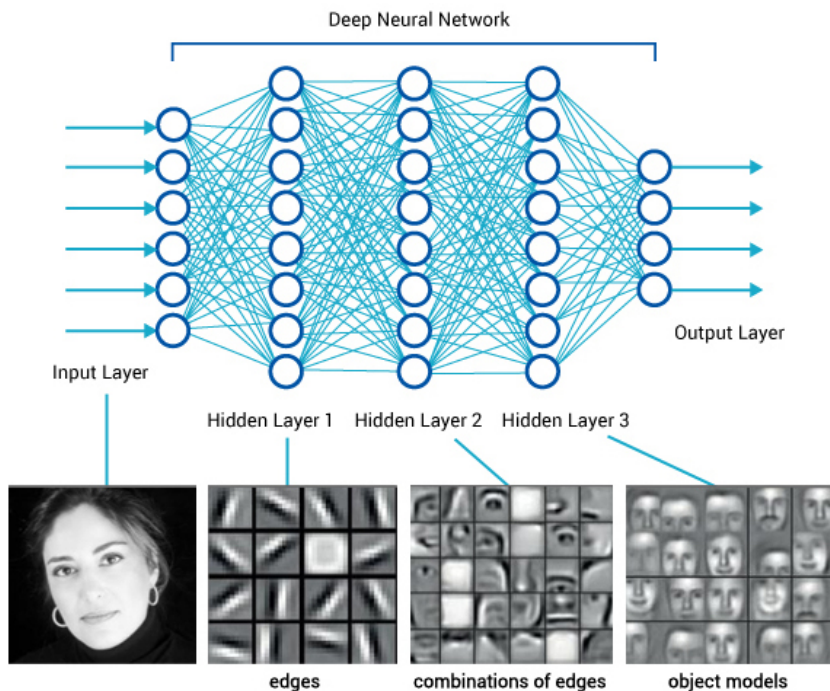


## Lesson 35: Image Embedding

Every data set so far came in the matrix (tabular) form: objects (say, tissue samples, students, flowers) were described by row vectors representing a number of features. Not all the data is like this; think about collections of text articles, nucleotide sequences, voice recordings or images. It would be great if we could represent them in the same matrix format we have used so far. We would turn collections of, say, images, into matrices and explore them with the familiar prediction or clustering techniques.

This depiction of deep learning network was borrowed from <http://www.amax.com/blog/?p=804>



Until very recently, finding useful representation of complex objects such as images was a real pain. Now, technology called deep learning is used to develop models that transform complex objects to vectors of numbers. Consider images. When we, humans, see an image, our neural networks go from pixels, to spots, to patches, and to some higher order representations like squares, triangles, frames, all the way to representation of complex objects. Artificial neural networks used for deep learning emulate these through layers of computational units (essentially,

logistic regression models and some other stuff we will ignore here). If we put an image to an input of such a network and collect the outputs from the higher levels, we get vectors containing an abstraction of the image. This is called embedding.

Deep learning requires a lot of data (thousands, possibly millions of data instances) and processing power to prepare the network. We will use one which is already prepared. Even so, embedding takes time, so Orange doesn't do it locally but uses a server invoked through the ImageNet Embedding widget.

For a start, we will use the image set of domestic animals that is available at <http://file.biolab.si/images/domestic-animals.zip>. Use Import Images and select a folder of the image files to load all the images from the folder.

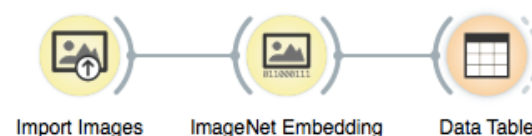


Image embedding describes the images with a set of 2048 features appended to the table with meta features of images.

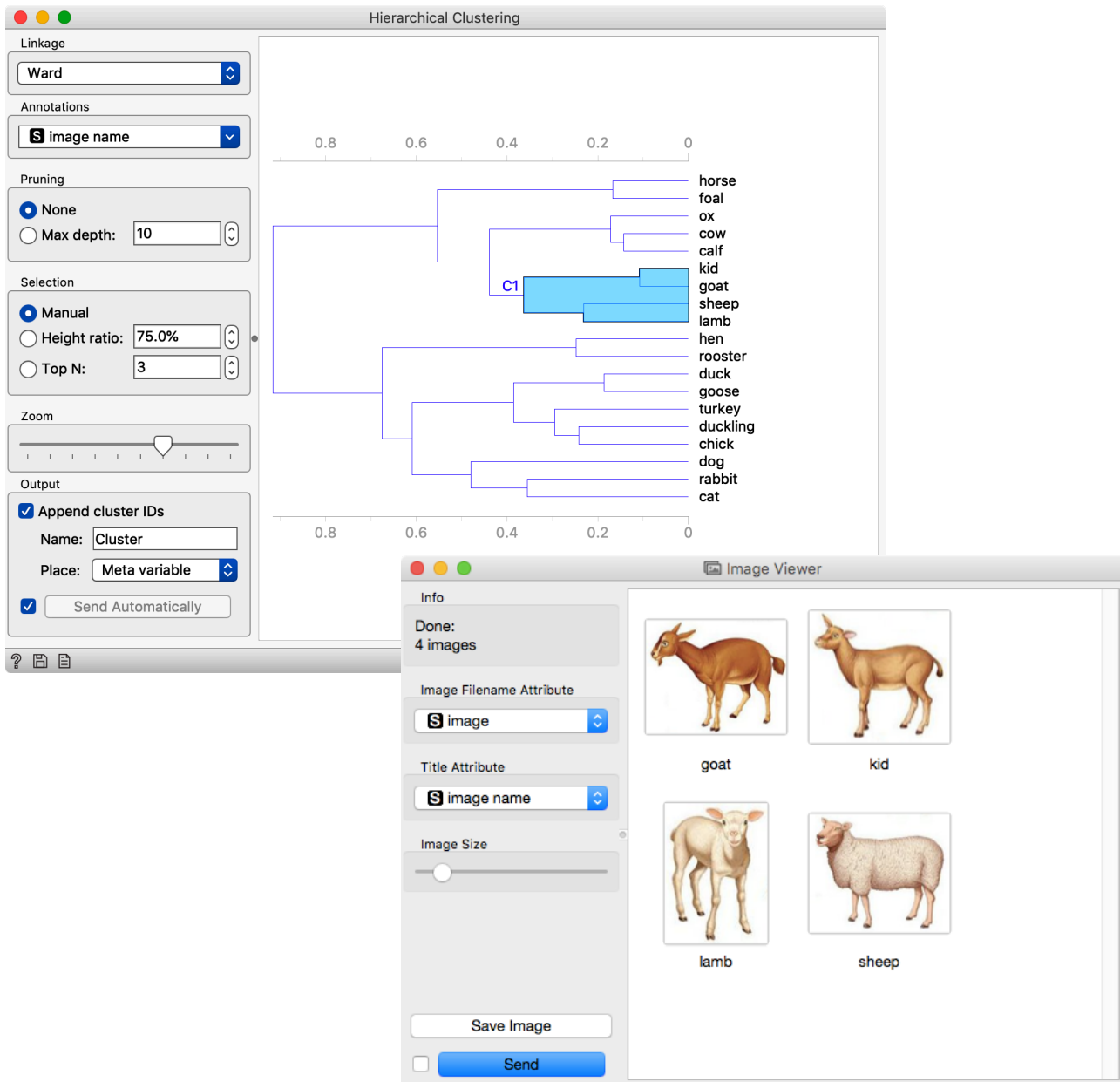
	image name	image image	size	width	height	n0	n1	n2	n3	n4	n5	n6
1	calf	/Users/bla...	45538	191	152	0.181	0.212	0.041	0.016	0.180	0.071	0.22
2	cat	/Users/bla...	22193	105	137	0.055	0.156	0.649	0.000	0.156	0.136	0.22
3	chick	/Users/bla...	14891	85	92	0.127	0.032	0.097	0.015	0.169	0.080	0.11
4	cow	/Users/bla...	62159	210	189	0.475	0.130	0.048	0.082	0.130	0.599	0.22
5	dog	/Users/bla...	28745	129	125	0.049	0.187	0.181	0.111	0.188	0.516	0.62
6	duck	/Users/bla...	39583	158	172	0.131	0.037	0.073	0.040	0.162	0.221	0.11
7	duckling	/Users/bla...	17109	99	119	0.068	0.050	0.033	0.055	0.184	0.189	0.11
8	foal	/Users/bla...	39210	147	177	0.061	0.252	0.040	0.155	0.481	0.348	0.11
9	goat	/Users/bla...	53039	221	179	0.265	0.124	0.017	0.019	0.176	0.110	0.22
10	goose	/Users/bla...	34442	141	202	0.355	0.246	0.159	0.000	0.422	0.374	0.11
11	hen	/Users/bla...	41716	134	168	0.389	0.062	0.037	0.083	0.429	0.218	0.11
12	horse	/Users/bla...	69109	285	195	0.280	0.229	0.084	0.095	0.387	0.295	0.22
13	kid	/Users/bla...	36290	170	160	0.131	0.140	0.024	0.067	0.130	0.030	0.11
14	lamb	/Users/bla...	35520	123	168	0.358	0.034	0.189	0.055	0.331	0.162	0.42
15	ox	/Users/bla...	56401	191	189	0.520	0.003	0.096	0.106	0.139	0.235	0.22

We have no idea what these features are, except that they represent some higher-abstraction concepts in the deep neural network (ok, this is not very helpful in terms of interpretation). Yet, we have just described images with vectors that we can compare and measure their similarities and distances. Distances? Right, we could do clustering. Let's cluster the images of animals and see what happens.



To recap: in the workflow about we have loaded the images from the local disk, turned them into numbers, computed the distance matrix containing distances between all pairs of images, used the distances for hierarchical clustering, and displayed the images that correspond to the selected branch of the dendrogram in the image viewer. We used cosine similarity to assess the distances (simply because of the dendrogram looked better than with the Euclidean distance).

Even the lecturer of this course was surprised at the result.  
Beautiful!



# Lesson 36: Images and Classification

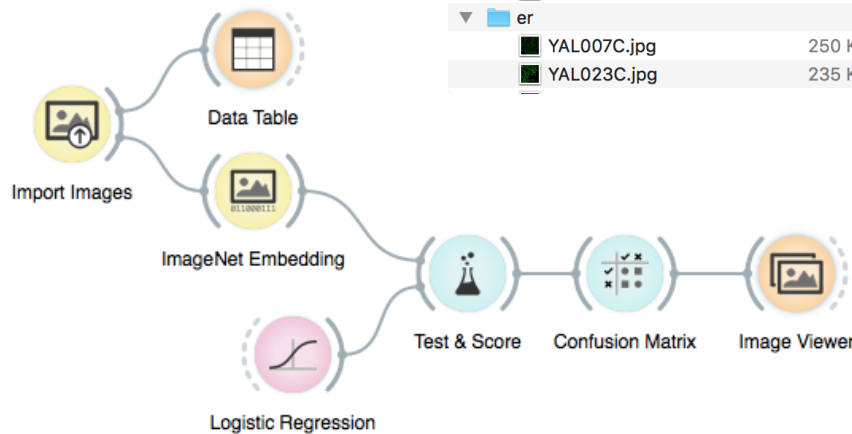
In this lesson, we are using images of yeast protein localization (<http://file.biolab.si/files/yeast-localization-small.zip>) in the classification setup. But this same data set could be explored in clustering as well. The workflow would be the same as the one from previous lesson. Try it out! Do Italian cities cluster next to American or are

We can use image data for classification. For that, we need to associate every image with the class label. The easiest way to do this is by storing images of different classes in different folders. Take, for instance, images of yeast protein localization. Screenshot of the file names shows we have stored them on the disk.

Name	Size	D
cytoplasm	--	Ti
YAL005C.jpg	163 KB	2
YAL011W.jpg	269 KB	2
YAL012W.jpg	256 KB	2
YAR019C.jpg	256 KB	2
YAR071W.jpg	276 KB	2
YBL001C.jpg	162 KB	2
YBL008W.jpg	256 KB	2
YBL016W.jpg	180 KB	2
YBL019W.jpg	41 KB	2
YBL036C.jpg	256 KB	2
YBL039C.jpg	298 KB	2
YBL051C.jpg	184 KB	2
endosome	--	Ti
YBL017C.jpg	224 KB	2
YBR097W.jpg	185 KB	2
YDR323C.jpg	184 KB	2
YDR456W.jpg	213 KB	2
YGR206W.jpg	211 KB	2
YJL053W.jpg	223 KB	2
YJR044C.jpg	233 KB	2
YLR025W.jpg	232 KB	2
er	--	Ti
YAL007C.jpg	250 KB	2
YAL023C.jpg	235 KB	2

Localization sites (cytoplasm, endosome, endoplasmic reticulum) will now become class labels for the images. We are just a step away from testing if logistic regression can classify images to their corresponding protein localization sites. The data set is small: you may use leave-one-out for evaluation in Test & Score widget instead of cross validation.

At about 0.9 the AUC score is quite high, and we can check where the mistakes are made and visualize these in an Image Viewer.



## For the End

The course on Introduction to Data Mining at Baylor College of Medicine and its installment in 2019 ends here. We covered quite some mileage, and we hope we have taught you some essential procedures that should be on the stack of every data scientists. The goal was not to turn you into one but to get you familiar with some basic techniques, tools, and concepts. Data science is a vast field, and it takes years of study and practice to master it. You may never become a data scientist, but as an expert in biomedicine, it should now be more comfortable to talk and collaborate with statisticians and computer scientists. And for those who want to go ahead with data science, well, you now know where to start.