

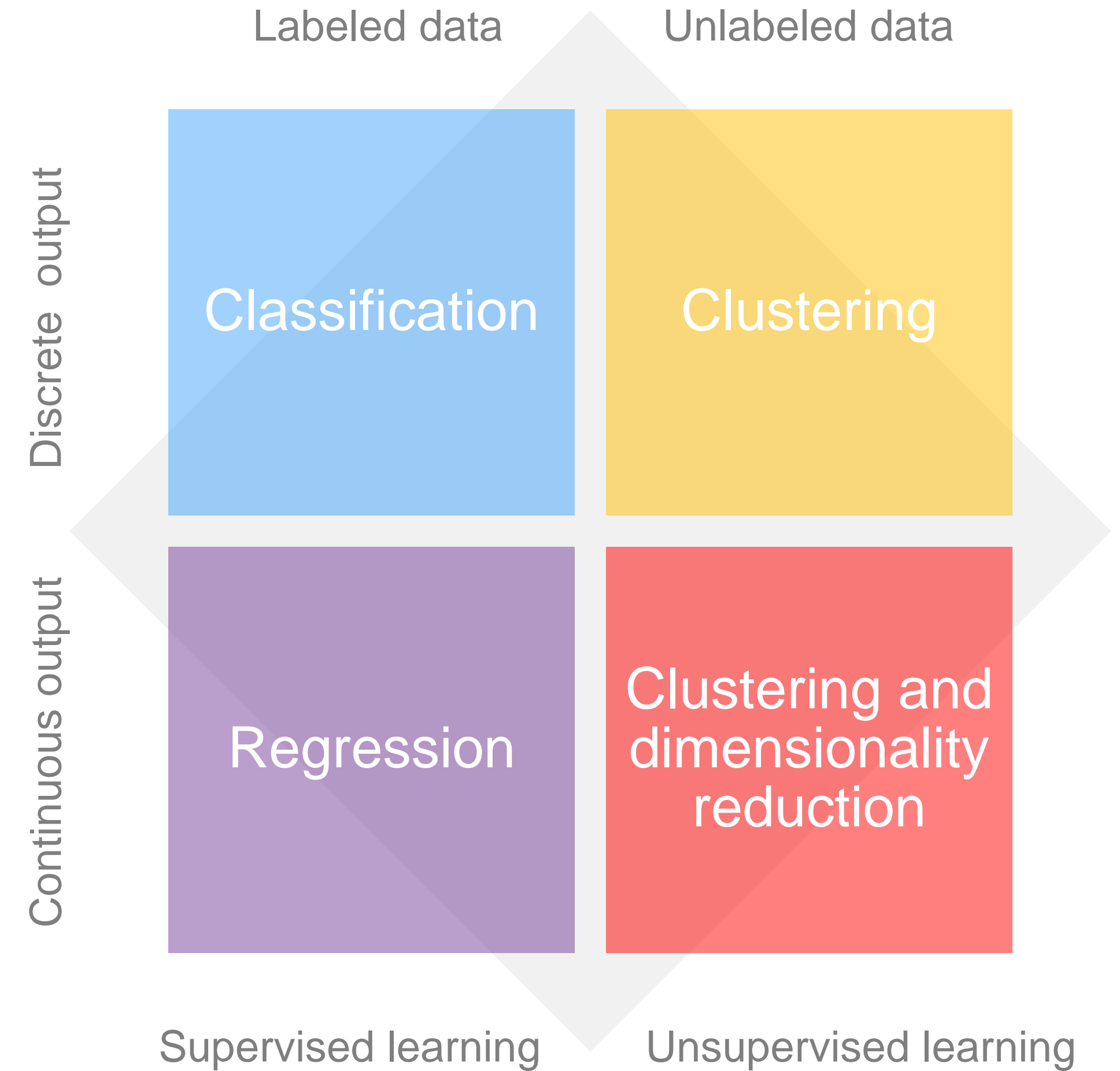
CREATIVE DATA MINING

Introduction to Python and Programming II

16.10.2017

Dr. Varun OJHA

Danielle GRIEGO



What we'll cover today

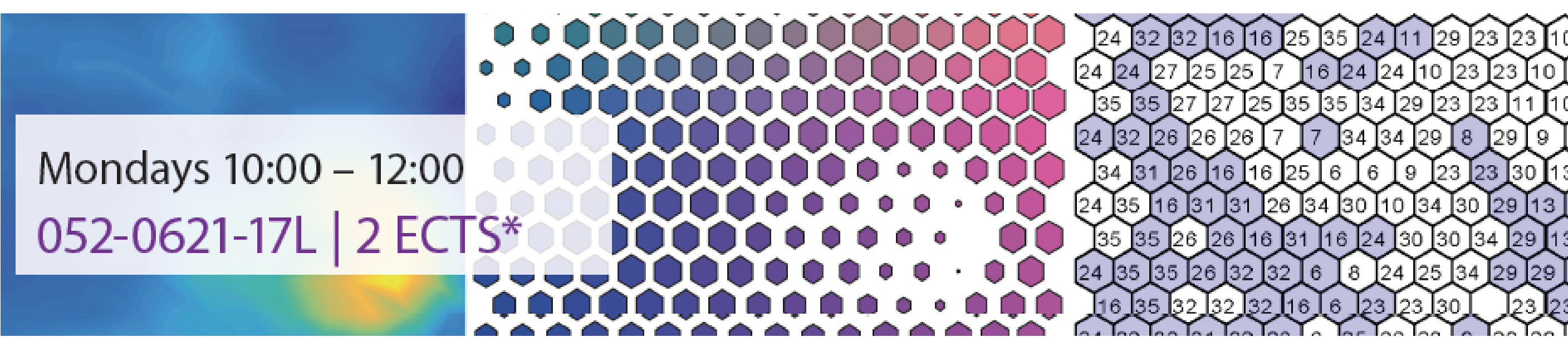
- Overview of exercise 1
- Revisit course schedule and objectives
- Introduce final project
- Recap from last week
- Python tutorial II
- Functions, modules, reading and writing files

Exercise 01

There are 100 students in a degree program; each student may get a grade: bad, average, good, and excellent in mathematics, physics, and chemistry tests. Each student is given a total 180 minutes to finish all three tests. A teacher calculates each student intelligence quotient (IQ) based on the student's test results, the time taken by the student to finish the three tests, and the age of the students.

From the given paragraph, please identify the following:

1. What are the variables/features in the given problem? Distinguish discrete and continuous variables.
2. What are the input variables?
3. What is the target variable?
4. How many examples are in the given problem?



Mondays 10:00 – 12:00
052-0621-17L | 2 ECTS*

Creative Data Mining

The Creative data mining course aims to provide aspirants a hands-on experience on machine learning (ML) tools and techniques for data processing and analysis. Since future technologies increasingly rely upon the ML, urban systems and architecture shall adopt it and aspirant should learn creative ways to apply ML to better understand urban systems. The course covers a wider range ML techniques including supervised and unsupervised learning methods for data analysis and pattern recognition that help to better understand urban system for improving urban life.

All methods taught in the course will be applied to a common project to evaluate various dynamics of the urban environment. Students will work with time-series and geo-referenced data including temperature, relative humidity, illuminance, noise, people density, and dust particulate matter. Subjective impression survey data will also be integrated into the student projects to further explore influencing factors of the urban environment on our perceptual experiences. A selected neighborhood in the city of Zurich will be used as the case study and each student will present the findings of their research question in a final project.

Additionally, there are two of non-architectural skills the participants can develop during this course. First is an introduction to programming where at a minimum they can successfully copy and paste code-snippets to customize the computational tools presented in the course. Second, how clustering methods like PCA or K-Means could be applied in an architectural context.

Where
HIT H 31.4 (Video wall)

25.09.2017	Introduction to the knowledge discovery process
02.10.2017	Fundamentals of supervised machine learning
9.10.2017	Introduction to python & programming I
16.10.2017	Introduction to python & programming II Introduce final projects
23.10.2017	Seminar Week- No Lecture
30.10.2016	Supervised learning problem solving in python (MLP & SVM)
06.11.2016	Fundamentals of unsupervised learning (K-means, DBSCAN, PCA)
13.11.2017	Unsupervised learning problem solving in python
20.11.2017	Review data and examples for final projects
27.11.2017	Project proposal discussions
04.12.2017	Exploration of Real-World problems Q&A Workshop I
11.12.2017	Q&A Workshop II
18.12.2017	Final Critique

Course objectives and outcomes

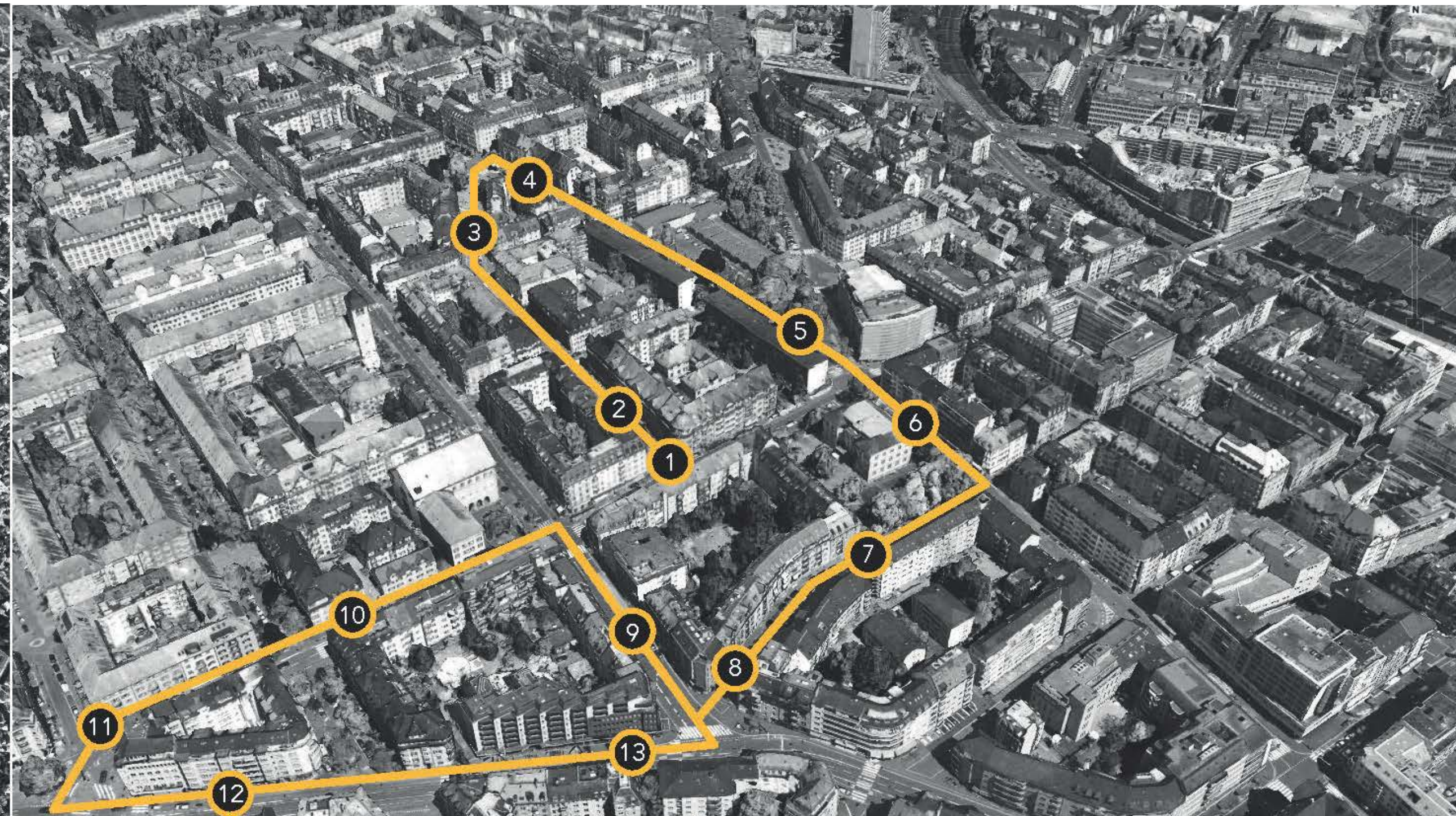
1. Become familiar with programming and integrating new tools in your work
2. Develop an interesting “research” question and learn how to answer it by:
 - Selecting appropriate data
 - Applying relevant analysis and visualization techniques
 - Interpreting and refining your results

Final projects: Possible dataset

ESUM- Analyzing trade-offs between Energy and Social performance of Urban Morphologies



Location Wiedikon Zürich



14 survey checkpoints along experimental path

Example

ESUM- Analyzing trade-offs between Energy and Social performance of Urban Morphologies

(Raw) data from 37 participants:

- Investigate impact of static (urban morphology) and dynamic features (environmental sensors) of the built environment on perception (using surveys and biofeedback data)

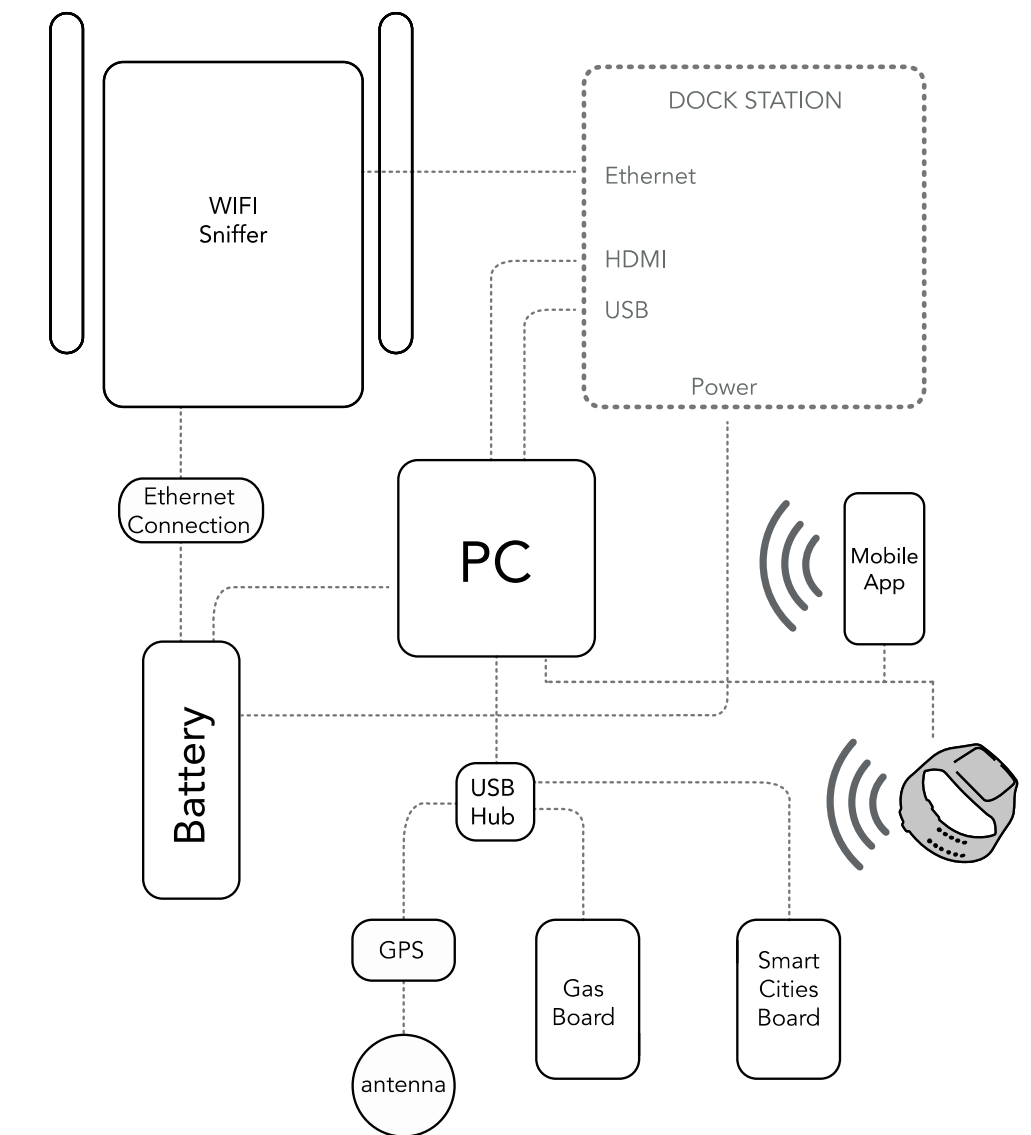
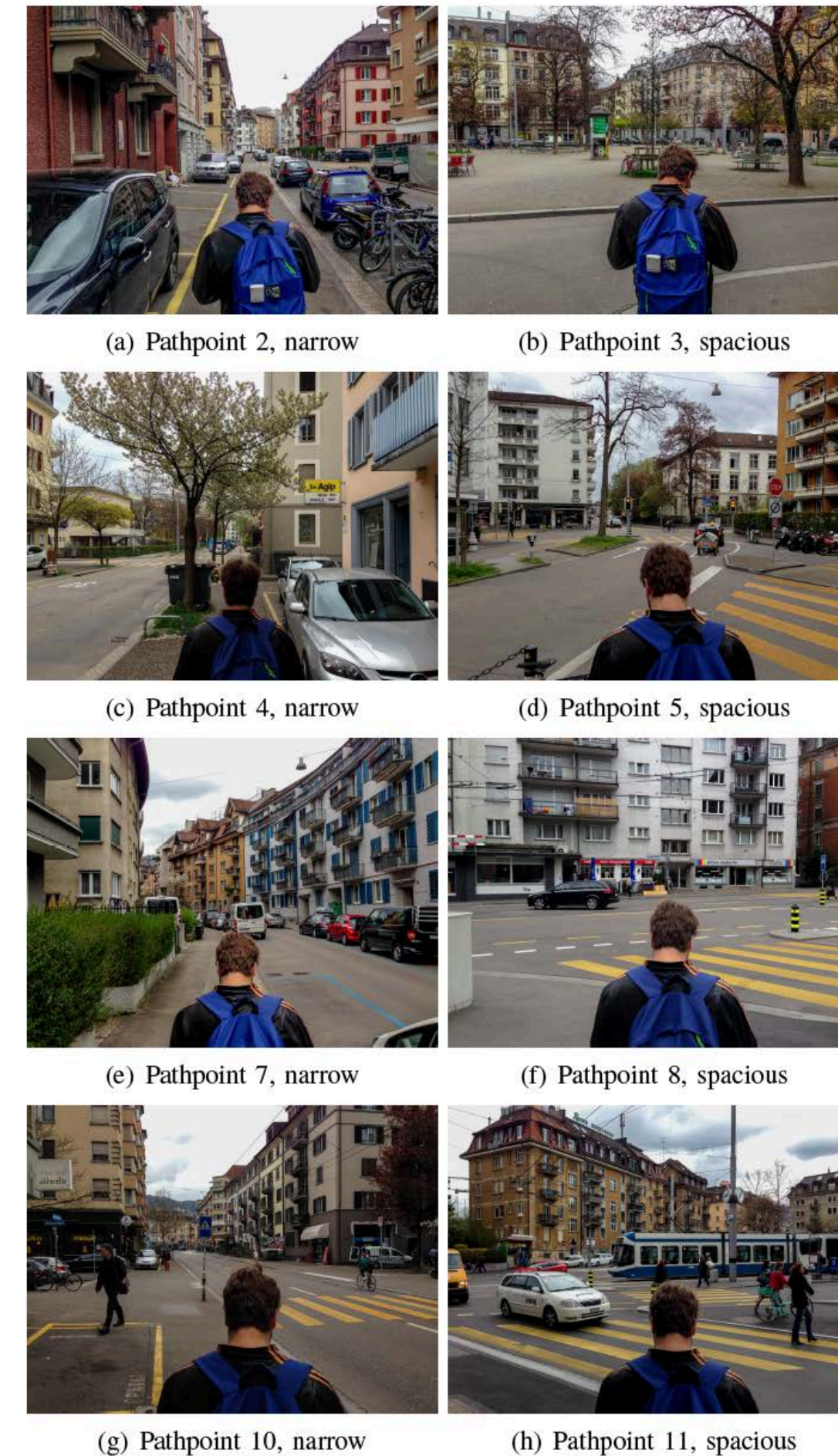
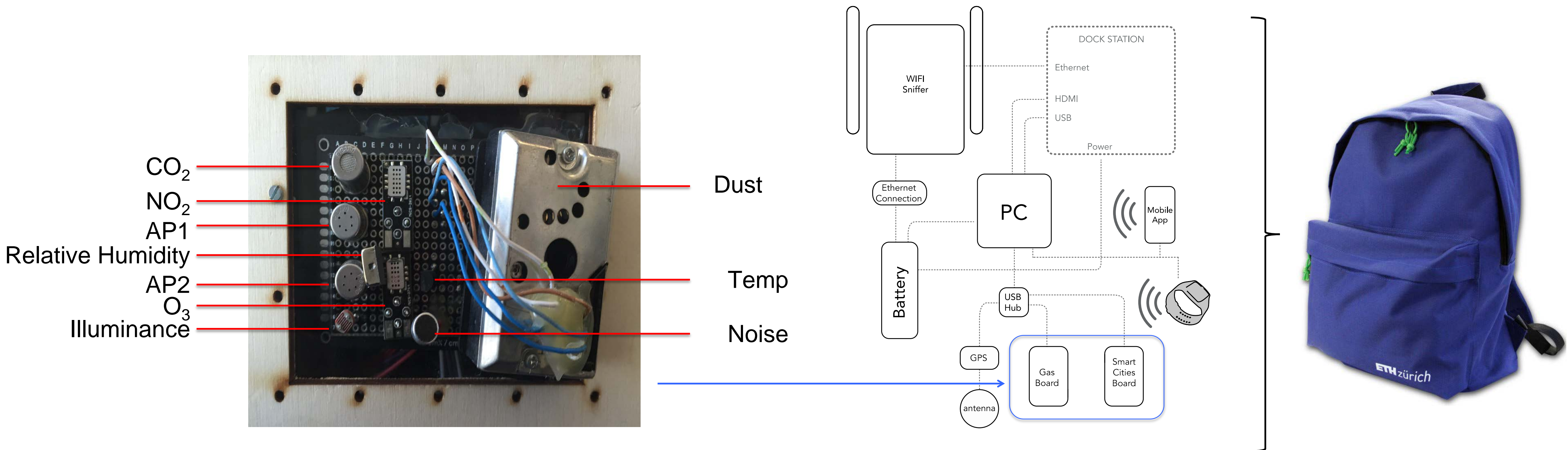


Fig. 2. Four instances of narrow-spacious spatial configurations and their corresponding pathpoints along the select path.

Mobile sensor equipment


Sensorbackpack with environmental and position sensors




Mobile Sensor equipment


Biofeedback wristband

E4 Sensors







PPG Sensor
Photoplethysmography Sensor - Measures Blood Volume Pulse (BVP), from which heart rate, heart rate variability (HRV), and other cardiovascular features may be derived




3-axis Accelerometer
Captures motion-based activity




Event Mark Button
Tags events and correlate them with physiological signals



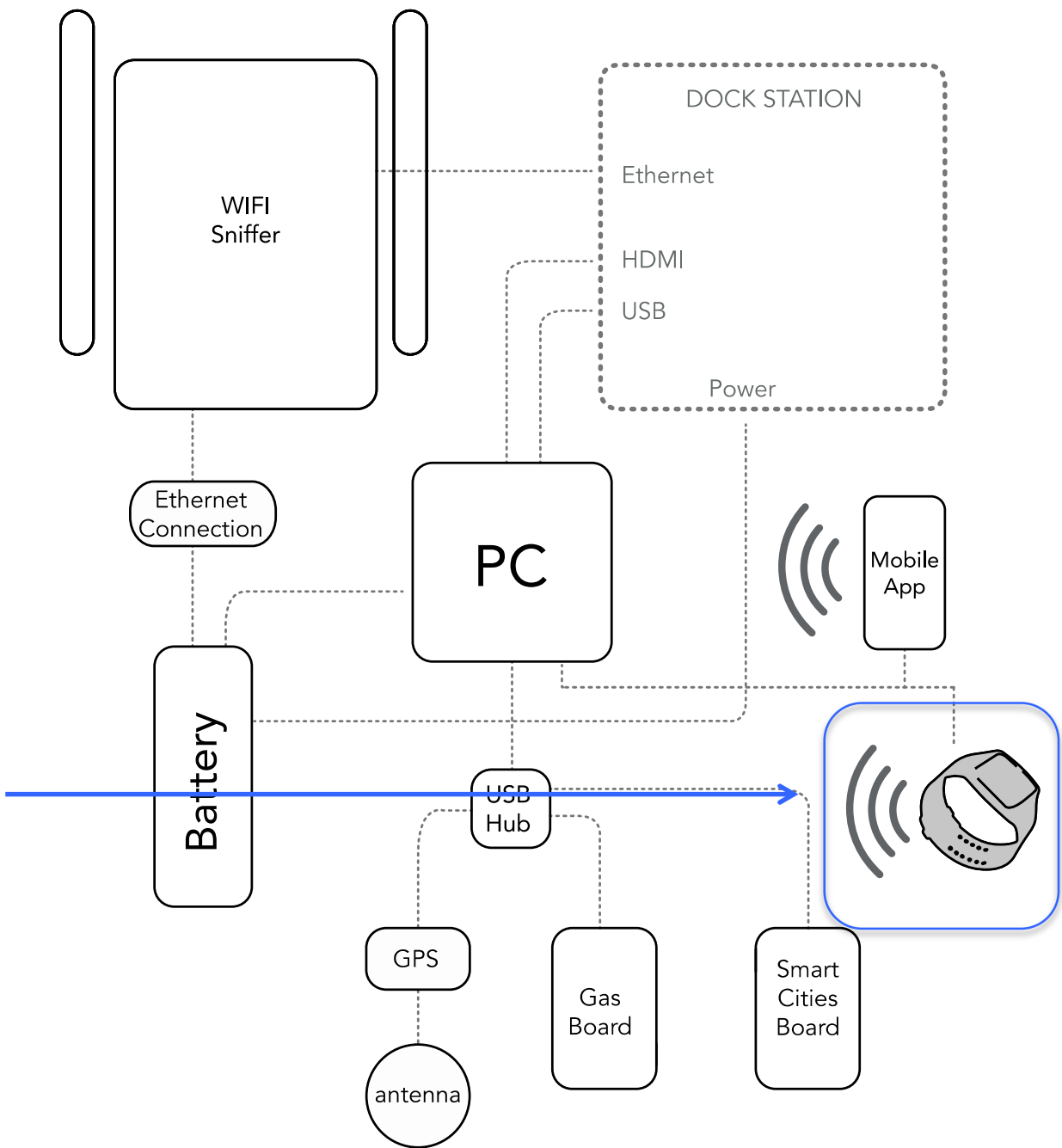
EDA Sensor (GSR Sensor)
Electrodermal Activity Sensor - Used to measure sympathetic nervous system arousal and to derive features related to stress, engagement, and excitement.



Infrared Thermopile
Reads peripheral skin temperature

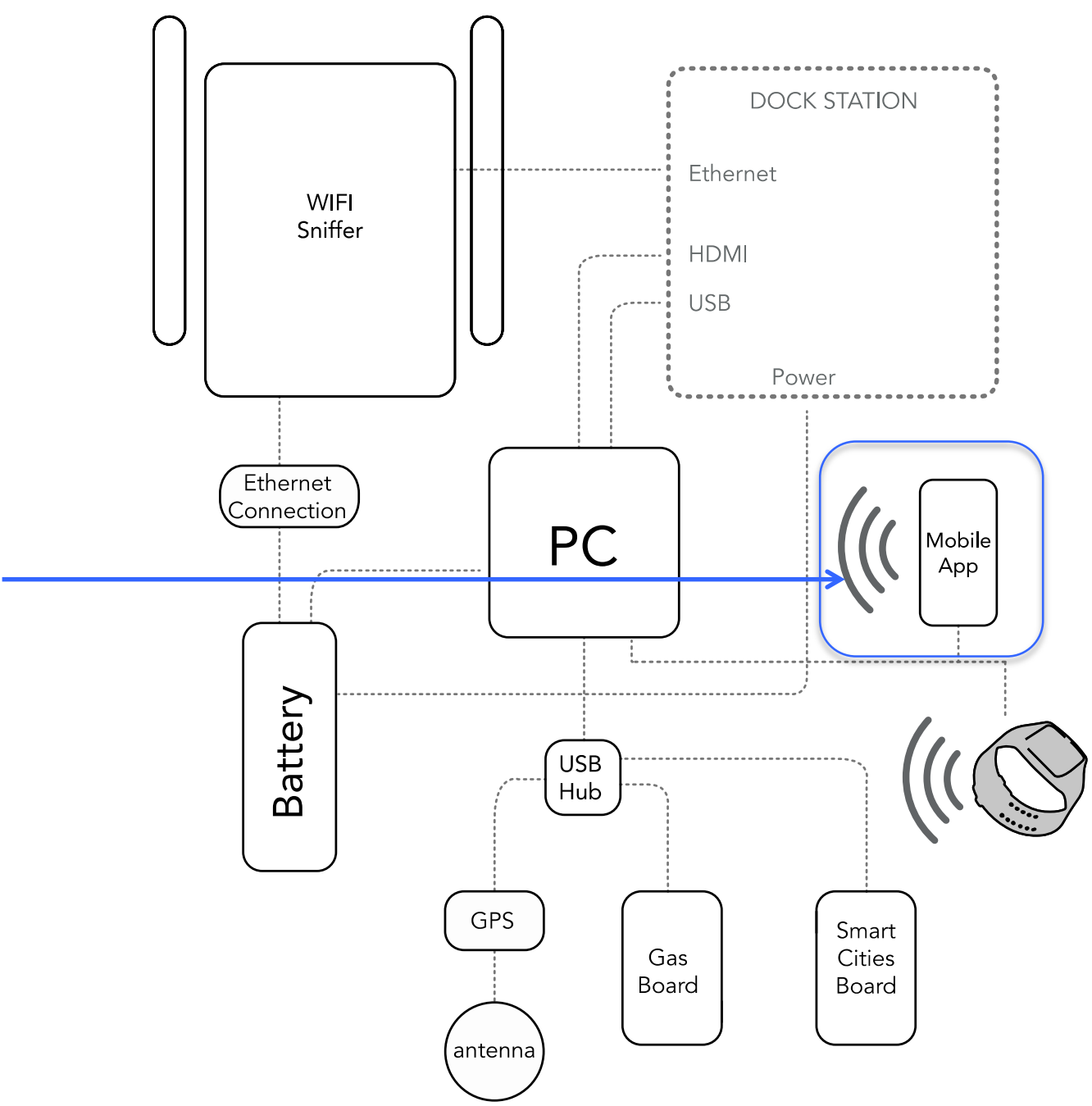
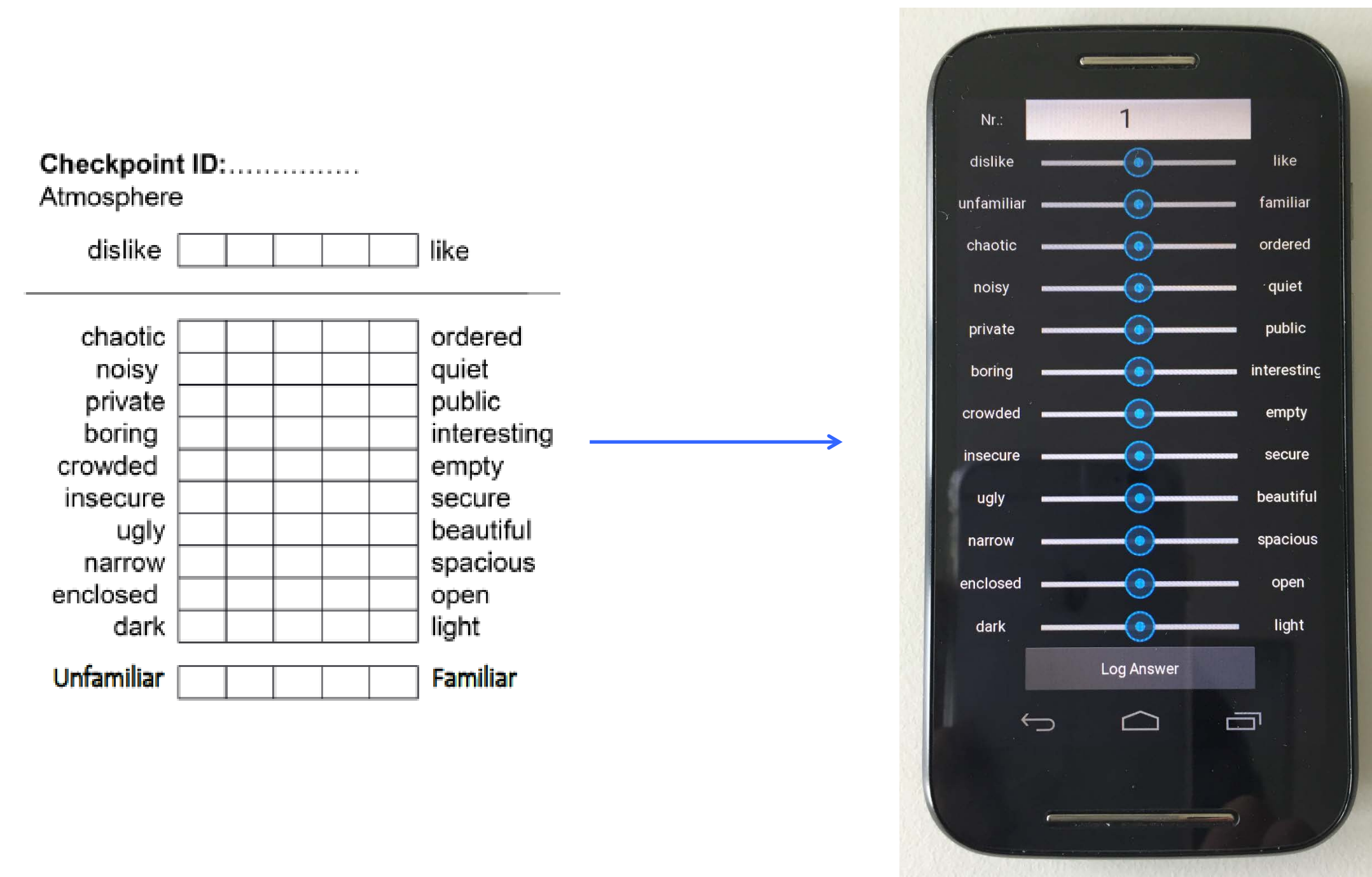


Internal Real-Time Clock
Temporal resolution up to 0.2 seconds in streaming mode



Mobile Sensor equipment

Biofeedback wristband



Experimental data-set

Environmental features



Noise (soundscape) sensor: Smart Cities Board
Manufacturer: Libelium Comunicaciones Distribuidas S.L., Zaragoza (Spain)
Measurement range: 50-100 dB
Frequency: 0.4 Hz
Accuracy: ± 2.5 dB
Physical meaning: [Visit](#)



Dust sensor: Smart Cities Board
Manufacturer: Libelium Comunicaciones Distribuidas S.L., Zaragoza (Spain)
Measurement range: 0.5V/(0.1 mg/m³)
Frequency: 0.4 Hz
Accuracy: Operating supply voltage 5 \pm 0.5V



Temperature sensor: HOBO U12 Logger
Manufacturer: Onset Computer Corporation, Bourne, MA, USA
Measurement range: -20°C -70°C
Frequency: 1 Hz
Accuracy: $\pm 0.35^\circ\text{C}$ from 0°C to 50°C



Relative Humidity Sensor: HOBO U12 Logger
Manufacturer: Onset Computer Corporation, Bourne, MA, USA
Measurement range: 5%-95% (non-condensing)
Measurement unit: %
Frequency: 1 Hz
Accuracy: : $\pm 2.5\%$ typical, $\pm 3.5\%$ max



Illuminance sensor: HOBO U12 Logger
Manufacturer: Onset Computer Corporation, Bourne, MA, USA
Measurement range: 10-32,200 lux
Frequency: 1 Hz
Physical meaning: [Visit](#)

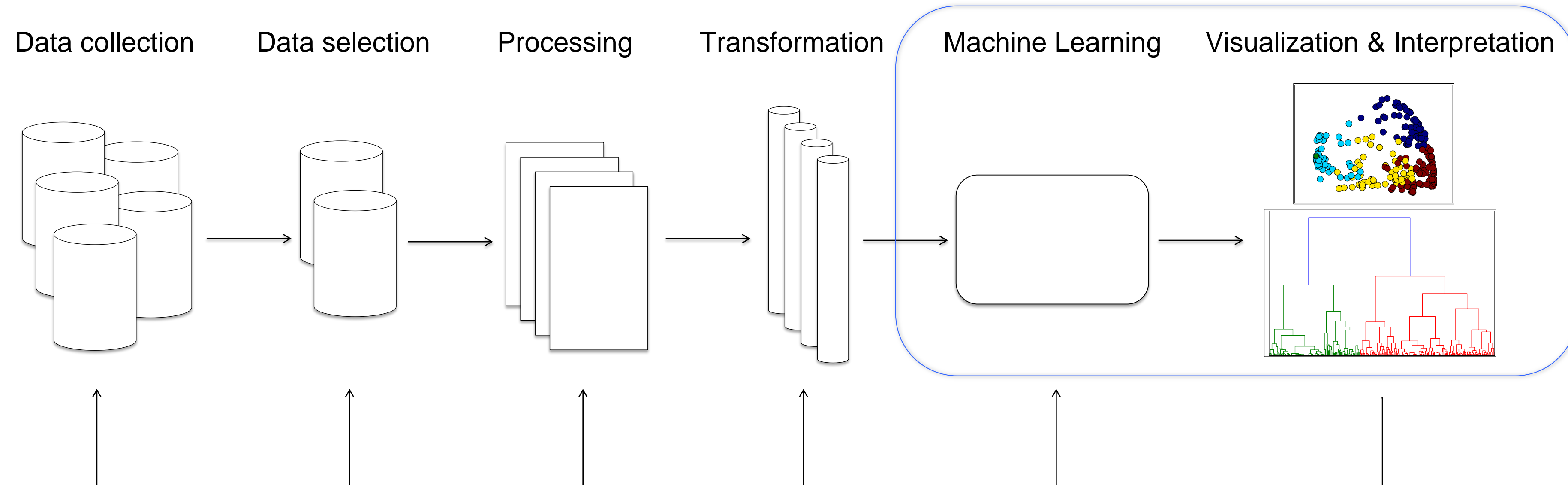


Global Navigation Satellite System (GNSS)/ Global Positioning System (GPS): Ducat 10
uBlox NEO-M8N, TW2410 Antenna
Manufacturer: Tallysman Wireless Inc.
Measurement: WGS84 spherical coordinates
Frequency: 1Hz
Accuracy: [Visit](#)

Background

How can data mining be creative?

What do we want to know?

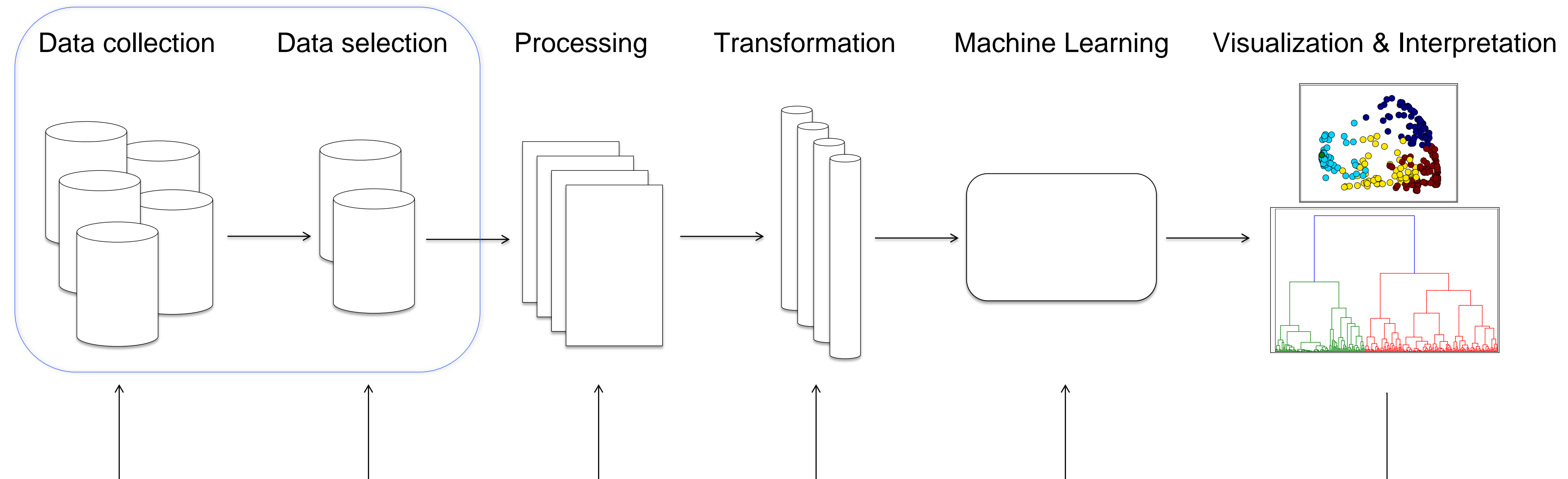


Typical Knowledge Discovery Diagram (KDD)

Background

How can data mining be creative?

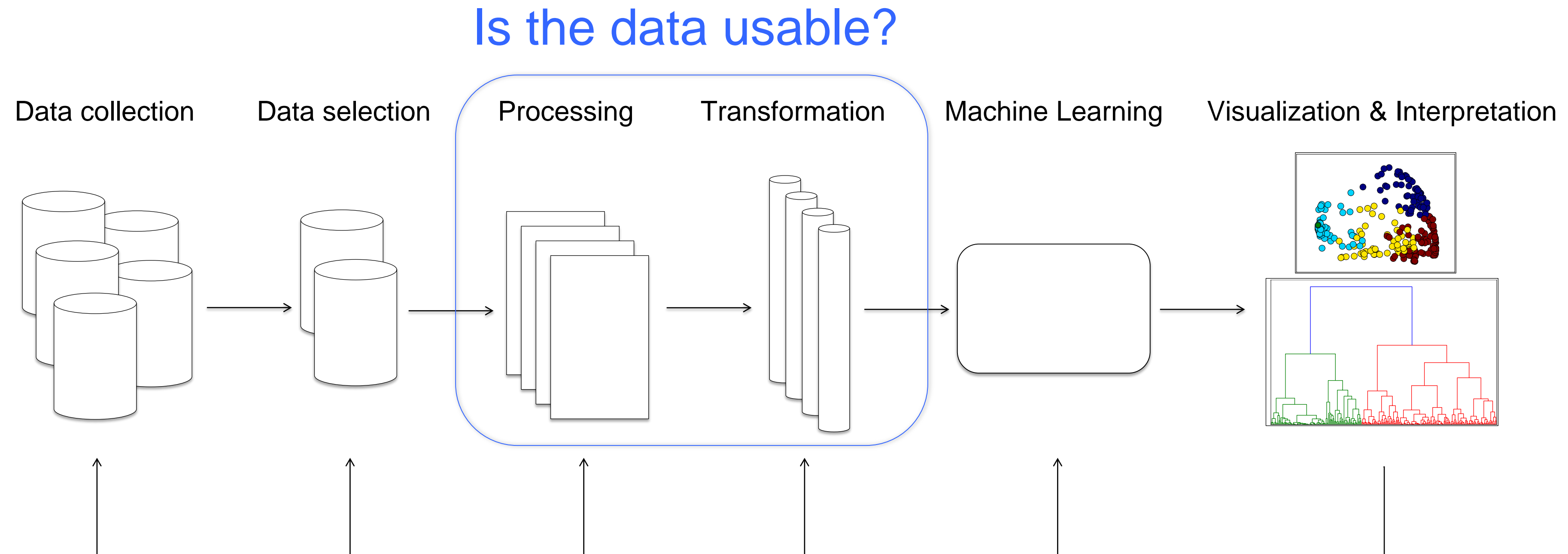
Domain specific data source(s)



Typical Knowledge Discovery Diagram (KDD)

Background

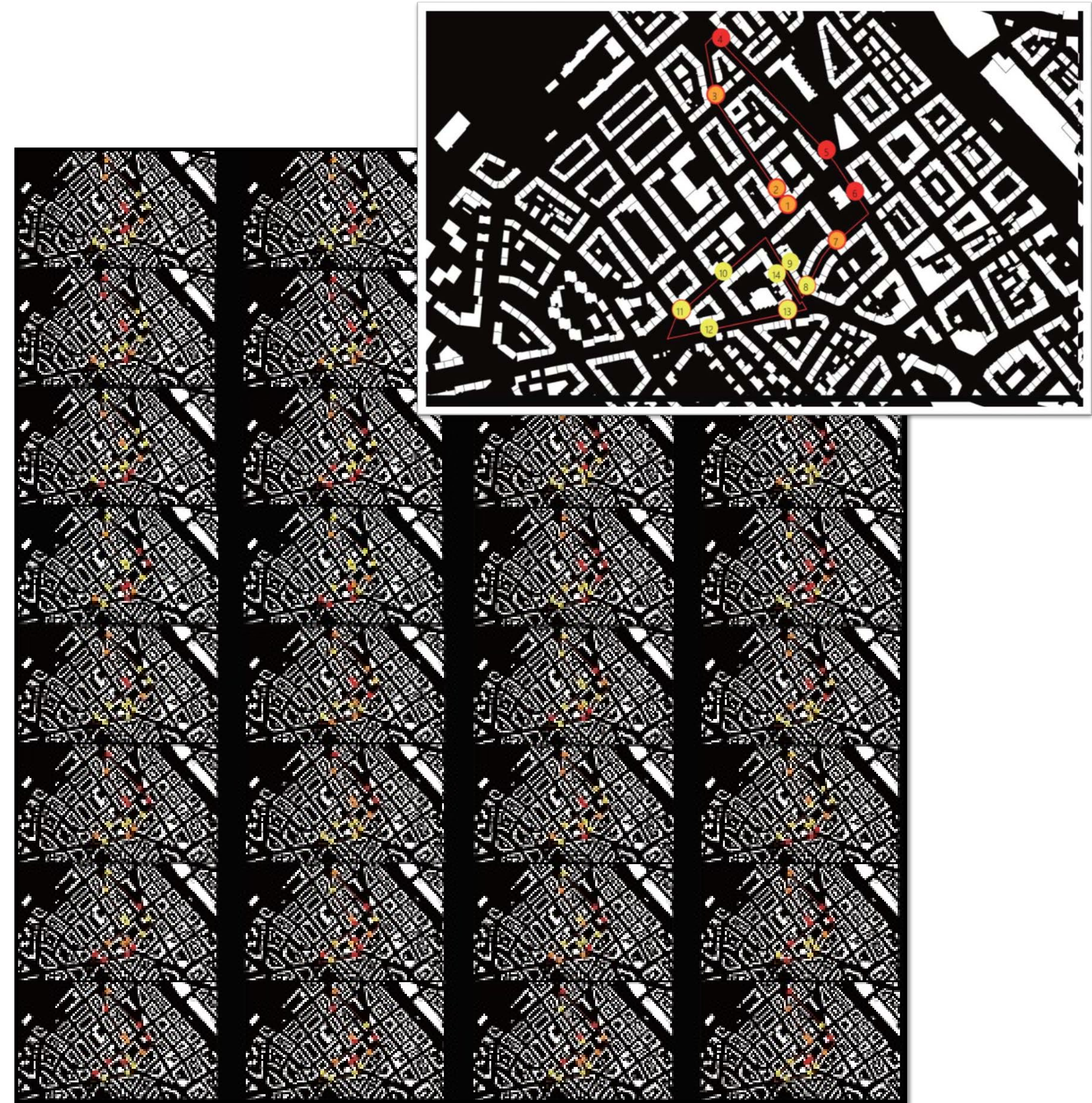
The not-so creative, but essential part of data mining



Typical Knowledge Discovery Diagram (KDD)

Final Project Description

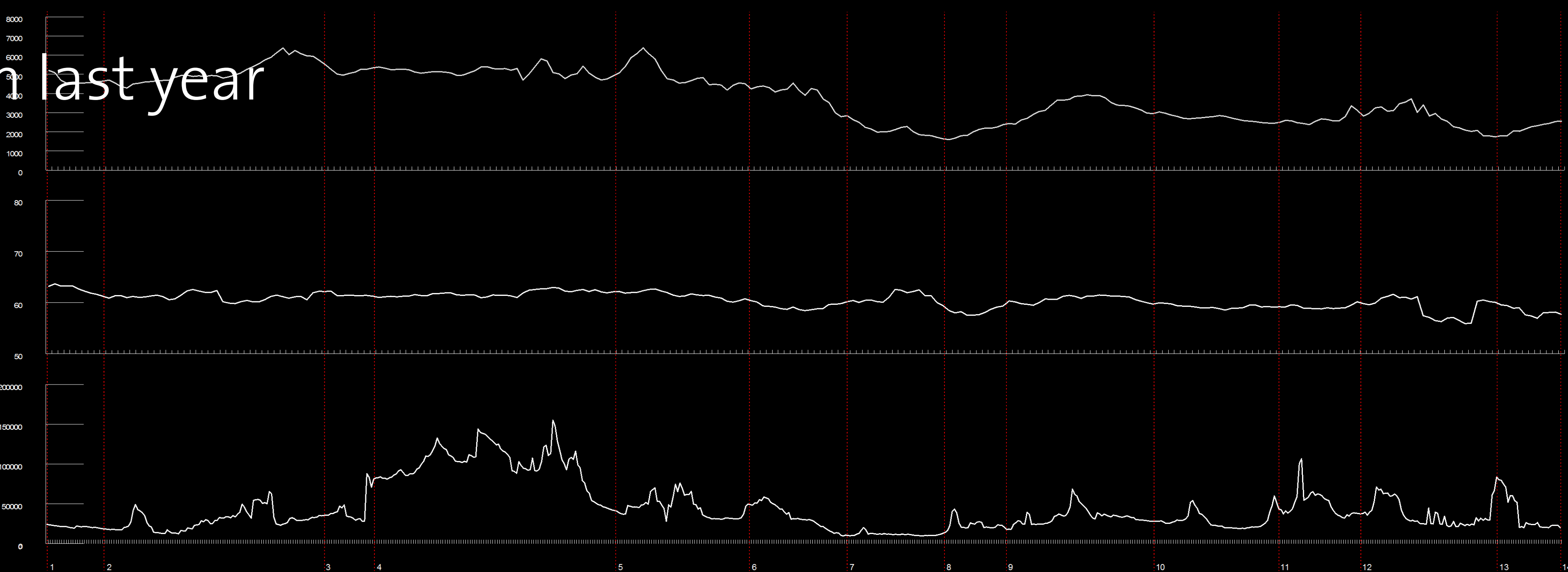
1. Formulate 1-2 specific question(s) of interest to you
2. State your hypothesis/expected outcome based on supporting literature (minimum one source) your expertise, and intuition
3. Answer that question through your analysis, for this:
 - Select the best available data sources for your question (min. of 2 data sources)
 - Include at least one supervised learning or unsupervised learning technique
4. Summarize your results
 - Does your analysis answer your question(s)?
5. Conclusions & lessons learned
 - Fore example if you had more time, data, resources, etc. how would you improve your study
6. Include motivation and references



Images: Final project CDM Spring 2017 by Biyu Wang & Jiani Liu



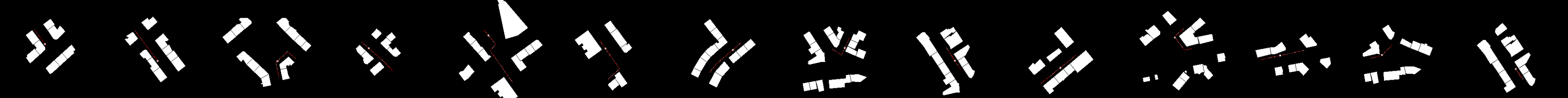
One example from last year



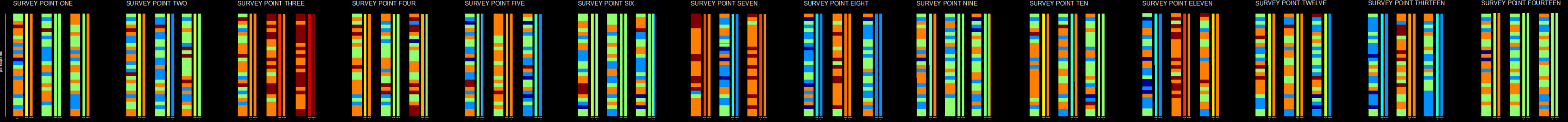
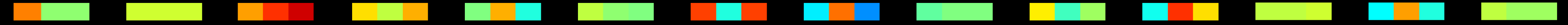
SURVEY POINT ONE SURVEY POINT TWO SURVEY POINT THREE SURVEY POINT FOUR SURVEY POINT FIVE SURVEY POINT SIX SURVEY POINT SEVEN SURVEY POINT EIGHT SURVEY POINT NINE SURVEY POINT TEN SURVEY POINT ELEVEN SURVEY POINT TWELVE SURVEY POINT THIRTEEN SURVEY POINT FOURTEEN



SURVEY POINT TYPOLOGIES

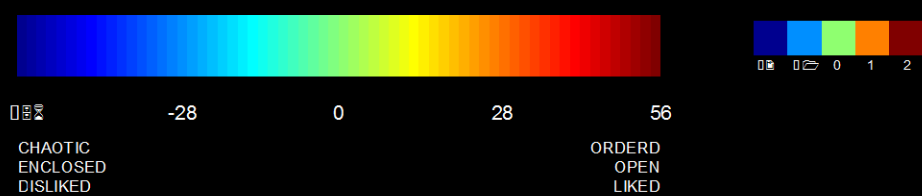


STRUCTURE | SPACIOUSNESS | PREFERENCE



DARCH

ia Chair of Information Architecture



Final project from Jochen Aarts and Stéphane de Weck Creative Data Mining FS2016

What is programming?

- Programming is about solving problems and puzzles. You describe a precise set of instructions which the computer follows exactly.

For example if we want to sort a list:

```
list=[23, 44, 5, 17, 8, 90, 102]  
list.sort()  
print(list)
```

```
output: [5, 8, 17, 23, 44, 90, 102]
```

Types and Values

- *Values* are for example: 5, 12.6, True or “Hello world”
- Each value is of a certain *type*
 - Numbers are: → numeric (integer or float)
 - True is: → boolean (only two values or states)
 - “Hello world” is: → a string (list of characters)

Operators

- Operators represent a value manipulation (+, -, /, %,)
- It is important to check that the data type is correct, otherwise you might get an unexpected result.

```
p=int(17/9)  
print(p)
```

Output :

1

```
p=round(17/9)  
print(p)
```

2

```
p=float(17/9)  
print(p)
```

1.888888888888

Data Structures

There are three basic data structures in Python: Example:

1. Lists- A sequence of values of any type

```
CDM=[ 'There are', [20, 17, 11, 7], 'students  
registered and', 2, 'instructors' ]
```

2. Lists provide several functionalities:

```
CDM[2]           #Accesses the third value  
CDM.remove(2)    #Removes the first occurrence of 2  
CDM.append(2)    #Adds 2 to the end of the list
```


Data Structures

There are three basic data structures in Python:

2. Tuples- a sequence of values of any type, however cannot be altered after it is initiated.

Example:

```
CDM= ( 'There are', [20, 17, 11, 7], 'students registered and', 2, 'instructors' )
```

Vs.

```
CDM= [ 'There are', [20, 17, 11, 7], 'students registered and', 2, 'instructors' ]
```

Data Structures

There are three basic data structures in Python:

3. Dictionaries- store data as key-value pairs.
Each key has an associated value.

Example

```
CDM={ 'Number of students so far': [20, 17, 11, 7]}
```

```
CDM[ 'Number of students so far' ]
```

```
# Returns [20,17,11,7]
```

```
CDM[ 'final number' ]='usually 1/3 of first day'
```

```
#Adds a new element with key 'final number' and value 'usually 1/3 of  
the first day'
```


Logical Operators

Logical Operators are expressions that evaluate to either *True* or *False*.

They are normally used in conditional statements.

Relational operators: `==`, `!=`, `<`, `>`, `<=`, `>=`, `and`, `or`, `not`

They stand for: equal, not equal, smaller, greater, smaller equal, greater equal, and, or, not

Examples:

`x < 2 or x >= 0`

`x == 0 or x != 0`

`not (x == 0 or x > 100)`

Functions

If you do not want to retype the same code over and over again, you can define a function that you can call over and over again

Structure:

```
def functionName(parameters):  
    instructions  
    return variable
```

Example:

```
def square(a):  
    a = a * a  
    return a
```

To get the square of 6, now just write:

```
square(6)
```


Built-in python Functions

The Python interpreter has a number of functions and types built into it that are always available. They are listed here in alphabetical order.

Built-in Functions				
<code>abs()</code>	<code>dict()</code>	<code>help()</code>	<code>min()</code>	<code>setattr()</code>
<code>all()</code>	<code>dir()</code>	<code>hex()</code>	<code>next()</code>	<code>slice()</code>
<code>any()</code>	<code>divmod()</code>	<code>id()</code>	<code>object()</code>	<code>sorted()</code>
<code>ascii()</code>	<code>enumerate()</code>	<code>input()</code>	<code>oct()</code>	<code>staticmethod()</code>
<code>bin()</code>	<code>eval()</code>	<code>int()</code>	<code>open()</code>	<code>str()</code>
<code>bool()</code>	<code>exec()</code>	<code>isinstance()</code>	<code>ord()</code>	<code>sum()</code>
<code>bytearray()</code>	<code>filter()</code>	<code>issubclass()</code>	<code>pow()</code>	<code>super()</code>
<code>bytes()</code>	<code>float()</code>	<code>iter()</code>	<code>print()</code>	<code>tuple()</code>
<code>callable()</code>	<code>format()</code>	<code>len()</code>	<code>property()</code>	<code>type()</code>
<code>chr()</code>	<code>frozenset()</code>	<code>list()</code>	<code>range()</code>	<code>vars()</code>
<code>classmethod()</code>	<code>getattr()</code>	<code>locals()</code>	<code>repr()</code>	<code>zip()</code>
<code>compile()</code>	<code>globals()</code>	<code>map()</code>	<code>reversed()</code>	<code>__import__()</code>
<code>complex()</code>	<code>hasattr()</code>	<code>max()</code>	<code>round()</code>	
<code>delattr()</code>	<code>hash()</code>	<code>memoryview()</code>	<code>set()</code>	

<https://docs.python.org/3/library/functions.html>

Modules

- Often, more complicated code is already implemented and can be used as modules. They help to be more efficient, because the programmer does not need to implement everything from scratch.
- When installed, modules can be integrated in the code using the keyword *import*.

Example:

```
import pandas as pd  
# imports the Pandas module, and now you can access the built in functions
```


Working with Files

Reading in a file using pandas:

```
d=pd.read_csv('filename.csv',
header=None)
```

Creating dataframes:

```
df=pd.DataFrame(data=d)
```

Appending column names:

```
columnNames=[ 'Date-Time' , 'path-point' , 'like-dislike' , 'familiar-unfamiliar' , 'ordered-chaotic' , 'quiet-noisy' , 'public-private' ,
'interesting-boring' , 'empty-crowded' , 'secure-insecure' , 'beautiful-ugly' , 'spacious-narrow' , 'open-enclosed' , 'light-dark' ]
```

```
df.columns = columnNames
```

23:30.4	1	1	-2	1	-2	-1	-1	0	2	-1	2	1	1
24:57.0	2	1	-2	1	-1	-1	1	1	2	0	1	1	1
29:03.1	3	2	-2	-1	2	1	1	-1	2	1	2	2	1
30:37.1	4	2	-1	1	0	1	0	2	2	1	2	2	1
34:24.9	5	-1	-1	-1	-2	2	1	-1	2	-1	1	2	0
36:37.1	6	1	0	2	-1	1	-1	1	2	1	1	1	1
38:38.9	7	1	-1	-1	-1	1	1	1	2	1	0	1	1
40:37.8	8	-1	1	-1	-2	1	-1	-1	2	-1	1	1	2
41:44.5	9	1	-1	1	-1	1	1	-1	2	1	-1	1	1
44:48.4	10	1	-1	1	0	1	-1	-1	2	-1	1	1	1
46:38.4	11	1	1	-2	-1	1	-1	-1	2	-1	1	1	0
48:04.0	12	1	-1	1	-1	2	1	-1	2	0	2	1	0
50:28.5	13	-1	1	-2	-2	2	0	-2	2	-1	1	1	1
51:56.2	14	1	2	1	-1	2	-1	1	2	2	2	1	1

Manipulating dataframes

For Example, from the previous dataset apply the Transpose Attribute:

```
transpose= df.T[1:14]
print(transpose)
```

Writing files into a folder:

```
import os
#hint the os module will also help with your hw

os.makedirs(r'./Single-survey/') #creates the directory
newpath = r'./Single-survey/'    #creates access

Transpose.to_csv(newpath + `transpose.csv` )
#now look in your working directory to find the new csv file
```

Attributes	
<code>T</code>	Transpose index and columns
<code>at</code>	Fast label-based scalar accessor
<code>axes</code>	Return a list with the row axis labels and column axis labels as the only members.
<code>blocks</code>	Internal property, property synonym for <code>as_blocks()</code>
<code>dtypes</code>	Return the dtypes in this object.
<code>empty</code>	True if NDFrame is entirely empty [no items], meaning any of the axes are of length 0.
<code>ftypes</code>	Return the ftypes (indication of sparse/dense and dtype) in this object.
<code>iat</code>	Fast integer location scalar accessor.
<code>iloc</code>	Purely integer-location based indexing for selection by position.
<code>is_copy</code>	
<code>ix</code>	A primarily label-location based indexer, with integer position fallback.
<code>loc</code>	Purely label-location based indexer for selection by label.
<code>ndim</code>	Number of axes / array dimensions
<code>shape</code>	Return a tuple representing the dimensionality of the DataFrame.
<code>size</code>	number of elements in the NDFrame
<code>style</code>	Property returning a Styler object containing methods for building a styled HTML representation for the DataFrame.
<code>values</code>	Numpy representation of NDFrame
Methods	
<code>abs()</code>	Return an object with absolute value taken—only applicable to objects that are all numeric.
<code>add(other[, axis, level, fill_value])</code>	Addition of dataframe and other, element-wise (binary operator <i>add</i>).
<code>add_prefix(prefix)</code>	Concatenate prefix string with panel items names.
<code>add_suffix(suffix)</code>	Concatenate suffix string with panel items names.
<code>agg(func[, axis])</code>	Aggregate using callable, string, dict, or list of string/callables
<code>aggregate(func[, axis])</code>	Aggregate using callable, string, dict, or list of string/callables

<https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.html>

Programming assignment

Data Processing: data transformation

Transform the survey data from each ESUM participant by survey question:

- The output from each participant's survey responses is a single dataframe (csv file); each participant answered 12 questions at each of the 14 check points
- Your task is to create a program which automatically generates a single dataframe (csv file) for each of the 12 questions; you should have 12 csv files (32 participants by 14 check points with the populated survey results)
- Write the files to a single folder and include this with the python file in a zip folder and email it to us by 9am October 30th

Conceptual assignment

Research other examples of urban data mining and make 2 slides about the most interesting project/application/research group(s) that you find. This will be presented at the beginning of next lecture

- Please send your slides to us by 9am on October 30th

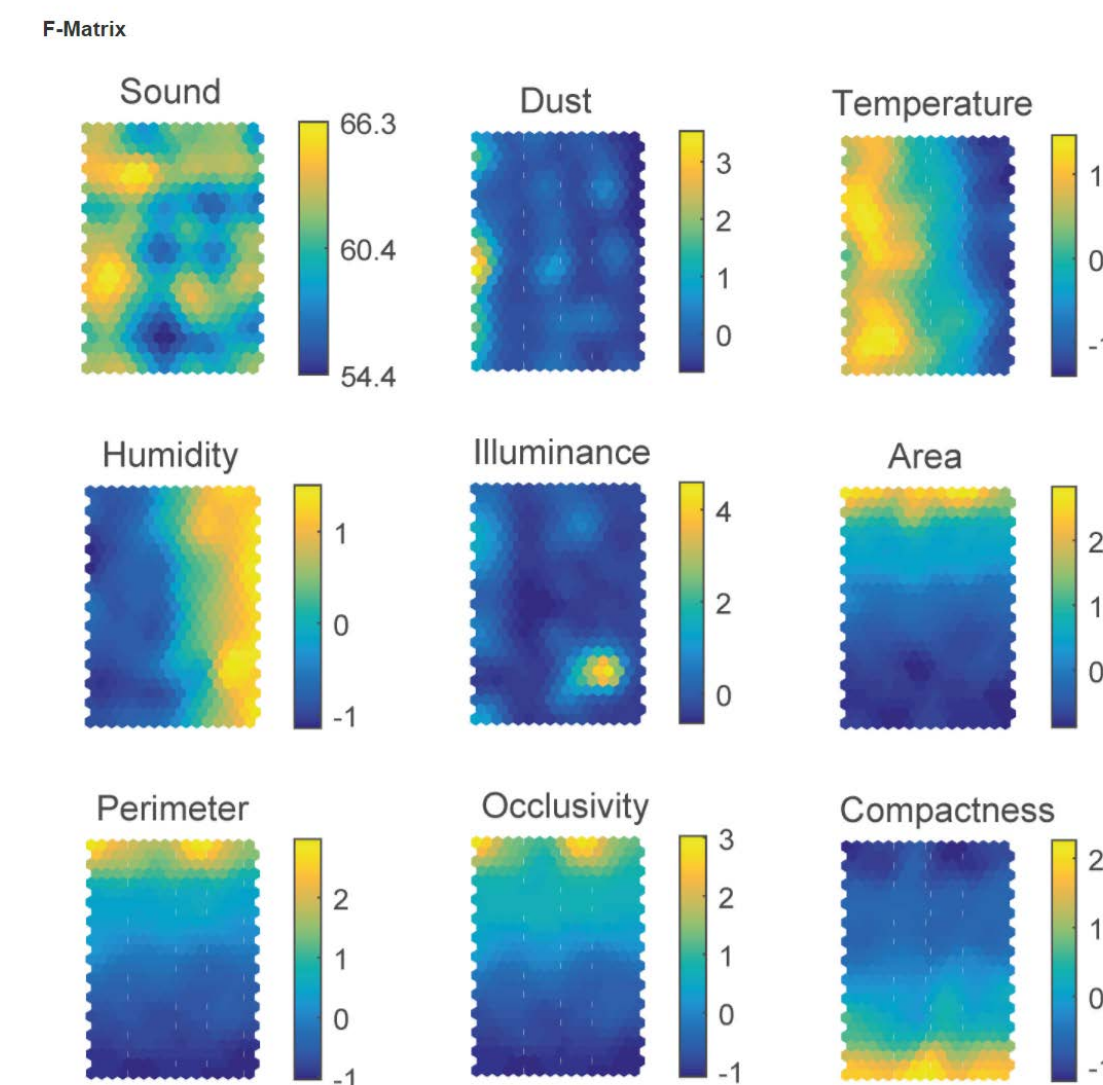


Fig Each feature was linearly scaled with a variance of one so that they have equal importance in computing distance and influence in clustering on the map.

Thank you!

Varun Ojha ojha@arch.ethz.ch

Danielle Griego griego@arch.ethz.ch