**Eckart Modrow**

# Machine Learning with
# λSnap! – ML.Sprites

© Eckart Modrow 2020
emodrow@informatik.uni-goettingen.de

The ML.Sprite-prototypes and this script can be loaded from

https://emu-online.de/MLSprites.zip   respectively
https://emu-online.de/MachineLearningWithSnapMLSprites.pdf

If this book is helpful for you and you would like to express your appreciation in form of a donation, you can do so at the following PayPal account:

emodrow@emu-online.de
Intended use: ML-Book

# Preface

This script describes the *ML.Sprite-Library* with *Snap!* blocks, which is intended for (relatively) fast processing of large amounts of data. "Large" data volumes are almost never used in schools and initial university education - because they were hardly freely available some time ago, and money is scarce in education. In the meantime, however, there are large amounts of data in abundance, be it as a data collection on the Internet or as image files, because they are also "large". Education thus has the chance to deal with relevant data and thus find numerous points of contact with the field of "computer science and society". In the long run they are more important than any programming tricks in terms of general education.

Especially for beginners it is important to "see" what they are doing with their programming attempts. *Snap!*'s fantastic visualization capabilities are complemented by the *ML.Sprite-Library*, which includes library functions for graphics and images that, like the *Snap!* tables, quickly display the results of operations. Speed is important in this area because it supports experimental work in trial and error style. If you must wait too long, you won't try that much. The *ML.Sprite-Library* supports this approach by implementing most time-critical functions in JavaScript. Besides, these blocks also show how text-based programming can be senseful integrated into a graphical development environment.

The *ML.Sprite-Library* contains blocks from the area of data visualization and table handling, which is supported by the introduction of the data type *table*. In addition, functions of linear algebra with the data types *vector* and *matrix*, the solution of linear systems of equations and interpolation by polynomials are available. SQL queries are integrated, neural networks from perceptrons can be easily created and trained, image operations can be quickly executed via kernels as well as through vector and matrix operations. The examples show how this can be done. But they always show only one way - invent others and better ones for yourself!

Unlike the first version, the *ML.Sprite-Library* is divided into six prototypes, which can be loaded as attached parts of the overall sprite called "*Arthur&Ina*" or as individual parts. This limits the number of visible blocks. You can limit yourself to the prototypes that are currently needed. Each prototype contains a small example script that illustrates its use.

This book is a translation from German. Unfortunately, I do not speak English well, so it will be bumpy. I apologize for that. Be strong and hold it! Many thanks for the wonderful help of the *DeepL*[1] translation program. I would probably never have finished without these.

I would like to thank Jens Mönig and Rick Hessman very much for their support and the numerous discussions.

I wish you a lot of fun working with *Snap!* and the prototypes of the *ML.Sprite-Library* from Arthur and Ina!

Goettingen, 13 January 2020

---

# Content

# 1    Artificial Intelligence and School

The term "artificial intelligence" is currently and for the foreseeable future more than current. In Germany, the Year of Science 2019 has been declared the "Year of Artificial Intelligence". In the field of "digitalization", the term is shaping discussions in the media, business and politics. Informatic didactic contributions are also increasingly being made on the subject.

In school informatics the topic is not really new. For three to four decades now, there have been examples of neural networks (NNs) suitable for use in schools, for example, which are developed and trained by the pupils themselves [Baumann] [Modrow1]. Such networks are clear and easy to understand, encourage students to work independently and then to discuss philosophical implications based on their professional experience [Modrow2]. Above all, however, they are small. This is exactly the difference to the current NNs: they are big. According to Ian Goodfellow [Deep Learning], one of the leading developers in this field, basically nothing has changed compared to the old small networks. The structure and methods have (almost) remained the same, but of course they have been improved. What has changed is the performance of the computers on which the NNs run and the amount of data available to train them. This, however, leaves older findings valid, such as Marvin Minsky's [Minsky] 1967 findings on the equivalence of NNs and finite automata. The result is no wonder, because the model of finite automata has its roots in the first NNs. However, such results help to classify a topic: if the term "learning" tends to bring brains to mind, the finite automaton tends to bring the field to the level of vending machines.

But we have a problem with that. In school, real applications are usually reduced to small model systems that still show some of the original properties. If, however, the property of being large in current neural networks is what makes the difference to earlier versions, the restriction to small networks is at least questionable. Such a thing could have been - and has been - done some decades ago. So, what is new about this topic?

The suggestions for treating large NNs in class often consist of training finished NNs using finished training data. Students then watch the net learn, slowly improving its results. You don't need a real NN for this experience, a video was enough. You can't see that the net is big and you can't see why this size is important from watching it. All you can see is that the results are improving. You don't learn anything from this experience alone from NNs. A discussion of the effects of NNs then is based on the information that they exist and that they can learn. Further technical basics are missing, so that this discussion could take place just as well in other subjects.

Let us compare the situation with an example from physics. The relatively new image of a black hole [SZ] shows that there are black holes and that they "swallow" matter. However, this information alone does not integrate the topic into the physics lesson, because a technical treatment of black holes is largely beyond the possibilities of the school. But within a subject area "gravitation", which contains numerous activities, historical and social references, typical problems of school physics, etc., the picture links school physics with "science after school", shows ways to a more profound occupation with it and, for example, encourages reflection on whether the learners see a personal perspective in this area - or not.

What do we learn from this?

The pure introduction of new technologies has no place in school - there are other channels for shows. The pure information that such technologies exist is also not enough to assign the topic to a specific subject. On the contrary, if you limit yourself to that, then it would be better to locate subjects in which, for example, the social or philosophical effects are discussed, and the topic is thus networked with other aspects. Only the didactic reduction of a question to a complexity level, on which the learners can work as independently and imaginatively as possible, makes the topic pedagogically fruitful.

> *In the field of artificial intelligence, it is not the passive observation of the learning of networks in schools that is important, but the active promotion of the understanding of human learners for the fundamentals and implications of this process.*

New for the school are the tools we can use today. The visualization possibilities on the one hand and the use of powerful libraries on the other hand make it possible for the learners to explore simple first approaches on their own and thus experience the consequences of this expansion. In this work the meaning of the terms used becomes clear and thus assessable. The term "learning" has, for example, in the field of machine learning largely the meaning of "parameter adjustment". This does not quite correspond to the common meaning - and thus its use in the media, for example. The number of parameters in e.g. small perceptron networks increases immensely with their expansion - and thus the time required for training as well as the number of training data required. Extrapolation to really large networks therefore raises the question where these resources come from. When it comes to training duration, the parallelizability of the algorithms and the speed of the computers is decisive, but when it comes to data volumes, their sources - and that is often us - are decisive. So, we can't be indifferent to the application of AI systems, it directly leads to problems with data protection and is therefore a current socio-political topic. This also applies to the data itself. Working on machine learning topics quickly teaches us that there is not much to do with the freely available data itself. It often only becomes interesting when different data sources are linked. However, this linkage is usually not made by statistical quantities, but by the individual sources themselves - i.e. us. For example, does the spread of a disease have anything to do with the behaviour of population groups? We can actually only answer this conclusively when we know whether the disease occurs more frequently in precisely these people than in others. Otherwise, we will not get very far beyond suspicions.

Working with machine learning requires learning time - and this is only available to a limited extent. The more powerful the instruction set used, the more time is left for the question of the consequences. The ML.Sprite-library is designed to free up this time.

One more remark on this: I think that besides the usual learning of technical contents and methods, especially in the field of computer science, a good deal of creativity belongs in the classroom. Computer science provides wonderful tools such as *Snap!* for this very purpose. If the school concentrates exclusively on teaching facts and data as well as practicing the application of calculations, there is the danger that the students never experience what it is like to discover and understand connections and backgrounds themselves or to find and test their own solutions for interesting problems. This would be a bit sad, because a chance to develop a creative personality, aware of its possibilities and limits, would be missed at least in this field. The goal of the ML.Sprite-library is therefore to provide a toolbox for the learners, which is suitable for their own projects in the field of machine learning. It is explicitly not the goal to provide ready-made solutions for certain problems.

# 2    Machine Learning

The term "machine learning" is often used as a synonym for "artificial intelligence" or "neural networks". However, this limitation is not true. For example, the definition found on the SAP page [SAP] is more precise:

> *Machine learning technology teaches computers to perform tasks by learning from data instead of being programmed for the tasks.*

"Learning from data" can be understood as adapting the parameters of a function. A data set (image, table, character string, ...) is presented as input vector $E$ to a machine. It calculates an output value $k$ from this, which assigns the input to a category ("It is a cat", "Feature present" (or not), "The word 'car'", ...).

$$f(E) = k$$

This assignment can take place in very different ways. For example, you can adjust the parameters of a polynomial, search for similar input values ("k-next neighbours"), work with decision trees, use Bayesian filters, ... - or even train an NN. All these methods have in common that the "machine" contains a set of parameters that can be changed. The machine "learns from data" by repeatedly reading in a data set, calculating the output value from this using the current parameter set, and then comparing this output with the "desired" output value using some method. If there is a deviation, it changes the parameters so that the output at least approaches the "desired" value. "Desired" values may be known in advance ("The image is a cat image"), may come from outside e.g. from a "trainer" ("supervised learning") or may be generated by the machine itself ("unsupervised learning"), e.g. by extracting features from many training data ("clustering"). In all cases, the machine does not "learn" anything, but adapts parameters according to a given procedure.

This approach, too, has long been widespread in schools. "Learning Nimm-games" etc. can already be found in the first computer science textbooks. What is new again is the scope of the required training data. A large NN can have billions of parameters that need to be trained - and this requires "a lot of" training data. Another new feature is that these data are available on the net. So, if applications available "for free" are paid "with data", then we now also know how and why this happens.

If you look at common textbooks on machine learning [Grus] [Albon], you won't find much about NNs, but a lot about data handling. These must be normalized, for example, in order to make the many input data, which can come from very different sources, compatible. For example, if we photograph many dogs with an older digital camera and many cats with a newer one, then an NN would very likely learn from these images that dog images are smaller than cat images.

The preparation of data now is a very manual activity. It can be done step by step, tested and then automated with simple algorithms. Testing is greatly facilitated if the structure of the data is easy to visualize, e.g. in tables or as a graph. And algorithms are simple if they have a clear structure, e.g. if, after some preparation steps, they consist of a loop in which some alternatives with the corresponding instructions are enumerated. The power of the developed scripts does not depend so much on the algorithmic structure as on the power of the available commands. Or vice versa: if you have enough powerful commands, you

can do a lot with simple programs. The parameters then can be adjusted in one of the usual ways. If the appropriate tools are available, the preparation of data is a very suitable topic for schools. The *ML.Sprite-Library* is intended as such a tool.

The instruction sets of the *ML.Sprites* contain solutions for a number of typical beginner problems, e.g. sorting, drawing a graph or displaying an image in false colors. This does not mean that there is nothing left to do - you can only tackle more complex problems. So instead of sorting a list, you can look for solutions to problems that require sorting, among other things. Part of the work will consist of assembling sequences of library functions, testing them, and then making them available as a new block that can be used by other sprites to perform other tasks. The access to resources of other objects can serve as a training ground for object-oriented programming - but it does not have to. Instead, simple procedures such as data exchange using global variables can be used.

# 3    The Structure of ML.Sprites

The structure of *ML.Sprites* is based on the idea of documented data sets consisting of two parts: the *metadata*, which describes the structure and context of the data (e.g. number format, image dimensions, recording device, recording date, ...) and the associated pure *data* segments. Metadata usually consists of dictionaries - names with assigned values (e.g. "Recording date: 24.12.2018").  Examples for this structure are FITS files [FITS], which are standard in astrophysics but are also used in the Vatican Library, or JPEG images from mobile phones. Also, here there are meta data (image size, compression degree, date of acquisition, often also GPS coordinates). Without these an image generation would not be possible. It is important that the image generation does not change the original data.

We adapt this structure by giving an *ML.Sprite* three local variables, each containing the data (*myData*), the data description (*myProperties*) and a collection basin for (error) messages of called blocks of the sprite (*myMessages*). These variables can be filled by importing data from different sources (SQL queries, text file, CVS file, JSON file, FITS file, direct assignment, ...), whereby the properties *myProperties* have to be adapted to the respective data. On the other hand, this can also be done "by hand". With the help of these properties, data can be converted into graphical representations (graph, data plot, histogram, image, ...), whereby either *myData* or another suitable table is selected as the source.

It is important that the image generation does not change the original data. If, for example, an image of Jupiter is used to determine the distances between its moons, then these must at least be visible in the image. Therefore, after adjusting some parameters, e.g. a false color image can be generated. In this image Jupiter itself will appear rather unstructured. If you want to examine the "eye" of the planet in more detail, the parameters must be selected completely differently, so that the moons are hardly visible. All these changes must be done in the pixels of the current *Snap!* sprite costume without affecting the image data itself.

Because tables can be displayed very nicely in *Snap!*, this display format is not additionally implemented. Therefore, the data type *table* is implemented with many of the operations commonly used in *data science* (table operations, correlation calculation, affine transformations, solving linear systems of equations, ...), which can handle larger amounts of data sufficiently quickly.

Since the library (currently) contains 118 new blocks, these have been grouped according to their functionality and distributed into a total of six sprites that can serve as prototypes for various tasks: a *DataSprite* for handling the actual data, an *ImageSprite* for image processing, a *PlotSprite* for graphical representations, a *NeuralNetSprite* for Perceptron nets, a *SQLSprite* for database queries and a *MathSprite* for linear algebra operations. In order to use them independently, some blocks were created in slightly different variants. A "pool" called *Arthur&Ina* contains these prototypes as parts. If you want to use them individually, they can be solved by *Arthur&Ina* and stored individually.
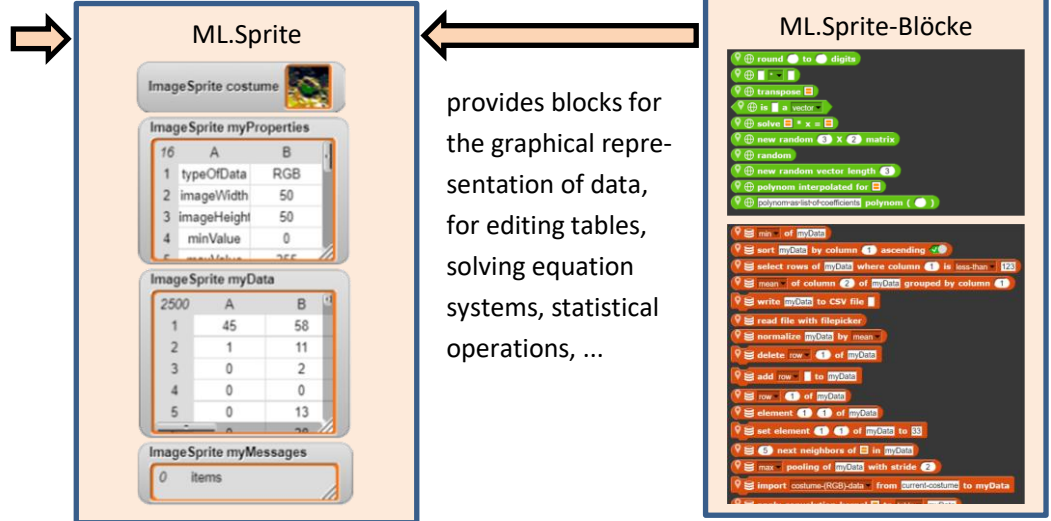
The overall structure is as follows:

imports data from
…

- image files
- text files
- SQL-queries
- JSON files
- CSV files
- …



provides blocks for the graphical representation of data, for editing tables, solving equation systems, statistical operations, ...
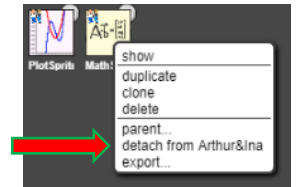
Most blocks get their parameters (image size, value ranges, colors, ...) from the dictionary *myProperties*. The properties preset with **set properties** make it possible to use blocks for creating graphics, diagrams, ... without too many parameters. If the values do not match, the properties are changed either individually (with **set property**) or in groups (e.g. with **set line attributes**).

The blocks of each sprite have a different symbol on the front to quickly distinguish them from each other and from the standard *Snap!* blocks. As soon as you can create your own new palettes in *Snap!*, the library blocks should be placed there. If this would already be done in an own *Snap!* variant, then you would either must adapt it continuously or decouple from the *Snap!* development. Both would be more unpleasant than the solution chosen now.

# 4    Working with ML.Sprite-Prototypes

## 4.1    Versions of ML.Sprites

As experienced data scientists, *Arthur&Ina* have the full range of methods at their disposal. The corresponding prototypes buzz around their heads, and if you double click on them, you will reach the corresponding sprite.  As we are professionally differently oriented than they are, we will hardly ever need all versions of the *ML.Sprites* at the same time - but which one we need depends on the chosen problem. Therefore, we can also solve all or single sprites of *Arthur&Ina* (right click on the sprite and select "*detach*") and save them individually ("*export*"). If we are lucky, someone has already done this, and we only import the needed prototypes (either with "*import*..." from the file menu or by "pulling" the file into the *Snap!* window). This is much faster than loading everything at once. After that we should set the stage to a suitable size (depending on the screen resolution) with "*Stage size...*" in the tool menu of *Snap!*, e.g. 800x600 pixels.

In each of the loaded sprites we find a small example that shows how they can be used. The new blocks are distributed on the usual palettes, mostly on *Looks*, *Operators* or *Variables*. They are located down there and have a common symbol to distinguish them from the standard blocks.

If we want, we can...

- ... work directly in the prototype, i.e. use the new blocks without further formalities.
  **Example**: The current costume, which was inserted from outside, for example, is read into the *myData* area and the corresponding properties are set. It is then printed in false-color and logarithmic representation, using the maximum value just determined during import.

- ... right click on the sprite to create a copy or a permanent clone of the prototype and work with it in the same way.

- … create new temporary clones from the Control palette with the **new clone of ...** block and use them in the following.
  **Example**: First a temporary clone of the *ImageSprite* is created. This is asked to first read astronomical FITS data, whereby the file to be read is chosen by the user. Then the image is displayed as shown in the other example.

## 4.2    Access to ML.Sprites

If you don't find the functionality you are looking for in the current *ML.Sprite,* the question arises how to access the data and/or methods of another *ML.Sprite* "from outside". Since these variables are all local, you cannot see them directly "from outside". On the one hand, this has the advantage of keeping the number of currently visible data and methods clear, but on the other hand it makes access more difficult. We will go through the different access options in sequence. We have to distinguish clearly between local and global data and methods. As an example, we choose the situation where a "normal" *Snap!* sprite wants to use the possibilities of two *Arthur&Ina* sprites, a *PlotSprite* and a *DataSprite.*

| global data and blocks |
| --- |
| global variables, e.g. **newSprite, data** |
| global blocks, e.g. **move \<n> steps** **x position** |

| *Snap!*-Sprite | *PlotSprite* | *DataSprite* |
| --- | --- | --- |
| local variables, e.g. **test** | local variables, e.g. **myData** | local variables, e.g. **myData** |
| local blocks, e.g. **commandA** **reporterB** | local blocks, e.g. **set properties** **copy of costume \<…>** | local blocks, e.g. **set properties** **copy of \<data>** |

**1. Access to the data of another sprite**

Example:     The *Snap!* sprite needs the data of the *DataSprite.*

Solution 1: In the *DataSprite* you assign the value of the data (**myData**) to a global variable (**data**). This can be accessed directly everywhere, even in the *Snap!* sprite. (But this is not very elegant!)



Solution 2: You can access directly the data of the *DataSprite* by using the **of**-Block. First select the desired sprite in the right input field, here: the *DataSprite.* Afterwards, by clicking on the left input field the local variables (here: **myData**, ...) and the local methods (here: set property, ...) are listed next to some standard attributes like position, size, ... (see next page). From these you can select the desired one.

In our case it is the variable **myData**.

`myData ▼ of DataSprite ▼`

You can work with this as usual, e.g. by assigning the first element to another variable.

`set test ▼ to item 1 ▼ of myData ▼ of DataSprite ▼`

**Example**:    The *Snap!* sprite needs the data of a clone of a *DataSprite.*

Solution:    You bind the clone to a local or global variable. With this variable you can access the clone later.

`set newSprite ▼ to ⚲ ≋ new clone of DataSprite`

In the following, we proceed as in the previous example, because only the properties of "named" sprites can be listed. Finally, replace the name of the selected sprite (here: *DataSprite*) with the variable that points to the clone.

`set test ▼ to item 1 ▼ of myData ▼ of newSprite`

For methods we choose the following convention: global methods as well as methods of other sprites <u>without parameters</u> are called by *tell* and *ask* blocks, methods of other sprites <u>with parameters</u> are called by *run* and *call* blocks. (But sometimes we do it differently. 😉)

1. **Execution of a global method (command) by another sprite**

**Example**: The *Snap!* sprite tells the *PlotSprite* to move a little.

Solution: The *Snap!* sprite asks the *PlotSprite* to move by dragging the global block or blocks into the script area of a *tell* block. On the left side of the *tell* block the addressed sprite (here: the *PlotSprite*) is selected. The blocks are executed in the context of the other sprite, e.g. with its current position and direction.



If the addressed sprite is a clone, proceed as above: replace the name of the plot sprite with the variable.



2. **Execution of a global method (command) with parameters by another sprite**

**Example**: The Snap! sprite tells the *PlotSprite* to move differently.

Solution: The *tell* block is expanded to the right by as many fields (small right arrow) as there are open parameters. The corresponding input fields for the methods must be completely empty!



Note: There are other ways to do this and, above all, much more differentiated. Please read the Snap! manual.

3. **Call of a global method (reporter) by another sprite**

**Example**: The *Snap!* sprite asks for the properties of the *PlotSprite*.

Solution: The *ask* block is used instead of the *tell* block. Otherwise as described above.

4. **Calling a local method by another sprite**

**Example**:    The *Snap!* sprite changes the size and color of the *PlotSprite*.

Solution:    The **run** block is used with a local method that is selected using the of block.



**Example**:    An *ImageSprite* is asked to determine the total brightness around the specified point (100|50) in a radius of 5 pixels.

Solution:



**Example**:    A *MathSprite* is asked to provide a 4x3 matrix with random numbers.

 Solution:



5. **Calling the code of a local method by another sprite**

**Example**:    The *PlotSprite* wants to execute the method of the *ImageSprite* to draw a circle on its own costume.

Solution:    If local methods do not depend on local variables and/or other local methods, you can export the code and execute it in another context. This is typically the case with JavaScript functions. The drawing operations of the *ImageSprite*, which are sometimes also needed by other prototypes, serve as examples here. Since we cannot store global methods for the chosen way of binding functionality to sprites, the drawing methods of the *ImageSprite* are available in a second version, which has been marked as "*exportable*". Since the properties of the *ImageSprite* can no longer be accessed here, the number of required parameters "explodes". But for this you can export the code. 😉

In *PlotSprite* the code of the drawing function is accessible via the **of**-Block.

This code is executed in the context of the *PlotSprite*, whereby the required parameters must all be listed. Afterwards the changed costume is displayed again.

If you have to combine several such calls, it becomes a bit cumbersome. It is therefore recommended to write a local method within the addressed sprite that triggers the required actions. This method is then called from outside.

**Example**:    The specifications for a diagram are set in *PlotSprite*. In another sprite you can simply call this method.

## 4.3    Importing and exporting data

*Snap!* can import a range of data formats directly. This can be done by "*dropping*" the corresponding files on the *Snap!* window or by right-clicking on a variable watcher[2]. Both works well with text, CSV and JSON files. Other text file formats such as FITS can also be imported in this way, where you are asked if you are serious. Exporting works the same way. If you want to do the same programmatically, use the reporter block **read file with filepicker**. A file manager window appears, where you select the file as usual. Afterwards the data will be imported.

The main task is then to assign this data to the **myData** variable and set the corresponding properties in **myProperties**. This is done by the following block, which imports data from outside into the **myData** area. This can be image data, table data or the data of the current costume. The costume is stored as a table of RGB values.

**Example**: The *ImageSprite* saves an image (source: [NASA]) and displays it with false colors.

**Example**: Almost 600.000 data records from a CSV file are read in about 10 seconds. The properties are set.

**Example**: SQL-import

If we have access to an SQL server, we can also read in data from there. In our case we use an *SQLSprite* to import the results of a query into the variable **myData**. The data is converted into a table and its relevant properties such as number of columns and rows, ... are reset.

---

[2] You get a variable-watcher, if you set a checkmark in the box beside the variable.

**Example**: JSON-import

Again, the easiest way is to simply "*drop*" a JSON file into the *Snap!* window. But it also works automatically. First of all, we look for interesting JSON data and of course choose the statistics of baby names in New York City - what else. The appropriate block for this is again **import *<table data>* from *<read file with filepicker>* to myData**. The result is a list with two columns and two rows, the metadata and the actual data. Because we are interested in these, we replace the original data with the element (2|2) of the table. Of course, we looked at the individual elements in table form beforehand to check what we loaded. Of the many columns, we copy the three interesting ones into a new table, add column headings and import the result back into **myData**.

The result: 19419 baby names

Who would have thought it!

**Example**: Importing data with the mouse

In many cases it is advantageous, especially with images, to read in data with the mouse. The *Sensing* palette of the *ImageSprite* and partly also of the *PlotSpite* offers blocks such as **<...> by mouse**, with which image values, image coordinates, coordinates in the used coordinate system for graphs and/or data points, the data on a section through the image, starting and end point of a line, center and radius of a circle and the summed brightness values together with their number in a circle can be determined. As an example, the height of ancient columns is to be measured. Therefore, the costume image of the *ImageSprite* with the columns is imported and then measured with the mouse (yellow line).

**Example**:  Measuring distances on an image

**Example**: Measure the total brightness around a pixel in an image (Source: [HOU])

The export of data can be done directly from a Variable-Watcher.

For scripts there are two new blocks **write <table> to CSV file <filename>** and **write string <string> to file <filename>**. The results will end up in the download folder of the browser, as usual in *Snap!.* The two blocks allow you to automate data exchange with spreadsheet programs or text files, for example, to save the results of data processing.

## 4.4 The DataSprite

All *DataSprites* contain a block in the *Looks* palette that creates a costume of the specified size with the RGB values of the background color and switches to it. They also all have the same set of local variables: **myData** contains the actual data, **myProperties** its metadata, and **myMessages** any messages generated when the blocks are executed - usually error messages. If something does not work, you should look there.

The *DataSprite* is used to manipulate table data. The properties are adapted accordingly: they describe the current state of the stored table.

The block **import <data> from <source> to myData**, which is often used together with **read file with filepicker**, is used to read data. Existing data in table form can be written to CSV files with **write <table> to CSV file <filename>** for further processing. If they are available as a character string, they can be saved as a text file with **write string <string> to file <filename>**.

In the *Operators* palette we find three more blocks: **random** provides random numbers from the usual range between 0 and 1 with full accuracy, **regression line parameters of <source>** calculates the parameters of a regression line through the given data set. The predicate **is <source> a <type>** tests the input for whether it is a table, a vector or a matrix, the latter being allowed to contain only numerical values. The block is mainly used to intercept errors.

The main functions of the *DataSprite* can be found in the *Variables* palette. There are two blocks for generating test data: **<n> random points with ranges <xmin><xmax> and <ymin><ymax>** generates random points from the given range that spread around a straight line, **<n> random points near <term> between <xmin> and <xmax>** corresponding to points distributed around the given function graph. The block **new <n> X <m> table** creates a new empty table of the specified size and **copy of <source>** returns a copy of the passed data. This is sometimes necessary to prevent the operations on the data from altering the original data itself.

The next group of blocks refers to the elements of the specified table: **delete <row/column> <n> of <table>** deletes a column or row of the specified table, **add <row/column> <data> to <table>** appends the passed data as a new row or column to the table. With **<row/column> <n> of <source>** you can copy individual rows or columns of a table. The last two blocks allow access to individual table elements.

**select rows of** <table> **where column** <n> **is** <predicate> <value> allows to select table rows with certain properties.

The block **<property> of <data>** determines the minimum or maximum of a vector and their positions, the number of elements or their sum, the mean value or median as well as the variance or standard deviation.

Relationships between two columns of a table can be determined with **<property> of column <n> and <m> of <table>**. **ranges** determines the value ranges that are needed for graphical representations, for example, the **covariance** and the **correlation coefficient** are used for statistical investigations. The next block groups the data of a column and calculates the specified values of the respective groups.

The block **normalize <data> by <option>** is intended for graphical representations and the comparison of data. Also, for neural networks it is sometimes needed with the option **softmax**.

The last four blocks are mainly needed for examples from the field of machine learning. The reporter block *<k> next neighbors of <point> in <data>* determines the **k** nearest neighbors of a point in a point set. It can be used, for example, for clustering problems. ***sort <data> by column <n> <ascending/descending>*** allows you to sort a table by a column to be specified. The ***pooling*** block reduces the amount of data with a step size to be specified, e.g. for convolutional neural networks, and the actual convolution can be performed for images or neural networks with ***apply convolution kernel <kernel> to <table/image>***. Examples of this are given in the next section.

**Example**: Income data from the US Census income dataset (Source: [Census])

We want to dig a little bit in data and therefore download the *Census Income Dataset*[3] from the net.  The corresponding CSV file can be loaded from the storage directory into **myData** and displayed immediately. It contains 32562 datasets. A right click on it and selecting "*open in dialog...*" shows all columns.

What connections could now become apparent in?

Our *DataSprite* blocks do not help much at first, because they mostly process numerical data. If we want to use them, we have to scale the columns so that numerical contents result. In the simplest case, we just replace texts with numerical values - and we should think carefully about the consequences this might have in terms of interpretation.

Let's start with the last column: The income values are given for only two ranges: less than or greater than $50000. We assign the values *1* and *2* to these ranges. (Or 0 and 1, or -1 and +1, or 0 and 100, or ... Would these changes have consequences)? In order not to change the original values, we create a variable **income** and store the changed values there by copying column 15 into this variable, deleting the first value (the heading) and then using the ***map...over...*** block to change the contents. At least that's what we're trying to do. Since there are quite a lot of values, we right-click on the map block and choose *just-in-time compilation* with "*compile...*". A small lightning symbol appears in front of map. Unfortunately we only get the unchanged column 13 when we look at the result as a table again.

What's up? We look at the first element of **income** and check if it is a string. That is the case, but it is longer than we thought:

---

So, we have to throw out the leading spaces first. That works: our variable *income* only contains the values *1* and *2*, as we can quickly check by looking at it.

What does this income now depend on?

Maybe from age? We combine column 1 (age) and our modified income column into a new table called **testdata**.

The correlation between age and income is described by the *correlation coefficient*. The calculation is simple:

And what does that mean?

**Tasks**:

1. Find out about the meaning of the correlation coefficient and the interpretation of the value obtained. What does the value "0.2340..." mean?
2. Does the correlation coefficient in this case depend on the type of numerical scaling of the data (1 and 2, -1 and 1, ...)? Check that.
3. Determine other correlation coefficients, e.g. between education and income, country of origin and income, marital status and income, country of origin and occupation, ...
4. Find out if and when scaling of non-numerical data can have an influence on the result.

**Example**: New York Citibike Tripdata (Source: [NYcitibike])

Let's go see who actually rides a bike in New York. To do this, we download the NYCitibike rental data of one month onto our computer, that is the already mentioned almost 600000 data records. We take a closer look at them.

| 577704 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | tripduration | starttime | stoptime | start station | start station | start station | start station | end station | end station | end station | end station | bikeid | usertype | birth year | gender |
| 2 | 695 | 2013-06-01 ( | 2013-06-01 ( | 444 | Broadway & | 40.7423543 | -73.9891507 | 434 | 9 Ave & W 1 | 40.7431744 | -74.0036644 | 19678 | Subscriber | 1983 | 1 |
| 3 | 693 | 2013-06-01 ( | 2013-06-01 ( | 444 | Broadway & | 40.7423543 | -73.9891507 | 434 | 9 Ave & W 1 | 40.7431744 | -74.0036644 | 16649 | Subscriber | 1984 | 1 |
| 4 | 2059 | 2013-06-01 ( | 2013-06-01 ( | 406 | Hicks St & M | 40.6951284 | -73.9959506 | 406 | Hicks St & M | 40.6951284 | -73.9959506 | 19599 | Customer | NULL | 0 |
| 5 | 123 | 2013-06-01 ( | 2013-06-01 ( | 475 | E 15 St & Irv | 40.7352427 | -73.9875856 | 262 | Washington | 40.6917823 | -73.9737299 | 16352 | Subscriber | 1960 | 1 |
| 6 | 1521 | 2013-06-01 ( | 2013-06-01 ( | 2008 | Little West S | 40.7056925 | -74.0167768 | 310 | State St & S | 40.6892694 | -73.9891286 | 15567 | Subscriber | 1983 | 1 |
| 7 | 2028 | 2013-06-01 ( | 2013-06-01 ( | 485 | W 37 St & 5 | 40.7503800 | -73.9833898 | 406 | Hicks St & M | 40.6951284 | -73.9959506 | 18445 | Customer | NULL | 0 |
| 8 | 2057 | 2013-06-01 ( | 2013-06-01 ( | 285 | Broadway & | 40.73454567 | -73.9907414 | 532 | S 5 Pl & S 5 | 40.710451 | -73.960876 | 15693 | Subscriber | 1991 | 1 |

Of course, we still have to find out from the source what the data actually means - i.e. look at the metadata. For *gender* we learn that *0: unknown*, *1: male* and *2: female*. For the columns "*tripduration*" and "*gender*" we get some data:

The average loan period, based on gender:

| 4 | A | 2 |
|---|---|---|
| 1 | value | mean |
| 2 | 0 | 1753.298818681775 |
| 3 | 1 | 1063.548722541860 |
| 4 | 2 | 1233.249445298994 |

We figured as much!

For further calculations we delete the header line of the table …

… and see if they're any lazier on Broadway:

result 1380.553088413

I see. It's probably worse in Central Park!

result 2230.303350254

All right. All the prejudices don't have to be true. 😉

**Tasks**:

1. Maybe only women in Central Park ride bikes more. Check it out.
2. It's not like there's only one bike rental place in Central Park. Find out the appropriate mean values for the entire area.
3. Is there actually borrowing information for other parts of the city? Search and compare the results with Manhattan.
4. Determine the average borrowing times per weekday, overall and for individual stations. Are there any differences? What are the differences?
5. Above, the average borrowing duration was calculated based on gender. It could be done the other way around. Would that be complete nonsense or are there any questions where this would make sense?

## 4.5    The PlotSprite

The *PlotSprite* is used to create and display diagrams. Therefore, it is mostly used together with a *DataSprite*. In addition to the three blocks for managing its properties in the *Variables* palette, which primarily contain presets for generating the diagrams, there is only one new block *<costume-/graph-coordinates> by mouse* in the *Sensing* palette for the "*remeasurement*" of values with the mouse and two new blocks in the *Operators* palette *convert <value> to <...>* for the conversion of screen to graph co-ordinates and vice versa as well as the already known predicate block for type checking *is <data> a <vector/matrix/table>?*.

| 33 | A | B |
|---|---|---|
| 1 | typeOfData | empty |
| 2 | backColorRe | 255 |
| 3 | backColorGr | 245 |
| 4 | backColorBl | 255 |
| 5 | leftOffset | 60 |
| 6 | upperOffset | 0 |
| 7 | lowerOffset | 20 |
| 8 | title | |
| 9 | titleHeight | 20 |
| 10 | xLabel | |
| 11 | xLabelHeigh | 15 |
| 12 | yLabel | |
| 13 | yLabelHeigh | 15 |
| 14 | xLeft | -10 |
| 15 | xRight | 10 |
| 16 | yLower | -10 |
| 17 | yUpper | 10 |
| 18 | lineStyle | continuous |
| 19 | lineWidth | 1 |
| 20 | lineColorRec | 0 |
| 21 | lineColorGre | 0 |
| 22 | lineColorBlu | 0 |
| 23 | datapointSty | square |
| 24 | datapointWic | 5 |
| 25 | datapointCo | false |
| 26 | datapointCol | 255 |
| 27 | datapointCol | 0 |
| 28 | datapointCol | 0 |
| 29 | scalesPrecis | 3 |
| 30 | scalesXtexth | 12 |
| 31 | scalesYtexth | 12 |
| 32 | scalesNumb | 10 |
| 33 | scalesNumb | 10 |

The actual new tools can be found at the bottom of the *Looks* palette. On the one hand, you can use them to change the preset properties and thus adapt them to the current problem, e.g. the diagram labels or the value ranges. On the other hand, you can create different types of diagrams.

*set labels ...* set the values for the diagram title and the labels of the axes. Since the distances of the diagram axes from the edges depend on them, these values are also determined by *set offsets from edges*. While we're at it, we can also reset the value ranges with *set ranges ...* and the display of the numbers on the axes with *set scale attributes ...*. *set pretty ranges* automatically sets the number ranges with "*pretty*" limits. *add axes and scales* then draws the frame of the new diagram.

Since we often need several diagrams, we first create a clone of the *PlotSprite* and pass the necessary data to it. This clone will display the tastefully designed new diagram on its costume.

Three of the blocks in the *Looks* palette are for programming convenience: **properties of group <groupname>** combines several properties and thus facilitates the data transfer to *JavaScript functions*. **ranges of <data>** determines the value ranges of a two-dimensional table and **copy of costume <costume>** is needed if you want to switch quickly between two versions of a costume.

The diagrams themselves are generated with the remaining new blocks. The properties of line diagrams are set with the block **set line attributes ...** Function graphs are drawn with these attributes using the **add graph ...** block, for example. The function term can be passed either as a list of the coefficients of a polynomial or as a "*ringified*" *Snap!* term.

**Example**: Drawing a function and its derivatives in different colors and line styles.

**Tasks**:

1. Plot different types of functions (trigonometric functions, logarithms, polynomials, ...) as a graph on a *PlotSprite.*
2. Extend the function graphs with their derivatives.
3. Choose different value ranges, accuracies of number representation, text sizes and labels for the representations.

If we want to display the contents of a data table graph-
ically, this can be done with the block **add dataplot with
numeric scales ...** Scales and labels are again added with
the block **add axes and scales**. The way of displaying the
data points can be adjusted very precisely with the block
**set datapoint attributes ...**. The lines in between keep the
line attributes.

**Example**: Representation of a point set

We ask a *DataSprite* to provide 100 random points scat-
tered around a straight line with the slope m=0.5 and the
intercept b=0. The obtained points are displayed in a dia-
gram.

**Example**: Additionally the regression line is now drawn.

**Example**: Display of mixed data

Text data is often combined with numerical data. An example would be the turnover data of different representatives in one year in one area. If we want to display them graphically, the x-axis, for example, must be labelled with text data, while the y-axis is treated as before. To create the diagram we use the block **add dataplot with text and numerical scale**.



As the last of the new blocks of *PlotSprite* we consider **add histogram of <data> with <n> groups**. Histograms can be generated and displayed directly from data sources.



**Example**: An RGB image is loaded, decomposed into grayscale, and the normalized distribution of image values is displayed as a histogram on a new *PlotSprite*. We find the actual image as a costume of an additional sprite called "*thePicture*".

First of all we load the image into the data area of a *DataSprite*:



We get 172800 RBG values.

We convert these into grey values. We use the compiled version of the **map** block.

Then we switch to the *PlotSprite* and copy the loaded data of the *DataSprite*.

Now we can display the data as histogram e.g. on a new *PlotSprite*.

**Tasks**:

1. Search for different amounts of data in the network. Display these or parts of them graphically.

2. Automate the creation of histograms with a new block **histogram of <costume>**. Compare the histograms of typical image types. To what extent is it possible to compare images in this way, or where could difficulties arise?

3. Represent the three colors of an RGB image in the same diagram by graphs and/or histograms.

## 4.6   The ImageSprite

The *ImageSprite* is used to display and manipulate im-
ages. For this purpose, it contains in the *Variables* palette,
in addition to the obligatory blocks for managing the prop-
erties, the possibility to load images and determine the
minimum and maximum values of the image data.

In the *Looks* palette we find the already known block **add
image ...**, which allows to display data on the costume as
gray or false color image. With the help of the block **copy
of costume** you can duplicate costumes if necessary.

In addition to the usual block for creating a
new costume, available for all *ML.Sprites*,
there are two blocks for setting line prop-
erties and fill color. With these and six
blocks for drawing figure outlines and filled
figures on the costume, you can easily cre-
ate random graphics, for example.

| ImageSprite myProperties | | |
|---|---|---|
| 16 | A | B |
| 1 | typeOfData | empty |
| 2 | imageWidth | 50 |
| 3 | imageHeight | 50 |
| 4 | minValue | not set |
| 5 | maxValue | not set |
| 6 | backColorRe | 225 |
| 7 | backColorGr | 225 |
| 8 | backColorBl | 255 |
| 9 | lineStyle | continuous |
| 10 | lineWidth | 1 |
| 11 | lineColorRec | 0 |
| 12 | lineColorGre | 0 |
| 13 | lineColorBlu | 0 |
| 14 | surfaceColor | 180 |
| 15 | surfaceColor | 180 |
| 16 | surfaceColor | 180 |

**Example**: Generation of random graphics





**Tasks**:

1. Search the web for pictures of Piet Mondrian. Try to create similar random images on the *ImageSprite*.

2. Use a "*vanishing point*" to create images in which objects appear to move "*from back to front*". Try it.

In the Sensing palette of a *ImageSprite* there are some new blocks for accessing pixel values. The block *... by mouse* has already been described above for data import. With *image value ... at ...* and *set image value of ...* individual pixels can be read or changed if they are located in the data area *myData*. *RGB at ...* and *set RGB at ...* allow the same directly on the current costume. The block *brightness around ...* determines the total brightness around the specified point.

**Example**: Simulation of a planetary transit in front of the sun

We look for a nice image of the sun (source here: [Schul-Astro]) and load it as a costume of an *ImageSprite*. To make it look more like space, we enlarge the stage and color it black. If we still draw the planet, we get the image on the right.

The planet should pass in front of the sun as a black circle. If we paint such a circle, we change the actual image of the sun. From this we make a copy *newCostume* to draw on it. Our planet should move from the far left a little outside the image (*x*=-2r) to the far right (*x*=image width+2r) on the height *y*. We can also specify the radius *r* of the planet.

We can determine the current brightness of this arrangement without too many copying processes by subtracting the brightness of the pixels covered by the planet from the total brightness determined at the beginning. To do this, we initially import the image of the sun into the data area *myData* and determine the brightness around the image center in the radius "*half image width*" as well as the number of pixels involved. *brightness around* returns the summed gray values as well as its number. From these values we calculate the average brightness of the "*slightly darkened*" sun and store it together with the current position in the variable *transit data*. We package these operations in the new block *planet transit at <y> with r <r>*.

We now want to follow the planet transit "*live*" in a diagram. For this purpose, we load the *PlotSprite* and write a block for this plot to prepare the diagram parameters. This block can be called by the *ImageSprite*.

We insert the corresponding block at the beginning of our transit script and supplement it with two calls at the end of each loop to redraw the *PlotSprite* diagram with the new data.

The result corresponds roughly to one of the methods used to find exo planets.



In the *Operators* palette, we find the **is ... a ...** predicate, which is used everywhere, and two somewhat more sophisticated blocks. The first allows to perform *affine transformations* in an image by mapping three points to three others - and all other points accordingly. Additionally, the dimensions of the image must be specified.

**Example**: Let's flip an image vertically along the center line. We load the image - here: of a church - and select corresponding points on the edges. These points are combined to the two lists **source** and **target**.

After that we import the image into the data area **myData** and execute the affine transformation with the two point-lists. We save the result in the variable **newData**.

Finally, we create a clone of the *ImageSprite* and ask it to display the transformed image as a costume.

As last new block we find the **apply convolution kernel ... to myData** – block to perform convolutions on images.

**Example**: Edge detection

We load a picture of an ancient temple and import its data into the data area of **myData**. We apply the Laplace kernel $\begin{bmatrix} 1 & 1 & 1 \\ 1 & -4 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ to it and save the result in the variable **new-Data**.

Next, we create a clone of the *ImageSprite*, resize it and let it show the changed image as its costume.

Since there is sometimes a need to execute the graphics methods in the context of another sprite (see above), these methods are available in an "*exportable*" version that does not rely on local sizes. They are labeled by a left arrow ⇦ . The usage is described in chapter 4.2.

**Tasks**:

1.    Pictures are sometimes a bit "*weak*". This is because they do not use the full range of values for the three color-channels from 0 to 255.
   a: Develop a method to determine and display the value ranges of an image.
   b: Develop a method to use the full range of values, i.e. to map black pixels to 0, bright pixels to 255.
   c: Summarize the method as a new block with the costume of a sprite as input and the improved costume as output.

2. a: You can try to find "*faces*" on images by highlighting contiguous areas of a color range, e.g. "*orange*", and deleting the rest of the image. Try to develop a new block for this.
   b: You can use a kernel to isolate the edges of such areas. Find out about suitable kernels on the net and test them for the purpose mentioned.
   c: Faces are often "*oval*". Try to distinguish faces from other "*orange*" objects in this way.

3. a: Really artistic photos are black and white of course. If you don't have any, you can create grayscale images from RGB images. Do that.
   b: It's even more artistic if the photos are "hard", i.e. have a very strong contrast. Experiment a little!

## 4.7    The MathSprite

The *MathSprite* can create clones and a new costume and manage its - few - properties, but otherwise it only has new blocks in the *Operators* palette. It is intended for services for the other sprite types. Essentially, it adds linear algebra capabilities to the operators of *Snap!*.

Three of the blocks are quickly explained: ***random*** returns a random number of the usual type, i.e. between 0 and 1 with full accuracy. ***round*** rounds numbers with a fractional part after the decimal point to the number of digits specified. The predicate ***is <data> a <vector/matrix/table>*** is mostly used in scripts to catch errors and has ***been*** used in this script many times.

The next two blocks should also be self-explanatory. ***new random vector*** generates a new vector from random numbers, which is sometimes needed for testing purposes. ***new random matrix*** does the same for matrices.

The *MathSprite* works with numbers, vectors and matrices. Since vectors and matrices are sometimes needed in transposed form, we also find a block ***transpose***. This block can handle both data types.

More interesting are the next blocks.

The somewhat inconspicuous reporter block ***<a> operator <b>*** can perform the specified operation with numbers, vectors and matrices - if they are allowed. So it works between numbers, numbers and vectors, vectors and matrices, matrices and matrices, ...

Polynomials are processed by *MathSprite* as lists of their coefficients. The coefficient of the highest power is on the left. $x^2 - 2x + 1$ is written as `list 1 -2 1`. The value for a given argument is calculated by the block *<polynomial> polynomial (<argument>)*.

*solve <matrix> * x = <vector>* solves linear systems of equations - if there is a solution.

If a list of points is given, the block *polynomial interpolated for <points>* calculates a polynomial whose graph runs through these points.

Applications of these blocks follow.

**Example**: Curve through *n* points
We need three sprites: a *DataSprite* to generate the random points, a *MathSprite* to calculate the interpolation polynomial and a *PlotSprite* to plot the results. Everything will be controlled by a fourth sprite called *Control*. The points for the interpolation are to be selected by mouse.

Generation of the random data: In the *DataSprite* we write a function to generate the points. We call this function from Control.

Random data display: We write a method in *PlotSprite* to display the data. We pass the data to this method from *Control*.

Interpolation data collection with the mouse: We add an event handling method to *PlotSprite* that reacts to mouse clicks. With *<costume-coordinates> by mouse* we get the costume coordinates and store them in the variable *newPoint*. To see the points, we import the "*exportable*" *draw circle* method of the *ImageSprite*. Since it requires quite a lot of parameters and we only want to draw small red circles, we use a helper method *circle at <x><y>*, which only needs the current coordinates - the rest is already filled in. Next, we save the costume coordinates converted to graph coordinates in the list *points*. If we've clicked three times, we let the *MathSprite* calculate a polynomial through these points. We draw this polynomial a little bit thicker in red.



random data for interpolation

In *Control* we prepare everything in advance for the data acquisition.



**Tasks**:

1. a: Create "*point clouds*" which scatter around other fully rational function graphs.

   b: As in the example, define some points in these clouds through which an interpolation polynomial is to be drawn.

   c: Have these polynomials drawn.

2. a: Experiment with the number of selected points. Will the results be better if you select more points?

   b: Create "*point clouds*" which scatter around non-rational function graphs (trigonometric, ...). Can you also describe them by interpolation polynomials?

   c: Formulate a rule, when and how interpolation polynomials can be used meaningfully - and why just like that.

## 4.8    The SQLSprite

Similar to the *MathSprite*, the *SQLSprite* is intended as a service provider for projects. It encapsulates the functionality for database access and should therefore only be loaded when needed.

The few properties of the *SQLSprite* essentially contain the connection data to a sample database server and the current state of the connection (database used, current table, ...). It is possible to create clones of *SQLSprite*, but there will be few projects where this is necessary. Also new costumes are rarely needed. Most of the time the SQLSprite will just "lie around somewhere" and show by the color of its costume if the connection to the server is established.

| 12 | A | B |
|----|---|---|
| 1 | typeOfData | empty |
| 2 | connected | false |
| 3 | connection | https://snapextensions. |
| 4 | current-database | |
| 5 | current-table | |
| 6 | databases | E |
| 7 | tables | E |
| 8 | attributes | E |
| 9 | columns | 0 |
| 10 | rows | 0 |
| 11 | minValue | not set |
| 12 | maxValue | not set |

Outside the *Variables* palette, there is only one operator block to determine the part of a string - this is needed from time to time to evaluate the query results. Since no local data is required, the block is marked as "*exportable*" by the left arrow. It can therefore be executed by all sprites regardless of context. If this happens often, you should write a global method *substring* to make the call easier.

In the *Variables* palette you will find the actual SQL blocks. If you want to establish a connection, you can either use the default server or set a new connection with *set property<connection>* and the connection data. Then the block connect should establish the connection. If this works, the SQLSprite changes to the green costume. The available databases are listed with *read databases*. One of them can be selected with *choose database*.

The handling of tables is correspondingly. A specific one can also be selected and displayed.

In practice, the details of the tables in the database used are always needed. If you assign the table attributes to a



variable one after the other and let it be displayed with "*open in dialog*", you can place the corresponding table data in the *Snap!* window and start the queries.



SQL queries are composed as SELECT statements. There are two blocks for this: one for simple statements, one for (almost) complete statements.



The standard predicates and functions of SQL are used to compose the queries.

**Example**: a simple SQL query



**Example**: a more complex SQL query

The two SELECT blocks generate the text of an SQL query, but do not execute it. The reason is simple: you should be able to view the query. If you are satisfied with it, it is executed with *exec SQL-command*. In this block you can also enter other SQL commands, if you have the appropriate rights on the server. With the help of *import SQL-data ...* these are converted into tables and moved into the data area of the SQLSprite.





More complex instructions can be assembled and tested from these blocks. If the result meets the expectations, the query can be encapsulated in one block, which only contains the relevant parameters and can be used by other sprites without detailed knowledge of SQL.

**Example**: The course titles and grades for all courses of a learner, sorted by grade in descending order, are searched.





**Example**: For statistical purposes, the *schueler* table should be searchable according to different criteria.





**Example**: The 10 courses with the "*best*" results should be sorted in descending order.

## 4.9    The NeuralNetSprite

"*Deep*" neural networks are shaping the discussion on current "*artificial intelligence*". These are usually "*fully connected*" networks consisting of several perceptron layers. "*Fully connected*" means that all neurons of one layer are connected to all of the next layer. Each connection is assigned a *weight*, from which its influence on the connected perceptron is derived - but you'd better read this elsewhere.

Let's have a look at a net of three layers that receives the pixels of a current 20 M-pixel-photo as input, thus $2x10^7$ pixels. The input layer consists of $3x2x10^7$ MB numerical values between 0 and 255 (if we leave out the transparency byte). To the next layer, there are then $(6x10^7)^2$=$3.6x10^{15}$ connections - and then twice more. In total, $3x3.6x10^{15}$, i.e. about $10^{16}$ weights, would have to be determined - a completely utopian task for "*normal*" computers. So, we will have to restrict ourselves to somewhat smaller neural networks.

One way to train perceptron networks is to present them input vectors and the desired output right away. The net then calculates the output, which is the result of the existing, initially randomly selected weights, and determines the difference from the given result. Starting from the last result layer, it corrects the weights "*going backwards*" so that its output "*somewhat better*" matches the given result. The procedure is called **backpropagation**. You should also inform yourself about this elsewhere. The trained net results from many such corrections. So, "*learning*" means to adapt the parameters (the weights) by means of many examples. With the help of these parameters, the net determines an output vector from the input vector: it calculates a function value.

Our *NeuralNetSprite* can simulate and train such perceptron networks. Clones of the net can be created for this purpose - just like with the other *ML.Sprites*. For a *NeuralNetSprite* you can create costumes as usual and display the current state of the net on them. We have already described the new blocks in the *Control* and *Looks* palette.

Our weights together form a **tensor** with **m** layers consisting of **nxn** matrices. *NeuralNetsSprites* should therefore master linear algebra. To avoid having to ask the *MathSprite* every time, the *NeuralNetSprite* knows the most important operations itself. They can be found in the *Operators* palette. The only new feature is the **Softmax** function, which can be used to scale input vectors, for example. You should also inform yourself about this.

The actual *NeuralNet* blocks are therefore found in the *Variables* palette. On the one hand, the few properties of the sprite are managed there, on the other hand, the



weight tensor can be loaded as usual and saved again as a CSV file.
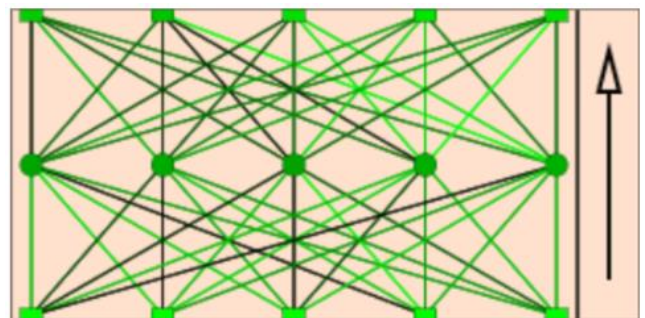


The dimensioning and initial assignment of the network is done in the block **add new weights**. With these blocks we can create and display a new neural net of any size. In this case, it has the **width** 5 and the **depth** 2. The result of the calculations with the input vector to be specified is displayed.



As the display of the many numbers would be rather confusing and also hardly informative, the connecting lines (the *edges*) are color-coded according to the values of the corresponding weights: from full green for large positive values to black for small amounts to red negative weights. As at the beginning only positive numbers are generated by random generator, new net is predominantly green.

The *nodes* of the net are color-coded like the edges. Below are the elements of the input vector as small rectangles. The inner layers form colored circles and the last layer is displayed as output layer again rectangular. The direction of the calculation from bottom to top is shown by the arrow on the far right. But since sprites can be easily rotated, the direction can be displayed differently.



Often you need the results of the last or even an inner layer of the network. These can be calculated using the **output of...** block for a given input. Since the color coding does not necessarily show the largest or smallest element clearly, it can be determined using the **of ...** block.

Now we have to train the net. In block *teach NN ...* this is done by backpropagation with a learning factor to be specified. At the beginning this factor may be a little bit bigger, but then it can be reduced.

**Example**: A neural network should learn to calculate the output vector *<0,1,0,0,0>* from the input vector *<1,2,3,4,5>*. If it did this successfully, the second output rectangle should appear light green, the rest darker.

So we create a new neural net as shown above and repeatedly apply the *teach* block. The state of the net is shown after each 20 repetitions.



And after about 200 more passes the pattern has clearly formed. The net maps the input vector to <0.1.0.0.0> reliably. You can see quite clearly that the first layer configures the inner layer in such a way that it can in turn strengthen (green) or inhibit (red) the output neurons with the correct weighting.



Of course, the net should not only recognize one pattern correctly. In training, it is given various other input vectors in random order with the desired outputs, with the weights being slightly modified so that the strongest output neuron has the correct place. Even such a small network has enough parameters for this. Examples of applications are given in the next section.

**Example**: Traffic sign recognition

We want to train a neural network (NN) to recognize 12 different traffic signs. For this purpose, we search for images of these traffic signs in the internet and reduce them to the format 100 x 100 pixels. You can display them on the screen, but the 10000 pixels are of course much too much for an NN.

To bring the amount of data within tolerable limits, we reduce the pixels to a 2x2 format by **mean-pooling**, i.e. we form the mean values of the color pixels in the four quadrants of the image. From the 30000 values of the traffic sign image we derive 12.

Since only the *DataSprite* can handle the pooling operation, we import it along with the *NeuralNetSprite* to the project with the *TrafficSignSprite*. As a fourth sprite we add a *ControlSprite* that should control everything.

For the start we give the NN-Sprite a new costume. Since we need to specify several parameters for size and color, we use the run block. Then we create the weights for a new (here: 12x2) net in NN. This net is drawn with a (still senseless) input. Next we send the NN to a well-chosen place in the upper middle and do the same with the traffic sign down there. Finally, we set some variables to 0. We need them later.

We only need the *DataSprite* for calculation. To use the **pooling** operation, the sprite must import the image data. Then it can convert it and return the result. We summarize everything in a new block **pooling of <costume>**, which we call from *Control*.

The color values of the reduced image are combined into an input vector using the new block of *Control*. The first two values of the **pooling** result specify the new image dimensions and are deleted. The vector from **get input data** thus contains the desired 12 values. These are modeled using the *NeuralNetSprite's* **Softmax** function to eliminate unfavorable input values.

Accordingly, the training vector in **get training data** can be determined with the output values of the NN you are looking for. In our case all values should be 0 except the one corresponding to the costume number of the traffic sign.

The NN gets two new methods **learn from...** and **test with...** During learning, the position of the place with the highest output value of the NN is determined and compared with the current costume number of the traffic sign. If these values do not match, learning continues.

For testing, the same operation is performed only once.

Now we have everything together to make *Control* work properly.

A teaching process consists of determining a random costume number with the corresponding costume change. Afterwards the learning process of the NN is started with new input and target data. The passes are counted.

The procedure for testing is similar: The costume is changed, and it is checked whether NN calculates the correct costume number. If this is the case, everyone is happy. Next, the percentage of correct attempts is determined.

Several learning and test runs can be easily triggered.

After about 50 training runs with a higher learning rate and another 50 with a low one for fine tuning, we reach detection rates of 100%.



**Tasks**:

1. Train a single-layer net with different learning rates and numbers of learning passes. Determine the recognition rate in percentage terms.
2. Plot the results from 1. graphically using a *PlotSprite*.
3. Experiment with multi-layer NNs. Will the results be better?
4. Increase the length of the input vector by changing the pooling. Do the results get better?
5. Increase the number of recognizable labels by allowing more than one 1 in the output.

# 5    Applications with ML.Sprites

## 5.1    Under- and Overfitting

Machine learning uses training data to adjust the parameters of a function so that other values are predicted well - if everything works out. So, you build a forecasting instrument, a kind of "telescope" for data.

One might think that the more customizable parameters a function contains, the better such a function is. But this is not so. On the one hand, (1.) more parameters also require more training data and training runs, i.e. more learning time; on the other hand, (2.) an "inappropriate" number of parameters can also prevent "good" solutions. For both we now give an example.

1: In the neural network for traffic sign recognition we have achieved very good results with one layer. If we increase the number of layers and leave the number of training runs the same, the recognition rate will deteriorate drastically.



2: If the training data are well reproduced by the function, this does not mean that this is also true for other data. It depends a lot on the type of function that is generated. As application we choose the example *polynomial interpolation*.

*The task is to use training data to adjust the coefficients of a polynomial so that OTHER data are predicted as well as possible.*

To do this, we have to generate data that can be used to calculate an interpolation polynomial. The functionalities for this are divided into three prototypes: the *PlotSprite* for plotting the graphs, the *MathSprite* for the polynomial interpolation and the *DataSprite* for generating random data that scatters around a given function.



So, we create a new project, enlarge the stage to 800x600 pixels and load the three prototypes. Then we create a variable **random data**. Now, the arrangement looks like the one on the right.

As "workhorse" we choose the *PlotSprite*. If necessary, we import required functionality from other sprites.

First of all, we ask the *DataSprite* to generate 20 random data scattered around the parabola $0{,}5 * x^2 - 3$. For this we write a function *<n> new points* in the *DataSprite*. We call this function from the *PlotSprite*.

This is enough to display the data. We write a block *show <data>*, which makes the necessary settings and does that.



We first try the interpolation with a regression line. The *DataSprite* can calculate the required parameters.



This looks actually quite nice, but on the sides it doesn't really fit.

So, we try a polynomial interpolation.

First of all, we select three random pairs from the training data, determine the interpolation polynomial from it and draw it. Because we want to experiment further, we generalize the solution to a polynomial by **n** points. We hope that everything goes well with the selection! The results depend on which points were hit (blue). Enclosed a bad and a quite good result.







Now we're getting brave! Instead of three points we vote for 5. After all, we want to do a good job! That works out great in the middle, and then – oops!



Maybe we just need to take more points. Let's try it with 10. The polynomials run through more points, but at the edges they "run away" - and in between mostly.

Well, then with all points!

You can see that with an increasing degree of the polyno-
mial there is more training data directly on the graph, but
in between only nonsensical values are "predicted" by the
wild oscillations of the polynomial.

The quality of what is learned therefore depends very
much on how we deal with deviations. We have to decide
which inaccuracies can be tolerated in detail so that the
forecast becomes reliable overall. If the degree of the po-
lynomial is too small, we speak of *underfitting*, if it is too
high, of *overfitting*.



### *Tasks*:

1. Discuss different ways to determine a "good" degree of the interpolation polynomial
   (i.e. its highest power).
2. Formulate your results so precisely that they can be realized as scripts.
3. Test the scripts on different data sets.

## 5.2    New York Citibike Tripdata [NYcitibike]

Even in New York, cycling has become "hip" and borrowing data can be loaded as CSV files. We do so and load the almost 600,000 data records from June 2013 into a table. We split the column headings to get a pure data table. Of course, this is done in a *DataSprite*. Since we also want to create graphics, we load the *PlotSprite* right away.

What did we actually find there?

The data legend provides the interpretation for the data: *Trip Duration (seconds), Start Time and Date, Stop Time and Date, Start Station Name, End Station Name, Station ID, Station Lat/Long, Bike ID, User Type (Customer = 24-hour pass or 3-day pass user; Subscriber = Annual Member), Gender (Zero=unknown; 1=male; 2=female), Year of Birth*

Since the geographical longitude and latitude of the rental stations are given, it is a good idea to use the *Word Map Library* from *Snap!*. We write a small block, which shows the surroundings of a rental station as a map.

Let's see where you can rent bicycles. For the overview we extract the rental stations from the complete list, e.g. by grouping them according to the name of the starting station (column 5) and selecting only this column as the result.

We get 337 stations after all.

Now we collect the data of a station …

… and build the coordinate list of the stations.

With these data we can send the sprite to the individual positions, where we leave circles with the *stamp*-block.

```
+show+all+citibike+stations+on+map+
warp
  clear
  switch to costume (circle)
  for (i) = (1) to (length of (coordinates of stations))
    go to x: (item (2) of (row (i) of (coordinates of stations))) y:
            (item (3) of (row (i) of (coordinates of stations)))
    stamp
  switch to costume (DataSpriteIcon)
```

At least in Midtown Manhattan, we don't have to worry about finding a rental station!

Now we want to have a closer look at the rental station Broadway - corner 41 Street (No. 55). To do this, we look for all records from the list that start or end at this station. That's 5005 events that day. Times are entered in this list together with the (same) date. We can throw this out (*split* with " ") and reduce it to the hour (*split* with ":"). We then have a numerical scale with the unit "hour". Now we can see what's going on at the station during the individual hours of the day. And we can graphically display this as usual using the *PlotSprite*.

```
+connections+to+or+from+station+(name)+
report
  keep items
    <(item (5) of (=) = (name)) or (item (9) of (=) = (name))>
  from (data)
```
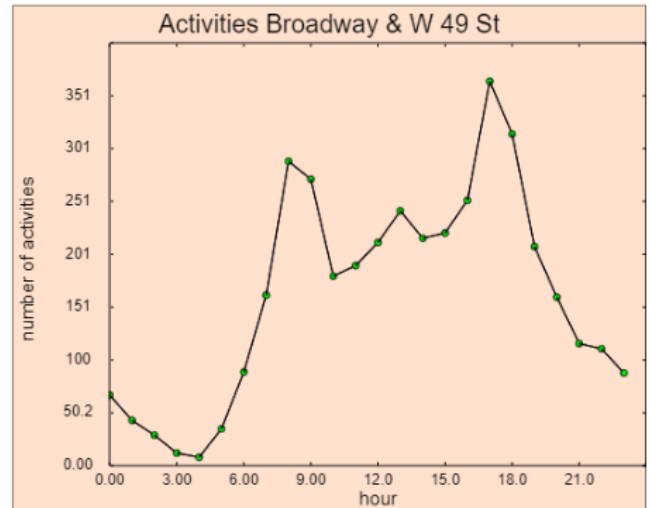
```
+reduce+time+columns+of+(table)+to+hours+
script variables (result)
warp connections to or from station (item (55) of (stations))
  set result to (list)
  add (column of (column (1) of (table)) to (result)
  add (column of
    map (item (1) of (split (item (2) of (split ( ) by ( ))) by ( ))) over
      (column (2) of (table))
    to (result)
  add (column of
    map (item (1) of (split (item (2) of (split ( ) by ( ))) by ( ))) over
      (column (3) of (table))
    to (result)
  for (i) = (4) to (15)
    add (column of (column (i) of (table)) to (result)
```

```
set (borrowing data) to (reduce time columns of (borrowing data) to hours)

set (plot data) to (number of column (1) of (borrowing data)
                    grouped by column (2))

+draw+diagram+of+activities+at+station+(n # = 1)+max=+
(max # = 100)+
  new costume width (500) height (400)
  color (255) (225) (205)
  go to x: (0) y: (0)
  set labels title (join (Activities) (item (n) of (stations))) x-label (hour)
    y-label (number of activities)
  set line attributes style (continuous) width (1)
    color (0) (0) (0)
  set datapoint attributes style (o_circle) width (5)
    connected (✓) color (0) (255) (0)
  set scale attributes precision (3) textheight x (12) y (12)
    number of x-intervals (8) number of y-intervals (8)
  set ranges to x [ (0) , (24) ] y [ (0) , (max) ]
  add dataplot with numeric scales of (plot data)
  add axes and scales

run (draw diagram of activities at station ( ) max= ( )) of (PlotSprite)
with inputs (55) (1.1) × (max of (column (2) of (plot data)))
```
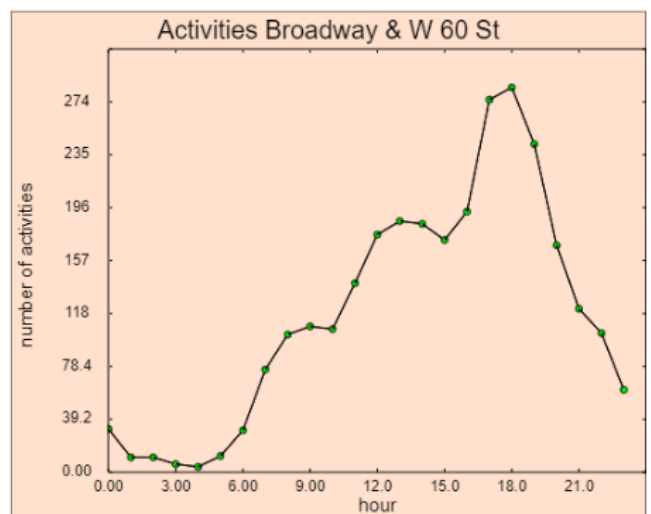
Activities Broadway & W 41 St

A few streets down the road, it looks similar.

Is that a general pattern?



Well, at Central Park the people get up later and the tourists are not there yet. But the museums always close at the same time.
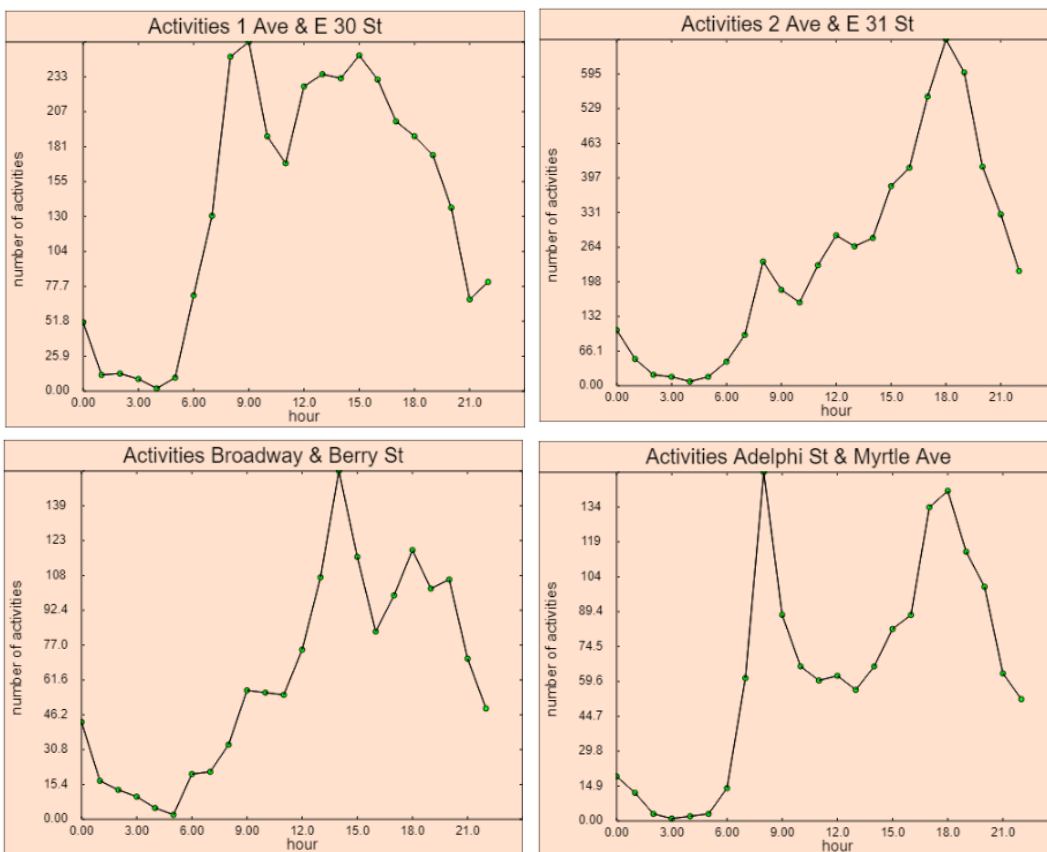


What can our programs learn from these data?

- *We could, for example, predict from the usual departures and arrivals as well as from the actual stock whether enough bicycles will be returned in time at a station or whether it would be better to transport some of them there.*
- *We could determine which rechargeable batteries are needed for eBikes from the average path lengths.*
- *We could determine whether women or men would rather borrow the bikes at a certain time of day and then make sure that the offer is right. We could do the appropriate thing for the age of the borrower.*
- *We could determine the borrowing data per bike and predict when repairs will be due. We could also do this, for example, depending on the location of the stands.*
- *We could try to generalize distributions from some stations in such a way that forecasts for others can be derived from them. So, when the museums close at Central Park, the program can "learn" from the old data in which districts the bikes will presumably be delivered and warn if there are not enough free slots available.*

etc.

**Tasks**:

1. Break down the activities of the stations according to arrivals and departures.
2. Write a forecast function that warns if there is a risk of a lack of bikes at a station in the next few hours.
3. For certain stations, display the connections to the most selected delivery stations graphically on the map using direct lines. Select the thickness of the lines according to the number of borrowing operations and the colors depending on the station. Are clusters formed?
4. Find out with the help of correlations-block whether there are correlations in rental behavior (e.g. with regard to times of day, location, ...) with gender, age, status of borrowers. You may have to replace the data with numeric data beforehand - similar to the times. Discuss possible consequences.
5. For a small section of Midtown (where everything is beautifully right-angled), find the coordinates of the street corners. Then develop a router that shows the shortest route to the nearest Citibike station.
6. The rental numbers depending on the time of day show quite a difference in different areas of Manhattan. Systematically examine similarities and differences and try to explain the results.

## 5.3    Star spectra [UniGOE]

Stars shine in different colors because they have different temperatures. In addition, the spectra differ in their absorption lines. We want to investigate this in more detail. We use a *PlotSprite* and a *DataSprite* as an assistant, e.g. for loading and preparing the data. Inside of this we start to work.

We get some star spectra (source: [UniGOE]) and save them as a text file. We read in such a file. In the first line we see the star name after the column captions. We isolate it and store it in the variable ***starname***.

We know the star's name now. If you search the Internet for it, you will find a wealth of information about it. For repeating the loading process with other data, we encapsulate it in a separate block. After its execution, the actual star data are available as a table. What's unpleasant about this is the very different order of magnitude of the data in the two columns. We therefore normalize the second column using the mean value and save the result as ***normalized data***.





With this the *DataSprite* has done its duty for now. We change to the *PlotSprite.*

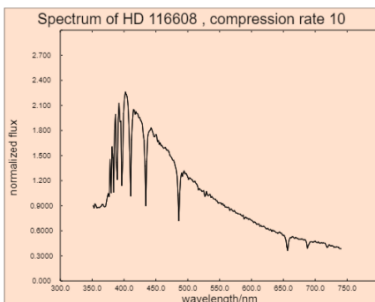With the normalized data you can quickly create a chart in *PlotSprite*.



The sloping course with some prominent absorption lines is clearly visible. But is it necessary to have all spectral data for this realization? Maybe it is enough to reduce the amount of data by averaging. We introduce a compression factor ***compression rate*** and complete the script before the diagram is created.
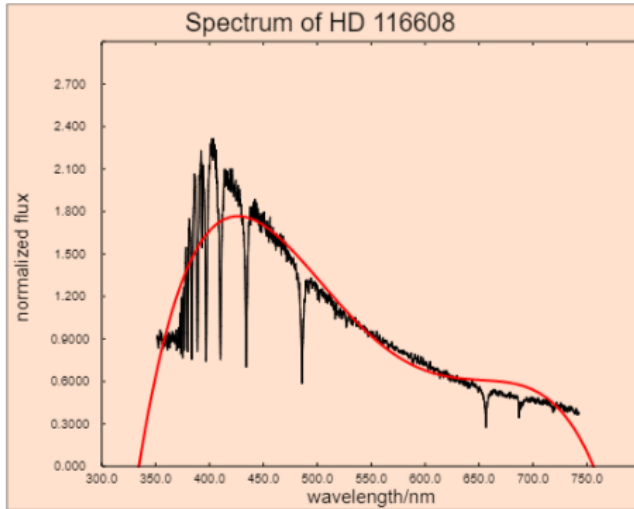


The factor 5 does not change much. So, let's keep trying.

One can see that the temperature-dependent course of the spectrum is hardly changed. Only the absorption lines are lost.

Thus, the type of spectrum should be describable by an interpolation polynomial, e.g. 4th degree. We load the *MathSprite* additionally and try it like this:
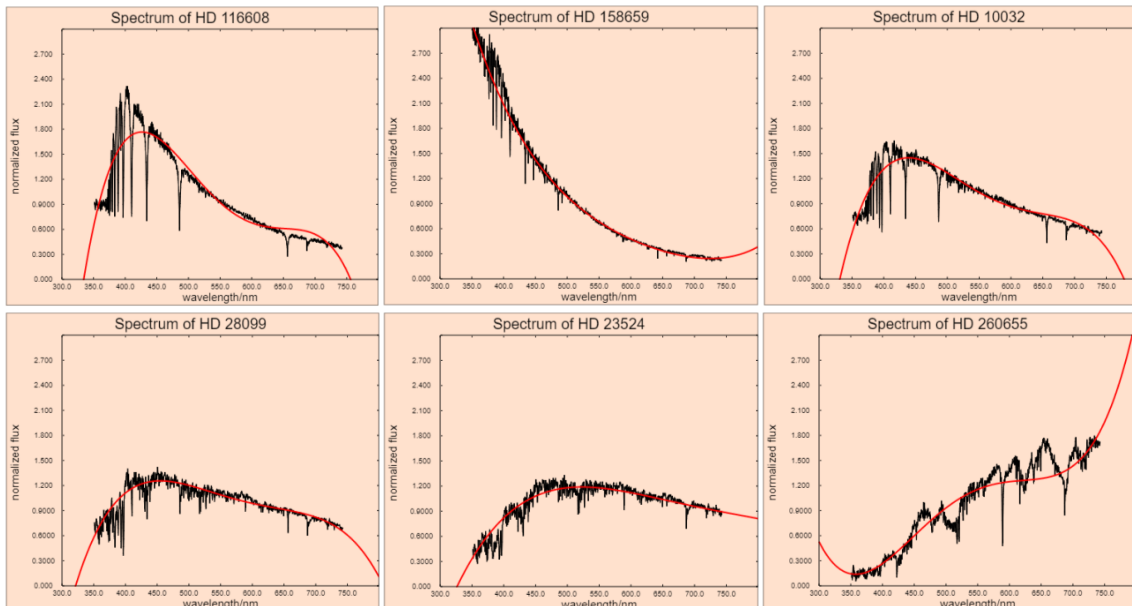


Spectrum of HD 116608

So, this works great! If we log the polynomial parameters during the test, we can easily distinguish the star types by means of the parameter ranges.

| 7 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | star name | a4 | a3 | a2 | a1 | a0 |
| 2 | HD 116608 | -1.2493580327340172e-9 | 0.0000028868621087800814 | -0.002459579857425176 | 0.9103088865090065 | 0.9103088865090065 |
| 3 | HD 158659 | 1.565259017017166e-10 | -3.963032080178107e-7 | 0.0003879661846290463 | -0.17622312866994078 | -0.17622312866994078 |
| 4 | HD 10032 | -7.27005023271818e-10 | 0.0000016949298479991264 | -0.001462425925103779 | 0.5500801501694278 | 0.5500801501694278 |
| 5 | HD 28099 | -4.0018935572381893e-10 | 9.399457129604694e-7 | -0.0008209889485783107 | 0.31410721917213127 | 0.31410721917213127 |
| 6 | HD 23524 | -8.18301248511472e-11 | 2.3253458278204257e-7 | -0.00024615800544876965 | 0.11348374829256708 | 0.11348374829256708 |
| 7 | HD 260655 | 6.248027476637483e-10 | -0.000001337322548726115 | 0.0010450333683869723 | -0.3486709605339992 | -0.3486709605339992 |



If you feed the polynomial coefficients to a neural network, it quickly learns to roughly assign a diagram to a star type. The program can therefore "learn" which parameter intervals belong to which star classes based on the old data. If you enter the data of a new star, it determines the coefficients of the polynomial and then makes a well-founded prognosis about what kind of star it could be.

**Tasks**:

1. Set up an interpolation polynomial of the lowest possible degree for the uncompressed spectrum data. Which points should be selected for this?  Are there any differences between these polynomials and the results of the method shown above?
2. Develop a script that assigns an unknown spectrum to one of the previously occurring types.
3. Develop a method to examine the most prominent absorption lines more closely. Enlarge them for stars of the same class and try to determine differences "automatically". Discuss your ideas before realization.

## 5.4    Classification of stars according to the kNN method

In the Hertzsprung-Russel diagram (see Wikipedia) the luminosity of stars is plotted above their star class. The result is a kind of line from top left to bottom right, the "main sequence". On this line stars like the sun are mostly located. Right above the main row we find the red giants, left below the main row the white dwarfs. That's enough for first. (Picture source: [HR])
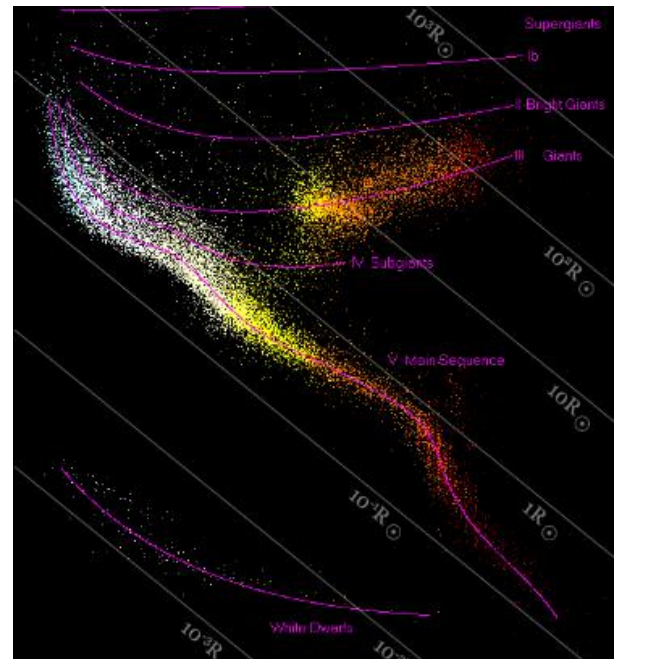
We want to classify new stars in this diagram using the k-next neighbor (kNN) method: As training data we generate a list of stars with their coordinates (simply as image coordinates in the diagram) and their type. If we want to classify a new star, we determine its position in the diagram and look for the nearest k (e.g. k=5) neighbors. Then we determine the most frequently appearing star type in this list. We assign it to the new star.

First of all, we need a picture of the Hertzsprung-Russel diagram ([HR]). We import it into Snap! as costume of an *ImageSprite* and generate the required data from it.

We generate the training data by specifying a star type and then clicking on some points in the diagram that correspond to this type.

Since we want to draw on the image, we work with a copy of the HR diagram so as not to alter the original.

Then we can classify new stars by clicking on them (here) and labeling them.

We set some properties for the representation ...

...and draw a circle at the location of the star.

Then we determine the five nearest neighbors and the number of occurrences of their type. In the result we delete the headings and sort the list in descending order. The type of the new star is then the first element in the first line. We write this next to the star.
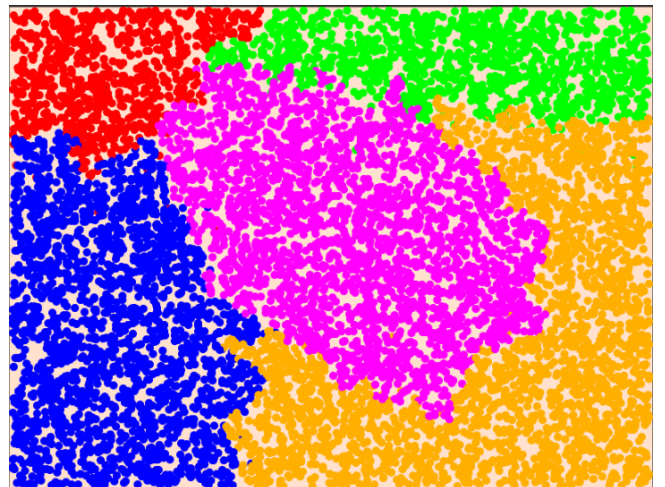
The result:



**Tasks**:
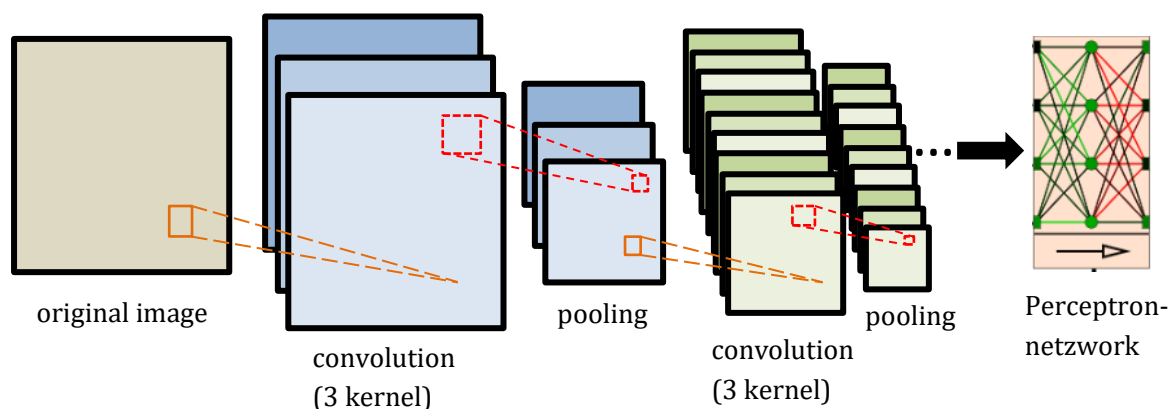
1. Add the newly classified stars to the sample list so that they are included in further classifications.
2. Draw differently colored dots at the correct places on the sprite for the different types of stars instead of labeling them.
3. Run the process for randomly selected points. Is the pattern always the same? Do completely different or similar patterns emerge? What does it depend on?

## 5.5    Character recognition with a convolutional network

The immense number of parameters in fully connected perceptron networks and the resulting need for huge amounts of training data has led to other network variants to drastically reduce this number. One of these is the *Convolutional Neural Network* (CNNs), where the amount of input data for the perceptron network is reduced. This type of network is used very successfully in image and speech recognition, for example.

CNNs reduce the amount of data by first applying several *kernels* in a multi-stage process, which filter out certain properties of e.g. an image (edges, oval surfaces, ...) and thus lead to several *feature maps*, which usually have the same size as the original. This first increases the amount of data. Afterwards, a non-linear activation function (*reLU*) is usually applied to the feature maps, followed by a *pooling operation* that reduces the amount of data again. This is usually called *Max-Pooling*, where the maximum value is determined from a section of the data. If you do this with a "window" that is moved across the feature map with a certain *stride*, each pooling step creates a value of the next, reduced feature map.



original image          convolution              pooling          convolution          pooling          Perceptron-
                        (3 kernel)                                (3 kernel)                             netzwork

As an example, let's take a kernel that filters vertical lines: it colors a point white if there is a second pixel next to the point, otherwise black. In the "folded" image we can see vertical lines of the original as bright spots. If it is not so important where exactly these lines are, we do not lose too much information during the pooling. A white point in a feature map after various pooling processes means then: "*In this area there was a vertical line somewhere*". Using such data from several feature maps, it can be deduced, for example, that there was also a horizontal line, i.e. a corner. If we had searched for "pink" areas and "oval" shapes, the chance of identifying faces would not be so bad.

We now want to build a model for such a CNN that can distinguish the handwritten digits *zero* and *one*. For this we use a *DataSprite* for auxiliary operations, an *ImageSprite* for the actual image and - of course - a *NeuralNetSprite* for the perceptron network at the end of the chain. Another, "normal" sprite called *Control* is supposed to control the operations. To make the model easier to use, we add some buttons and a pen to make the interface clearer. In the screenshot, the image to be analyzed is located at the top of the box, the neural network shows its result at the bottom. In between, the different intermediate layers are run through and displayed from top to bottom. As an addition, the model contains the possibility to draw your own numbers.
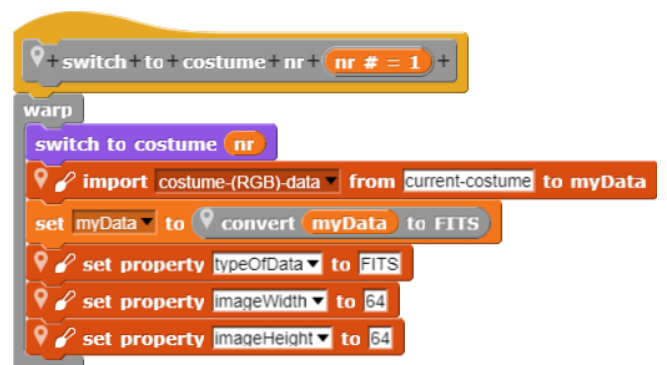
Our CNN is trained with 10 digits of 64x64 pixels each for the zeros and ones. Afterwards it should "recognize" these and other handwritten ones. Actually, we would have to train several kernels of our CNN especially for this task. Instead, we take only two known kernels for the recognition of vertical and horizontal lines, because by limiting the number of kernels to two, everything can be displayed on the screen and the results can even be interpreted halfway. (The recognition rate suffers severely from this!) So, we train only the perceptron network with four input values.

In the image above, after two stages of reduction, four feature maps of 16x16 pixels each are left over, each of which has undergone the *Convolution → reLU → Max-Pooling* operations twice: on the far left with the kernel for vertical lines, then with both kernels in different order, and finally twice with the kernel for horizontal lines. The numbers below indicate the average brightness value measured over the entire image. If we apply this to different digits, the possibility to measure differences between zeros and ones becomes apparent despite the very simple procedure.
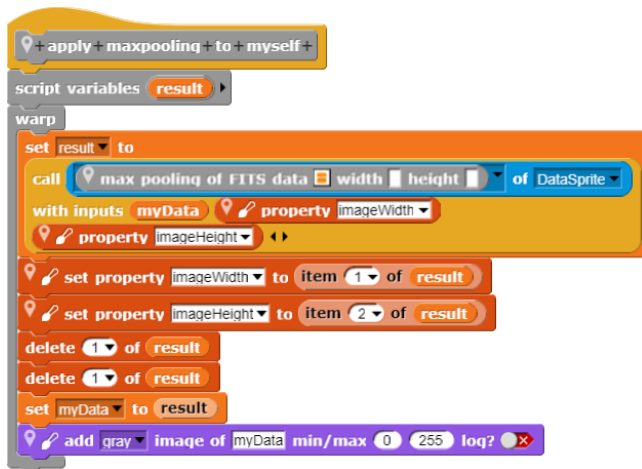


Let us look at the functionalities of the individual objects:

The *ImageSprite* should import the data of a new costume as gray values into its data area. This is very simple.
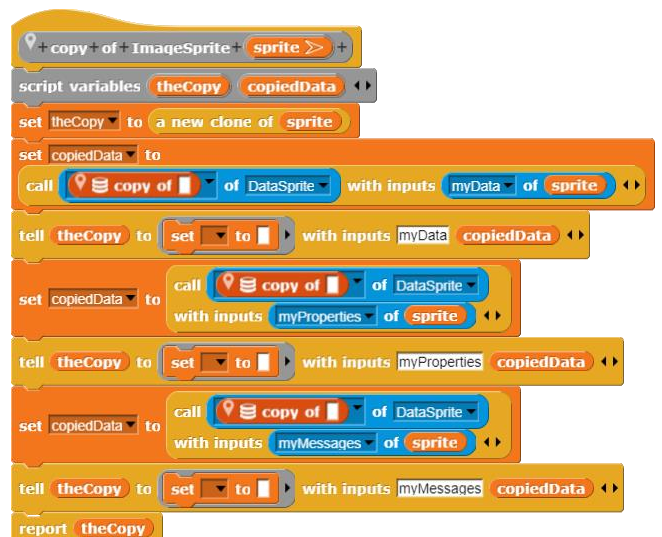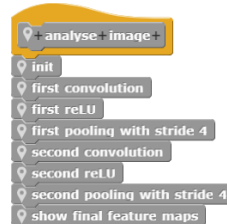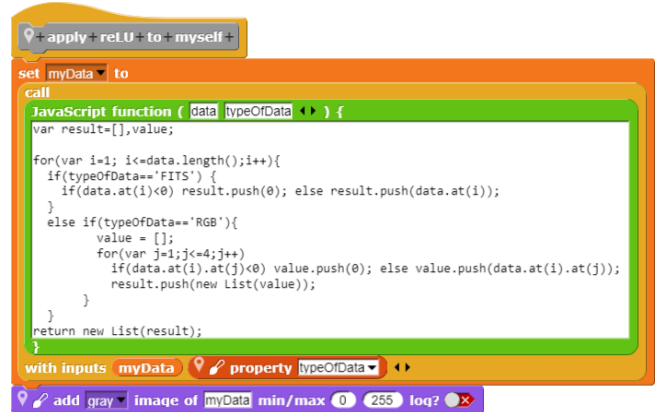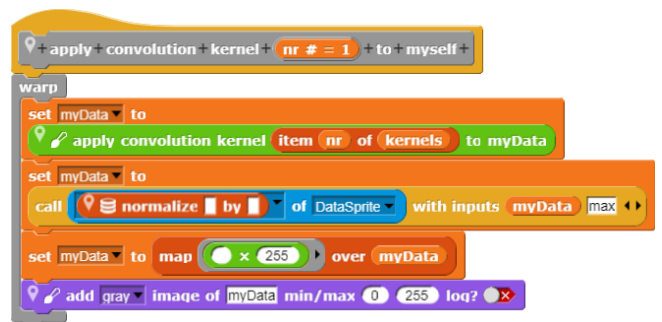
Furthermore, the sprite and its clones should be able to perform the three operations of a CNN. With some help of the *DataSprite* this is also possible.



The *ImageSprite* does not have to master more for our purposes.



The *Control* sprite has to ask the *ImageSprite* to change the costume and analyze it afterwards. In doing so, it strictly adheres to the specifications for CNNs. The **init** method only takes care of drawing the lines on stage. The other methods work with two layers of CNN, **first layer** and **second layer**, each containing the versions of the characters that appear on stage. So that they do not interfere with each other, copies of the *ImageSprite* are used, not clones. The *DataSprite* again helps with copying.



After the required copies have been made, *Control* asks them to perform the relevant CNN operation. Finally, the now rather small (4x4 pixel) clones of the last layer are displayed as "*final feature maps*", greatly enlarged. These are used to train the neural network.

The neural network in the form of a *NeuralNetSprite* should produce the largest output at output 1 when there are zeros and at output 4 when there are ones. This is of course completely arbitrary. The current output value is determined by the function **read output**. With its components the net can be trained, if we succeed in determining the mean values from the last layer of **second layer**. We model these aptly with the **softmax** function.
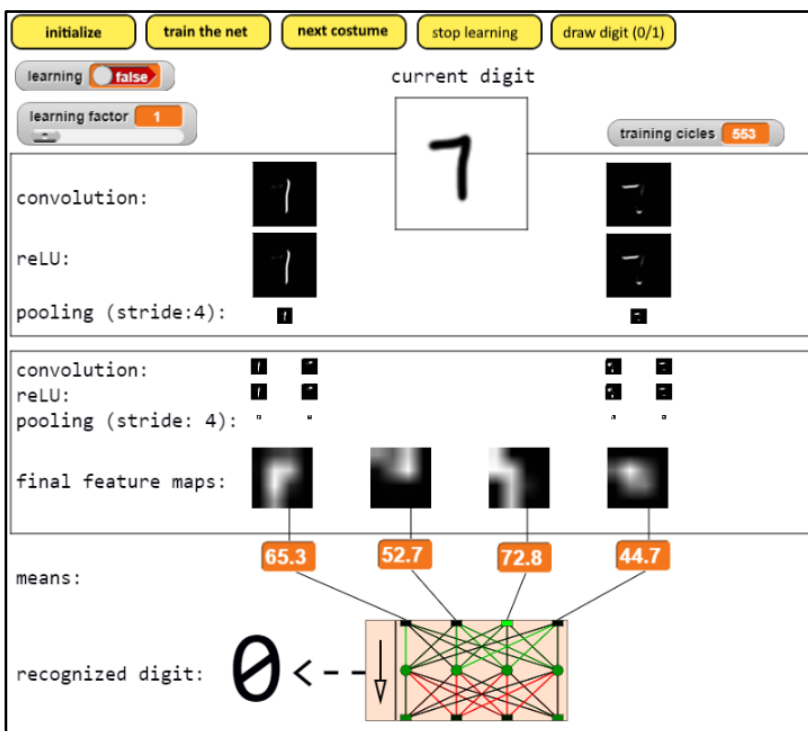




And - has the net learned something?



Well - there is still room for improvement!

# Latest hints

Machine learning consists to a large extent of preparing data - whether it is tabular data or images. The actual learning processes of the machines then consist of the parameter adjustments that result from the data. Since both can be visualized well, there is a broad field for beginners with many transitions to the area of "computer science and society".

Examples for the application of the operations of the *ML.SpriteLibrary*, especially the convolution with help of a kernel, can be found abundantly in [DBV].

# List of examples

| Subject | Seite |
|---|---|
| False colour representation | 12 |
| Image import from file | 12 |
| Access to local data of another sprite | 13 |
| Access to the data of a clone | 14 |
| Calling a global method by another sprite | 15 |
| Call of a global method with parameters by another sprite | 15 |
| Calling a local method with parameters by another sprite | 16 |
| Access to the code of a local method by another sprite | 16 |
| Importing data from the costume | 18 |
| Data import from CSV file | 18 |
| Data import from SQL query | 18 |
| Data import from JSON file | 19 |
| Importing data with the mouse | 19 |
| Measuring the overall brightness of an image area | 20 |
| Data export to CSV file | 20 |
| Data export to text file | 20 |
| Calculation of correlations (US census income dataset) | 23 |
| Data preparation (New York Citibike tripdata) | 25 |
| Function graphs | 27 |
| Diagram of a point set | 28 |
| Regression line | 28 |
| Mixed data chart | 29 |
| Histogram of an image | 29 |
| Random graphic | 32 |
| Planet Transit | 33 |
| Image reflection | 34 |
| Edge detection | 35 |
| Interpolation polynomial through n points | 38 |
| SQL Queries | 41 |
| Training of a neural network | 45 |
| Traffic sign recognition | 46 |
| Using the World Map Library | 53 |
| Diagram creation (New York Citibike tripdata) | 54 |
| Star Spectra | 57 |
| kNN method in the Hertzsprung-Russel diagram | 61 |
| Character recognition with a convolutional neural network | 63 |

# References and sources

[Albon]          Albon, Chris: Machine Learning Kochbuch,
                 O'Reilly, 2919

[Baumann]        Baumann, Rüdeger: Didaktik der Informatik,
                 Klett, 1990

[Census]         https://archive.ics.uci.edu/ml/datasets/census+income

[DBV]            Burger, W., Burge, M.-J-: Digitale Bildverarbeitung – Eine Einführung mit
                 Java und ImageJ, Springer 2006

[FITS]           de.wikipedia.org/wiki/Flexible_Image_Transport_System

[Goodfellow]     Goodfellow, I.; Bengio, Y.; Courville, A.: Deep Learning,
                 MIT Press, 2016

[Grus]           Grus, Joel: Einführung in Data Science,
                 O'Reilly, 2016

[HOU]            Hands-On Universe: handsonuniverse.org/

[HR]             https://studylibde.com/doc/2985884/hertzsprung-russell--und-farb-hel-
                 ligkeits

[JSON]           Popular Baby Names: https://catalog.data.gov/dataset/most-popular-
                 baby-names-by-sex-and-mothers-ethnic-group-new-york-city-8c742

[NYcitibike]     https://www.citibikenyc.com/system-data

[Minsky]         Minsky, Marvin: Computation: Finite and Infinite Machines,
                 Prentice-Hall, Englewood Cliffs, New York, 1967

[Modrow1]        Modrow, Eckart: Technische Informatik mit Delphi,
                 emu-online, 2004

[Modrow2]        Modrow, Eckart: Zur Didaktik des Informatikunterrichts – Band 2,
                 Dümmler, 1992

[Rojas]          Rojas, Rau´l: Neural Networks - A Systematic Introduction,
                 Springer Berlin, 1996

[SAP]            www.sap.com/germany/products/leonardo/machine-learning.html

[SchulAstro]     www.schul-astronomie.de

[SQL]            Modrow, Eckart: Informatik mit Snap!,
                 http://ddi-mod.uni-goettingen.de/InformatikMitSnap.pdf

[SZ]             Süddeutsche Zeitung: 10. April 2019
                 www.sueddeutsche.de/wissen/schwarzes-loch-bild-1.4404130

[UniGOE]         Institut für Astrophysik, Universitaet Goettingen