Hedysarum neglectum</i>	roots
			80	<i>Astragalus membranaceus</i>	roots



## 7 Confusion matrices for prediction of plant parts

In this section we provide confusion matrices measured on test1 part as medians of 5 times repeated 5-fold cross validation (25 runs in total). Columns: predicted labels; rows: true labels.

bark -	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
buds	0.00	0.90	0.00	0.10	0.00	0.00	0.00	0.00
flowers	0.00	0.00	0.90	0.00	0.10	0.00	0.00	0.00
fructus	0.00	0.00	0.00	0.61	0.11	0.14	0.00	0.14
leaves	0.00	0.00	0.01	0.01	0.93	0.02	0.01	0.01
roots	0.00	0.00	0.00	0.01	0.02	0.92	0.03	0.02
roots, rhizomes	0.00	0.00	0.00	0.06	0.03	0.28	0.62	0.00
seeds	0.00	0.00	0.00	0.10	0.12	0.04	0.00	0.73
	bark	buds	flowers	fructus	leaves	roots	roots, rhizomes	seeds

Figure S1.7(a): Median confusion matrix for Logistic Regression trained on autoencoded feature space.

bark -	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
buds	0.00	0.88	0.00	0.12	0.00	0.00	0.00	0.00
flowers	0.00	0.00	0.86	0.00	0.14	0.00	0.00	0.00
fructus	0.00	0.03	0.00	0.39	0.17	0.17	0.00	0.25
leaves	0.00	0.01	0.04	0.01	0.86	0.06	0.01	0.01
roots	0.00	0.00	0.00	0.03	0.07	0.86	0.03	0.01
roots, rhizomes	0.00	0.00	0.00	0.06	0.06	0.25	0.62	0.00
seeds	0.00	0.02	0.04	0.08	0.18	0.14	0.00	0.55
	bark	buds	flowers	fructus	leaves	roots	roots, rhizomes	seeds

Figure S1.7(b): Median confusion matrix for Naive Bayes classifier trained on autoencoded feature space.

bark -	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
buds	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00
flowers	0.00	0.00	0.83	0.00	0.17	0.00	0.00	0.00
fructus	0.00	0.00	0.00	0.55	0.16	0.16	0.03	0.11
leaves	0.00	0.01	0.05	0.02	0.83	0.05	0.01	0.03
roots	0.00	0.00	0.00	0.02	0.04	0.86	0.04	0.04
roots, rhizomes	0.00	0.00	0.00	0.03	0.03	0.35	0.55	0.03
seeds	0.00	0.00	0.02	0.04	0.20	0.14	0.00	0.59
	bark	buds	flowers	fructus	leaves	roots	roots, rhizomes	seeds

Figure S1.7(c): Median confusion matrix for Hybrid Bayesian Network trained on autoencoded feature space.

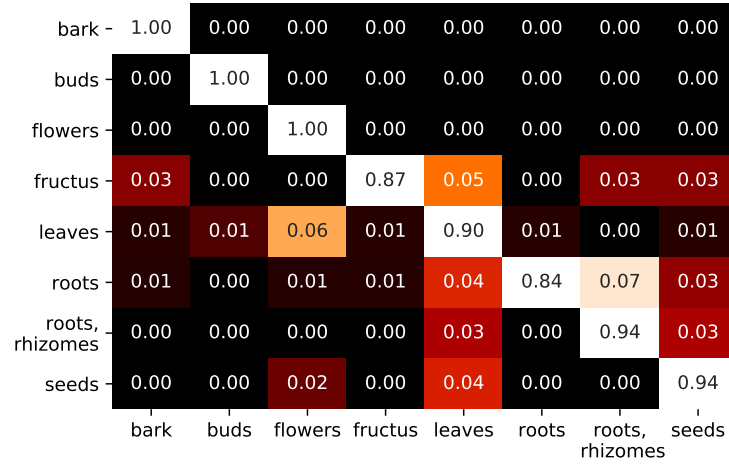


Figure S1.7(d): Median confusion matrix for classifier based on sparse non-negative Tucker decomposition.



Figure S1.7(e): Median confusion matrix for classifier based on sparse non-negative matrix factorization.

## 8 Github repository structure.

Directory	Description
data	Directory with data (csv, sif, npz)
models	Directory for storing computed models
notebook	Computational experiments
results	Directory for storing computed results
src	Python sources

Table S1.8(a): Structure of chemfin-plasp repository (github, link: <https://github.com/kharyuk/chemfin-plasp>)