

Research methodology - general principles of building a software system for data processing.

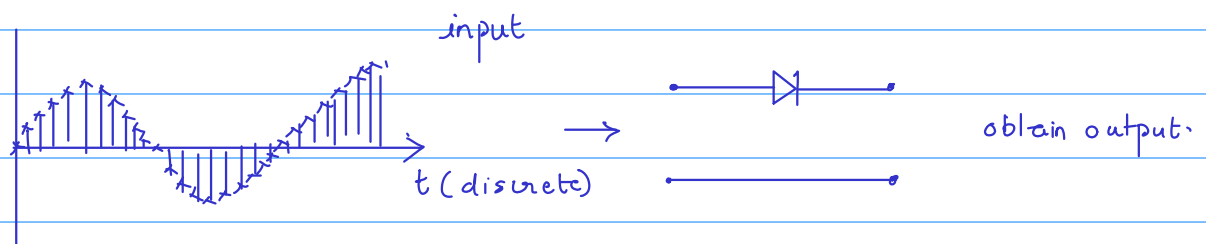
Course outline: Data processing is an essential component of any research venture today. Several data processing tools, software, libraries, and languages exist. In this course, we will explore some of the general principles of setting up a data processing system in a general purpose language (Python). Using these general principles in your own research would lead to well designed systems, with readable code, easier debugging, and reproducible research. This exploration is done via the detailed design of an example.

Lecture 1 Outline:

- Introduction to the course, Procedural programming
- Python (and Matlab) introduction.
- Variables and other data storage mechanisms
- Operations on variables / data processing
- Conditional flow
- Functions & libraries.

Assignment 1:

Simulating a rectifier operating on a sine wave - procedural programming.



```
1 import numpy as np
2 import scipy
3 import pylab as pl
4
5 # Sample times
6 t = np.arange(0, 2, 0.001)
7
8 # Input sine wave
9 f = 50
10 input = np.zeros(len(t))
11 for i in range(len(t)):
12     input[i] = np.sin(2 * np.pi * f * t[i])
13
14 # Simulate rectifier
15 output = np.zeros(len(t))
16 for i in range(len(t)):
17     if input[i] < 0:
18         output[i] = 0
19     else:
20         output[i] = input[i]
21
22 # Visualize output
23 pl.plot(t, input)
24 pl.plot(t + 0.001, output)
25 pl.show()
```

} libraries of functions which we can make use of. When we start out in making data processing systems, try to make use of standard libraries as much as possible without "re-inventing the wheel"

→ doing the same operation multiple times.

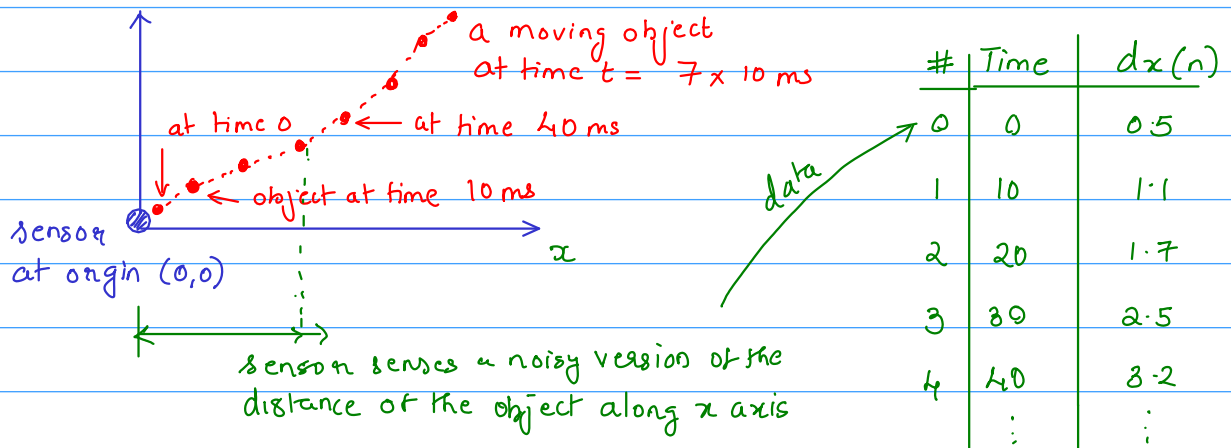
→ use of an inbuilt library function

} use of conditional flow.

} use of inbuilt functions for visualization.

Lecture 2 - an exploration of the general principles via an example.

We will use the following data processing example to illustrate the design of a software system that will help us to explore some good ideas to keep in mind.



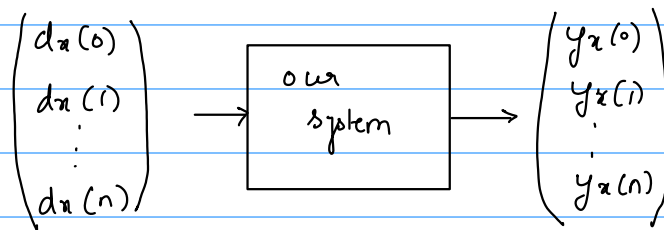
We need to build a system to clean this data. We will use the following algorithm to obtain the cleaned data which we will denote as y_x .

$$y_x(0) = dx(0)$$

$$y_x(n) = (y_x(n-1) + v \cdot \Delta t)(1 - \alpha) + \alpha \cdot dx(n)$$

(the motivation for this has been discussed in class).

Our system is a block that can be represented as



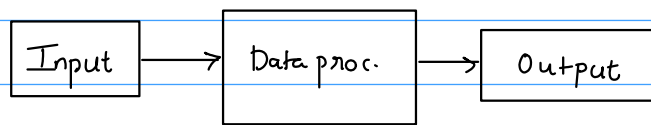
We can make a simple procedure to do this:

```
1 import numpy as np
2
3 v = 1
4 delta_t = 0.001
5
6 def clean_data(dx, alpha):
7     yx = np.zeros(len(dx))
8     yx[0] = dx[0]
9     for i in range(1, len(dx)):
10        yx[i] = (yx[i - 1] + v * delta_t) * (1 - alpha) + alpha * dx[i]
11
```

but is this the best design that we can have for our system?

In order to design a "good" system it is better to break up your system to subsystems, design those subsystems very well. In our case, it is actually better to use object oriented programming to design systems using this principle.

What are the subsystems you should design?



} design as subsystems so that you have lot of flexibility in doing experiments.

Each block can be represented by a class in Python.

```
1 import numpy as np
2
3 class Filter:
4     def __init__(self, alpha):
5         self.alpha = alpha
6         self.y = 0;
7
8     def initialize_y(self, d):
9         self.y = d
10
11     def process_one_data_point(self, d):
12         self.y = (self.y + 0.001) * (1 - self.alpha) + self.alpha * d
13         print self.y
14
15     def process_multiple_data_points(self, d):
16         for i in range(len(d)):
17             self.process_one_data_point(d[i])
18
19     def file_output(self):
20         do something for writing y to a file
21
22 f = Filter(0.1)
23 d = [0.1, 0.2, 0.3]
24 f.initialize_y(d[0])
25 f.process_multiple_data_points(d[1:3])
```

← Class for the data processing block.

} classes will have a constructor for initialization

← Use descriptive variable names, words separated by underscores or Capital Letters.

→ blocks will have parameters. These are useful for experimentation.

→ write functions that have a well defined objective.

} always test your subsystems

The output block represented as a class

```
1 import numpy as np
2
3 class Output:
4     def __init__(self, type_of_output, filename):
5         self.type_of_output = type_of_output
6         self.filename = filename
7
8     def produce_output(self, t, y):
9         if self.type_of_output == "csv":
10             np.savetxt(self.filename, [t, y], delimiter=",")
11         elif self.type_of_output == "excel":
12             pass
13         elif self.type_of_output == "fig":
14             pass
15
16 o = Output("csv", "temp.csv")
17 o.produce_output([0,1],[1,2])
```

} have multiple options for the output of your data processing system

} again test your subsystems as much as possible.

Lecture 3 – continuation of Lecture 2.

Let us now look at the input block. We can assume that input data might come in many formats. Also for testing our code and doing simulation experiments we would also need a simulator for simulating the motion of the object and producing the data as shown in the table. So how should the input block be implemented as a class?

- Should the Input block implement a file reader as well as a simulator for data?

maximum flexibility + intermediate data: Simulator → file → input block.
control flexibility by your workflow.

```

1 import numpy as np
2
3 class Input:
4     def __init__(self, type_of_input, filename):
5         self.type_of_input = type_of_input
6         self.filename = filename
7
8     def get_input(self):
9         if self.type_of_input == "csv":
10             pass
11         self.type_of_input == "excel":
12             pass
13         # return - what should be returned?
14     def get_one_input(self):
15         if self.type_of_input == "csv":
16             pass
17         self.type_of_input == "excel":
18             pass
19         # return - what should be returned?
20

```

} different file formats for inputs

} different ways of obtaining inputs.

The simulator for testing our code or for experimentation.

```

1 import numpy as np
2
3 class Simulator:
4     def __init__(self, v, delta_t, x0, sigmasq):
5         self.v = v
6         self.delta_t = delta_t
7         self.sigmasq = sigmasq
8         self.x = x0
9         self.times = []
10        self.dx = []
11        self.actual_x = []
12
13    def simulate(self, number_of_times):
14        time = 0
15        for i in range(number_of_times):
16            self.times = self.times.append(time)
17            self.dx = self.dx.append(self.x + np.random.randn() * sigmasq)
18            self.actual_x = self.append(self.x)
19
20            self.x = self.x + self.v * self.delta_t
21
22    def file_output(self, output_filename):
23        pass
24
25    def run_simulator(self, output_filename):
26        self.simulate()
27        self.file_output(output_filename)
28
29 # How to test this module?

```

} even within a block use simple functions for implementing the procedure

For simulated data it would be possible to actually see the performance of our algorithm, (because we know the ground truth). So we might need another module to do that. This metrics module should compare the actual data x -s with the filtered y -s. The comparison can be done in terms of sum of squared errors, mean squared error, or sum of absolute errors. How will you design a flexible Metrics module?

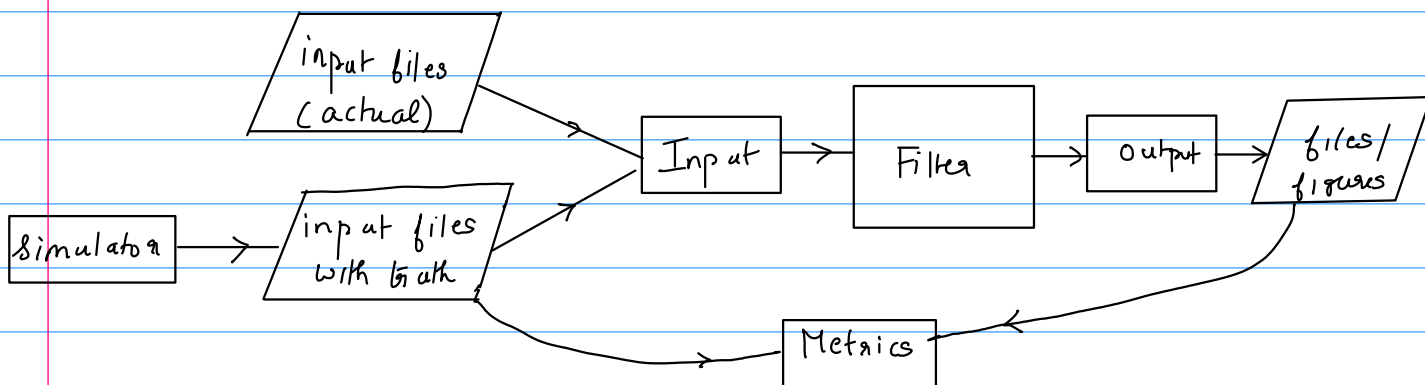
```

1 import numpy as np
2
3 class Metrics:
4     def __init__(self, metric_type):
5         self.metric_type = metric_type
6
7     def get_ground_truth(self, ground_truthfile):
8         pass
9         # what to return from this?
10    def get_output_data(self, output_file):
11        pass
12        # what to return from this?
13
14    def compute_metric(self, ground_truthfile, output_file):
15        truth = self.get_ground_truth(ground_truthfile)
16        data = self.get_output_data(output_file)
17
18        if self.metric_type == "mse":
19            self.compute_mse(truth, data)
20        elif self.metric_type == "sse":
21            pass
22        elif self.metric_type == "sae":
23            pass
24
25    def compute_mse(self, truth, data):
26        pass
27
28 # How to test this module?

```

again to have intermediate data available and for flexibility - choose input and output to be files.

Now our data processing system looks like this



Now we need to integrate all of this into a workflow.

```

1 from Simulator import *
2 from Input import *
3 from Filter import *
4 from Output import *
5 from Metrics import *
6
7 # Setup the parameters in your workflow
8 v = 1
9 delta_t = 0.001
10 x0 = 0
11 sigmasq = 1
12 input_data_file = "inputdata_file_experiment_1.csv"
13 infile_type = "csv"
14 output_data_file = "outputdata_file_experiment_1.csv"
15 outfile_type = "csv"
16 metrics_type = "mse"
17 alpha = 0.1
18
19 # Setup the components in your workflow or system
20 sim = Simulator(v, delta_t, x0, sigmasq)
21 inp = Input(infile_type, input_data_file)
22 metric_computation = Metrics(metrics_type)
23 outp = Output(outfile_type, output_data_file)
24 datafilter = Filter(alpha)
25
26 # Generate some test data
27 sim.run_simulate(input_data_file, number_of_times)
28 # A workflow which processes the data in blocks
29 data = inp.get_input()
30 datafilter.initialize_y(data[0]) # depends on the data type here
31 datafilter.process_multiple_data_points(data[1:])
32 outp.produce_output()
33
34 # A workflow which processes the data one at a time - good for interactive tuning, debugging etc

```

Other important facts to keep in mind.

- during experimentation and testing Great storage as being cheap. Multiple directories for experiments. Maintain temporary data.
- use version control for your code.

- git init
 - git add
 - git commit
 - git status
- } some common commands.

Evaluation.

Each group should take up a module and complete it

- complete all functions
- find and fix errors
- specify inputs and outputs
- comment the code
- test the module by writing test functions with sample inputs.

The modules are Filter.py, Input.py, Output.py, Simulator.py, Metrics.py
Experiment_1.py. Maximum marks is 10.

If more than two groups choose a module; say n groups choose ($n > 2$), then the maximum mark for each group will be $\left(\frac{10}{n-1}\right)$.