

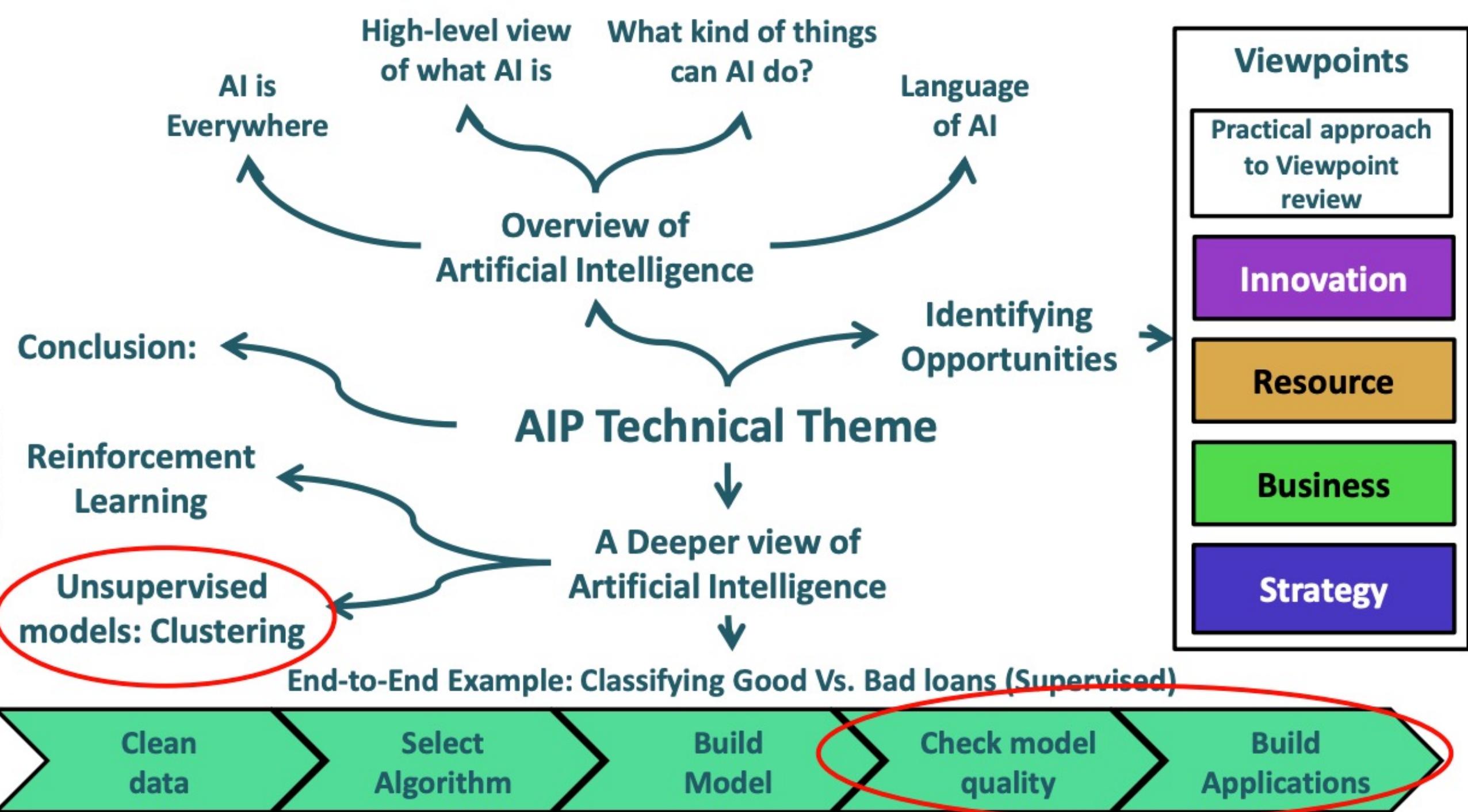
Unsupervised Learning

Elizabeth Savochkina | 7th June

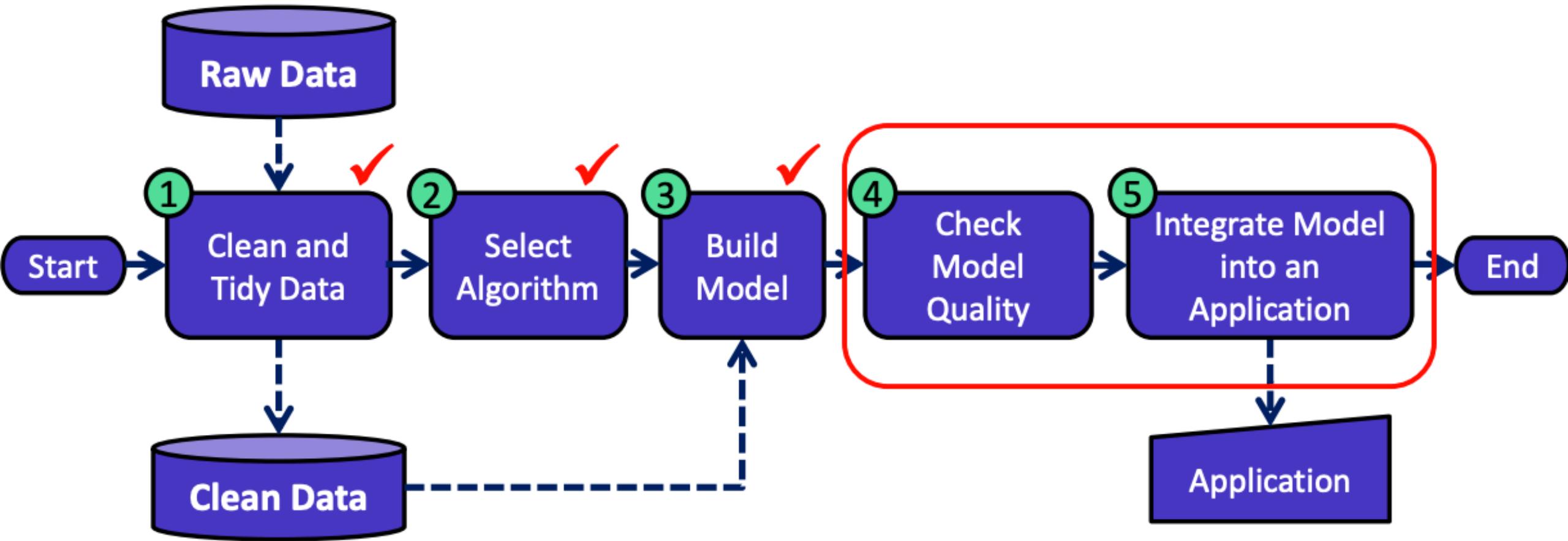


Shape of a Day

Registration 9.00am to 9.30am
Unsupervised Learning 9.30am to 11.30am
Break 11.30am to 13.30 pm
Project Discussion: Opportunity Identification & Ideation 13.30am to 15.30pm



Typical Workflow



Step 4 : Check Model Quality



Evaluate Model Quality

How can you be sure
you have a 'good'
model?

i.e. a model that
gives predictions that
are often correct?

Shoe Size	Hat Size	Body Weight	Supports Bayern Munich	Supports Manchester City	Supports Paris Saint-Germain	Likes pies for lunch	Likes Cats	Likes Dogs	Loan Default
35	30	102	0	0	1	0	0	1	1
41	23	106	0	0	0	0	0	1	0
37	18	136	0	0	0	0	1	0	1
41	16	182	0	0	0	1	0	0	0
44	22	19					0	0	1
42	26	7					0	0	1
36	30	12					1	1	1
36	13	150	0	0	0	1	0	1	1
45	10	105	0	0	0	1	1	0	0
32	27	134	0	0	0	0	1	0	0
37	22	137	0	0	0	1	1	1	1
36	13	149	0	0	0	1	0	1	0

Garbage
in

Garbage
Out



Different ways of being Wrong

The 'confusion matrix' is a useful way of reviewing model quality for category predictions

		Prediction	
		No Default	Default
Reality	No Default	Computer said "They will not default". We gave them a loan. They <u>did not</u> default. Good business!	Computer said "They will default". We didn't give them a loan. Lost business. Not good
	Default	Computer said "They will not default". We gave them a loan. They <u>did</u> default. Bad!	Computer said "They will default". We did not give them a loan. They <u>would</u> have defaulted in reality Good decision

Reviewing the Confusion Matrix

4

Check Model Quality

- In this case, a total of 155 predictions were correct (78%)
- A total of 45 were wrong ..
- Of which 15 were bad failures

		Prediction		
		No Default	Default	
Reality	No Default	75	30	105
	Default	15	80	95
		90	110	200

Model Quality Check : Practical Example



Check Model Quality

jupyter German Credit V2 (no Overplot across age (julian))

File Edit View Insert Cell Kernel Help Trusted Python 3

Step 4 Check Model Quality

In practice, before we build our model, we should split the data into two parts. One part will be used for building the model, the other part will be used to test how good it is at making predictions.

First we import a new library function that allows us to randomly split the data.

```
In [32]: from sklearn.model_selection import train_test_split
```

We use this function to split the data into two parts in the ratio 80%:20% as we don't need so much data for testing.

```
In [33]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

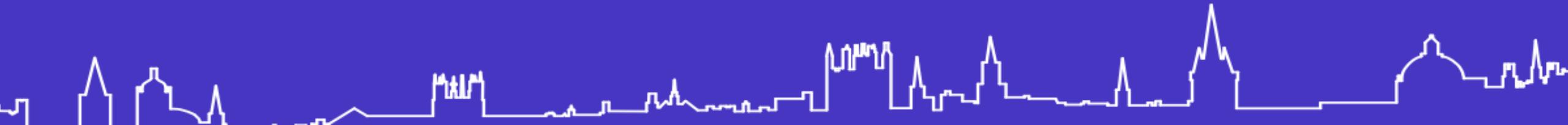
As above, we build a model, but this time based only on the training data.

```
In [34]: logModel = LogisticRegression(C=1000000, p=1)
```

Then generate a set of predictions based on the test data.

Demonstration: “04-01 German Credit Model Quality Check.mp4”

Step 5 : Building an Application



Options for Building Applications

5

Integrate Model
into an Application

- Building models often requires large amounts of computing power
- But running them as applications may need far fewer resources ..
 - Enabling ML apps to run on smartphones and laptop computers



Run the model directly on a desktop computer or smartphone



Run the model on a server – and access remotely



Building Models into Applications

5

Integrate Model
into an Application

- Regression models are ‘formulae’ that can be copied into software code and evaluated
- Decision trees might be implemented as a series of ‘If .. Then .. Else’ statements
- Library code is available to evaluate Neural Networks on smaller processors

$$Y = k + A \cdot x_1 + B \cdot x_2 + C \cdot x_3 \dots Z \cdot x_n$$

Variable Name	Constant	Number of Rooms	Number of Bathrooms	Garden area	Detached	Semi-detached	Terrace
Model Parameters	50,000	20000	3760	4.67	47577	32432	23543
Applicant data	1	6	2	1200	1	0	0
Result	50000	120000	7520	5604	47577	0	0

House predicted value

£230,701

‘linear’ (straight-line) model built by multiplying and adding terms



Build the Application

5

Integrate Model
into an Application

- In this case the model is simply a list of coefficients for a relatively simple equation ..

```
print(logModel.coef_)\nprint(logModel.intercept_)\n\n[[ 1.78294369e-01 -9.56571126e-01 -2.10532764e-01 -4.67743993e-01\n-5.32855166e-01 -8.92711403e-02 -1.48284746e-02 -4.13219760e-02\n 2.09198704e-01 -2.92896634e-01 -5.94578827e-01 -1.15009343e-02\n-4.76317246e-02 -6.19144098e-01 -4.00699728e-01 -1.99907975e-01\n-3.53242688e-01  2.75573634e-01 -1.03105672e-01  6.34713643e-02\n-6.38642449e-01  1.75219375e-01 -1.83209646e-01 -4.74459892e-01\n-1.95826594e-01  3.47136666e-01  1.02540582e-01 -9.87155603e-02\n-6.89921888e-01 -3.40316556e-01  3.29246259e-01 -1.40511096e-01\n-2.91350624e-03 -4.58538784e-01 -4.05559630e-01  5.88994290e-01\n-1.03965156e+00 -4.37697991e-01 -1.92838329e-01 -1.99360467e-01\n-3.14801370e-01  5.71761336e-01  5.62232677e-01 -3.39359854e-01\n 1.22444517e-01  3.07173187e-01  5.76979945e-01 -2.89642473e-01\n-3.09780167e-01 -9.63007250e-01  5.29671486e-01  1.45249041e-01\n-3.17664937e-01 -1.03553235e+00  3.03348301e-02  1.42773085e-04\n 3.22047587e-01  2.86347603e-02 -1.21770781e-02  1.80129775e-01\n-3.87951340e-02]]\n[-0.67827676]
```

Demonstration: “04-02 - Export model params.mp4”

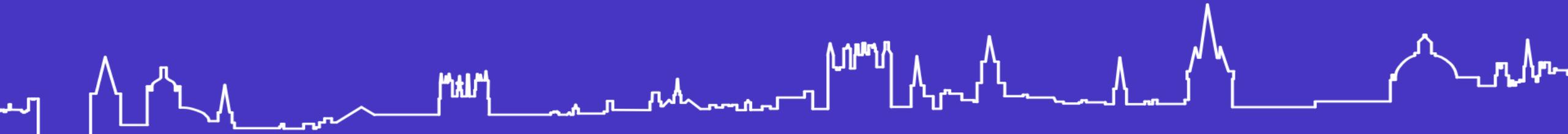
Embed the Model into an Application

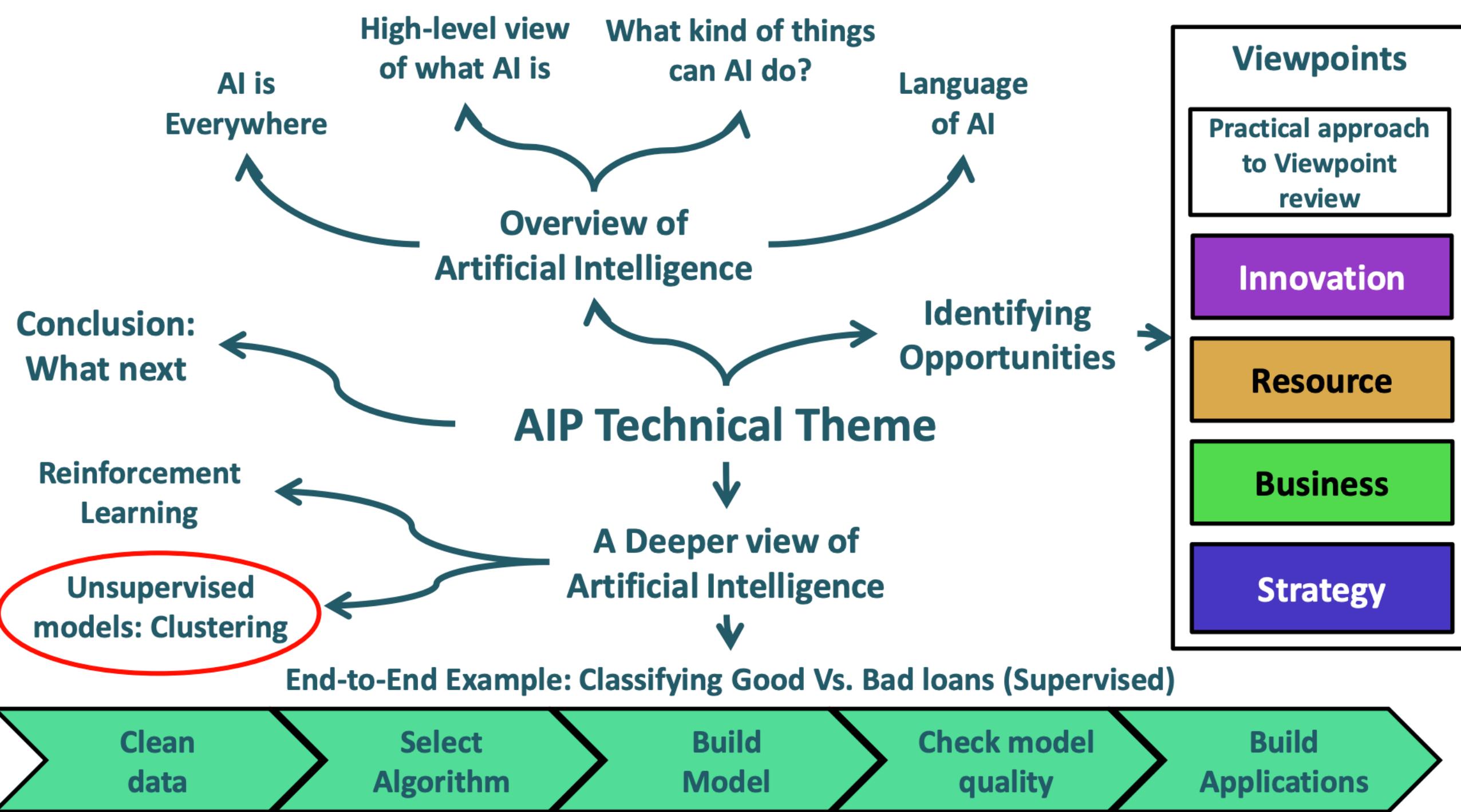
5 Integrate Model
into an Application

- Once you have the model, you can embed it into other software
- In this case, I have reproduced the model in Excel

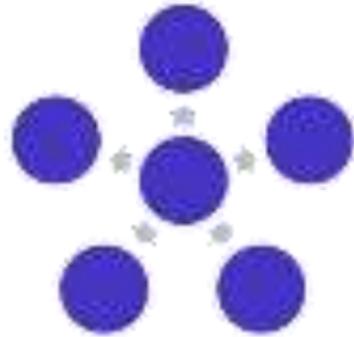
Original Factor Name	After expanding categories	Parameter value	Current Date	(x) (y)
		Constant	0.47627670	1 -0.7
Housing tenure		R111	-0.00000000	0 0.0
		R112	-0.537400	0 0.0
Tenure type		R113	-0.233104	0 -0.3
		R114	-0.288472	0 0.0
		R115	-0.388402	0 0.0
		R116	-0.999402	0 0.0
		R117	-0.999402	0 0.0
Age		R118	-0.888402	0 0.0
		R119	-0.333402	0 0.0
Browsing		R120	0.000000	0 0.0
		R121	0.000000	0 0.0
Other additional planes		R122	0.000000	0 0.0
		R123	0.000000	0 0.0
		R124	0.000000	0 0.0
		R125	0.000000	0 0.0
		R126	0.000000	0 0.0
Property		R127	-0.000000	0 0.0
		R128	-0.000000	0 0.0
		R129	0.234495	0 -0.4
Contact address		R130	0.765492	0 0.0
		R131	0.000000	0 0.0
		R132	0.000000	0 0.0
Marital status		R133	0.758492	0 0.0
		R134	0.000000	0 0.0
		R135	0.000000	0 0.0
Assessment Area		R136	0.000000	0 0.0
		R137	0.000000	0 0.0
		R138	0.000000	0 0.0
		R139	0.000000	0 0.0
		R140	0.000000	0 0.0
		R141	0.000000	0 0.0
		R142	0.000000	0 0.0
		R143	0.000000	0 0.0
		R144	0.000000	0 0.0
		R145	0.000000	0 0.0
		R146	0.000000	0 0.0
		R147	0.000000	0 0.0
		R148	0.000000	0 0.0
		R149	0.000000	0 0.0
		R150	0.000000	0 0.0
		R151	0.000000	0 0.0
		R152	0.000000	0 0.0
		R153	0.000000	0 0.0
		R154	0.000000	0 0.0
		R155	0.000000	0 0.0
		R156	0.000000	0 0.0
		R157	0.000000	0 0.0
		R158	0.000000	0 0.0
		R159	0.000000	0 0.0
		R160	0.000000	0 0.0
		R161	0.000000	0 0.0
		R162	0.000000	0 0.0
		R163	0.000000	0 0.0
		R164	0.000000	0 0.0
		R165	0.000000	0 0.0
		R166	0.000000	0 0.0
		R167	0.000000	0 0.0
		R168	0.000000	0 0.0
		R169	0.000000	0 0.0
		R170	0.000000	0 0.0
		R171	0.000000	0 0.0
		R172	0.000000	0 0.0
		R173	0.000000	0 0.0
		R174	0.000000	0 0.0
		R175	0.000000	0 0.0
		R176	0.000000	0 0.0
		R177	0.000000	0 0.0
		R178	0.000000	0 0.0
		R179	0.000000	0 0.0
		R180	0.000000	0 0.0
		R181	0.000000	0 0.0
		R182	0.000000	0 0.0
		R183	0.000000	0 0.0
		R184	0.000000	0 0.0
		R185	0.000000	0 0.0
		R186	0.000000	0 0.0
		R187	0.000000	0 0.0
		R188	0.000000	0 0.0
		R189	0.000000	0 0.0
		R190	0.000000	0 0.0
		R191	0.000000	0 0.0
		R192	0.000000	0 0.0
		R193	0.000000	0 0.0
		R194	0.000000	0 0.0
		R195	0.000000	0 0.0
		R196	0.000000	0 0.0
		R197	0.000000	0 0.0
		R198	0.000000	0 0.0
		R199	0.000000	0 0.0
		R200	0.000000	0 0.0
		R201	0.000000	0 0.0
		R202	0.000000	0 0.0
		R203	0.000000	0 0.0
		R204	0.000000	0 0.0
		R205	0.000000	0 0.0
		R206	0.000000	0 0.0
		R207	0.000000	0 0.0
		R208	0.000000	0 0.0
		R209	0.000000	0 0.0
		R210	0.000000	0 0.0
		R211	0.000000	0 0.0
		R212	0.000000	0 0.0
		R213	0.000000	0 0.0
		R214	0.000000	0 0.0
		R215	0.000000	0 0.0
		R216	0.000000	0 0.0
		R217	0.000000	0 0.0
		R218	0.000000	0 0.0
		R219	0.000000	0 0.0
		R220	0.000000	0 0.0
		R221	0.000000	0 0.0
		R222	0.000000	0 0.0
		R223	0.000000	0 0.0
		R224	0.000000	0 0.0
		R225	0.000000	0 0.0
		R226	0.000000	0 0.0
		R227	0.000000	0 0.0
		R228	0.000000	0 0.0
		R229	0.000000	0 0.0
		R230	0.000000	0 0.0
		R231	0.000000	0 0.0
		R232	0.000000	0 0.0
		R233	0.000000	0 0.0
		R234	0.000000	0 0.0
		R235	0.000000	0 0.0
		R236	0.000000	0 0.0
		R237	0.000000	0 0.0
		R238	0.000000	0 0.0
		R239	0.000000	0 0.0
		R240	0.000000	0 0.0
		R241	0.000000	0 0.0
		R242	0.000000	0 0.0
		R243	0.000000	0 0.0
		R244	0.000000	0 0.0
		R245	0.000000	0 0.0
		R246	0.000000	0 0.0
		R247	0.000000	0 0.0
		R248	0.000000	0 0.0
		R249	0.000000	0 0.0
		R250	0.000000	0 0.0
		R251	0.000000	0 0.0
		R252	0.000000	0 0.0
		R253	0.000000	0 0.0
		R254	0.000000	0 0.0
		R255	0.000000	0 0.0
		R256	0.000000	0 0.0
		R257	0.000000	0 0.0
		R258	0.000000	0 0.0
		R259	0.000000	0 0.0
		R260	0.000000	0 0.0
		R261	0.000000	0 0.0
		R262	0.000000	0 0.0
		R263	0.000000	0 0.0
		R264	0.000000	0 0.0
		R265	0.000000	0 0.0
		R266	0.000000	0 0.0
		R267	0.000000	0 0.0
		R268	0.000000	0 0.0
		R269	0.000000	0 0.0
		R270	0.000000	0 0.0
		R271	0.000000	0 0.0
		R272	0.000000	0 0.0
		R273	0.000000	0 0.0
		R274	0.000000	0 0.0
		R275	0.000000	0 0.0
		R276	0.000000	0 0.0
		R277	0.000000	0 0.0
		R278	0.000000	0 0.0
		R279	0.000000	0 0.0
		R280	0.000000	0 0.0
		R281	0.000000	0 0.0
		R282	0.000000	0 0.0
		R283	0.000000	0 0.0
		R284	0.000000	0 0.0
		R285	0.000000	0 0.0
		R286	0.000000	0 0.0
		R287	0.000000	0 0.0
		R288	0.000000	0 0.0
		R289	0.000000	0 0.0
		R290	0.000000	0 0.0
		R291	0.000000	0 0.0
		R292	0.000000	0 0.0
		R293	0.000000	0 0.0
		R294	0.000000	0 0.0
		R295	0.000000	0 0.0
		R296	0.000000	0 0.0
		R297	0.000000	0 0.0
		R298	0.000000	0 0.0
		R299	0.000000	0 0.0
		R300	0.000000	0 0.0
		R301	0.000000	0 0.0
		R302	0.000000	0 0.0
		R303	0.000000	0 0.0
		R304	0.000000	0 0.0
		R305	0.000000	0 0.0
		R306	0.000000	0 0.0
		R307	0.000000	0 0.0
		R308	0.000000	0 0.0
		R309	0.000000	0 0.0
		R310	0.000000	0 0.0
		R311	0.000000	0 0.0
		R312	0.000000	0 0.0
		R313	0.000000	0 0.0
		R314	0.000000	0 0.0
		R315	0.000000	0 0.0
		R316	0.000000	0 0.0
		R317	0.000000	0 0.0
		R318	0.000000	0 0.0
		R319	0.000000	0 0.0
		R320	0.000000	0 0.0
		R321	0.000000	0 0.0
		R322	0.000000	0 0.0
		R323	0.000000	0 0.0
		R324	0.000000	0 0.0
		R325	0.000000	0 0.0
		R326	0.000000	0 0.0
		R327	0.000000	0 0.0
		R328	0.000000	0 0.0
		R329	0.000000	0 0.0
		R330	0.000000	0 0.0
		R331	0.000000	0 0.0
		R332	0.000000	0 0.0
		R333	0.000000	0 0.0
		R334	0.000000	0 0.0
		R335	0.000000	0 0.0
		R336	0.000000	0 0.0
		R337	0.000000	0 0.0
		R338	0.000000	0 0.0
		R339	0.000000	0 0.0
		R340	0.000000	0 0.0
		R341	0.000000	0 0.0
		R342	0.000000	0 0.0
		R343	0.000000	0 0.0
		R344	0.000000	0 0.0
		R345	0.000000	0 0.0
		R346	0.000000	0 0.0
		R347	0.000000	0 0.0
		R348	0.000000	0 0.0
		R349	0.000000	0 0.0
		R350	0.000000	0 0.0
		R351	0.000000	0 0.0
		R352	0.000000	0 0.0
		R353	0.000000	0 0.0
		R354	0.000000	0 0.0
		R355	0.000000	0 0.0
		R356	0.000000	0 0.0

Unsupervised Models : Clustering





A Brief Reminder: What is Clustering?



Clustering

**Dividing a large group
into smaller groups
containing members that
are similar to each other**

- Segmenting customers into groups
- Grouping products into families that customers will relate to
- Grouping complaint types



A Simple and Fast Clustering Algorithm

- Here we have a simple example:
 - There are only two dimensions
 - We can easily see two groups
- But the following algorithm, whilst simple, can:
 - Deal with many dimensions ..
 - Divide data into as many groups as we decide
 - Make a decision even if the groups are similar



Step 1 : Take a Guess

- We want to divide this data into two groups which we will call 'A' and 'B'
- First, we select two, completely random points as a guess at the centre of each group

These starting points can be anywhere



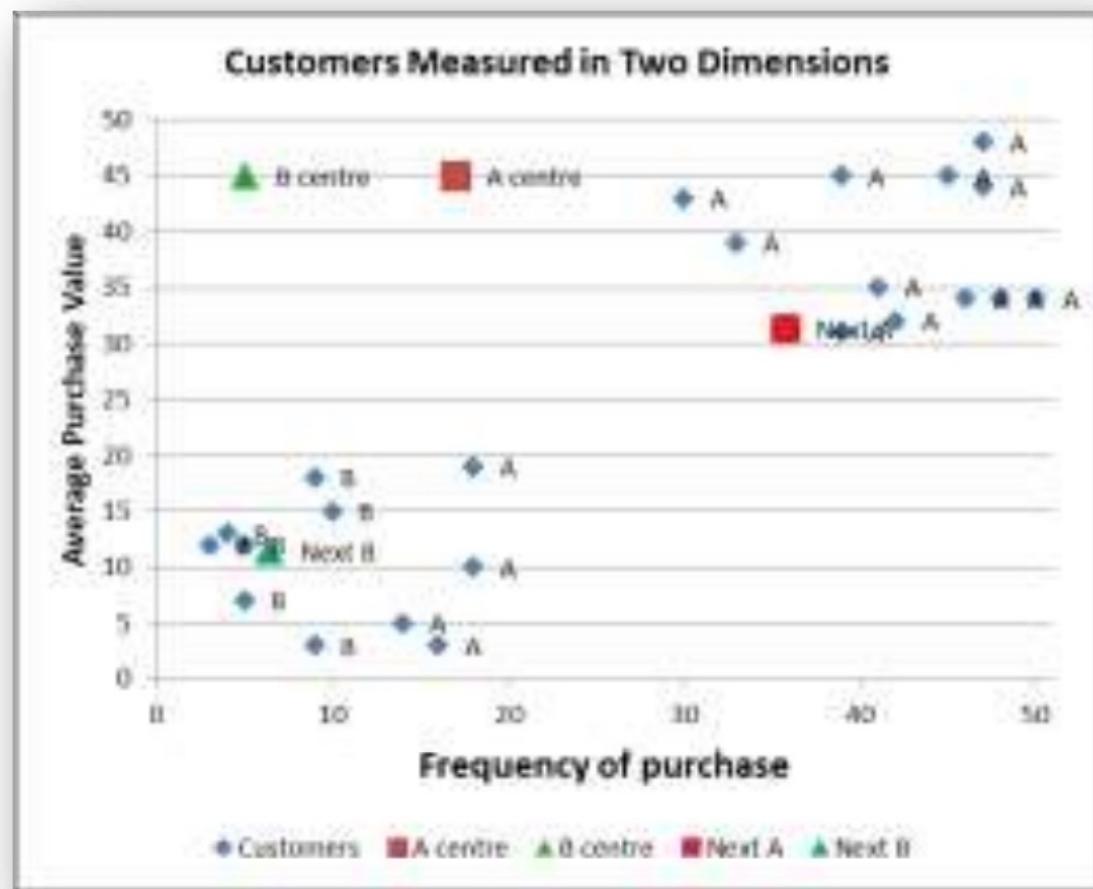
Step 2 : Which Guess is each Point Closest to?

- Measure the distance every point and each of the guesses
- If a point is closest to guess ‘A centre’ – then label it ‘A’
- If the point is closest to guess ‘B centre’ – then label it ‘B’



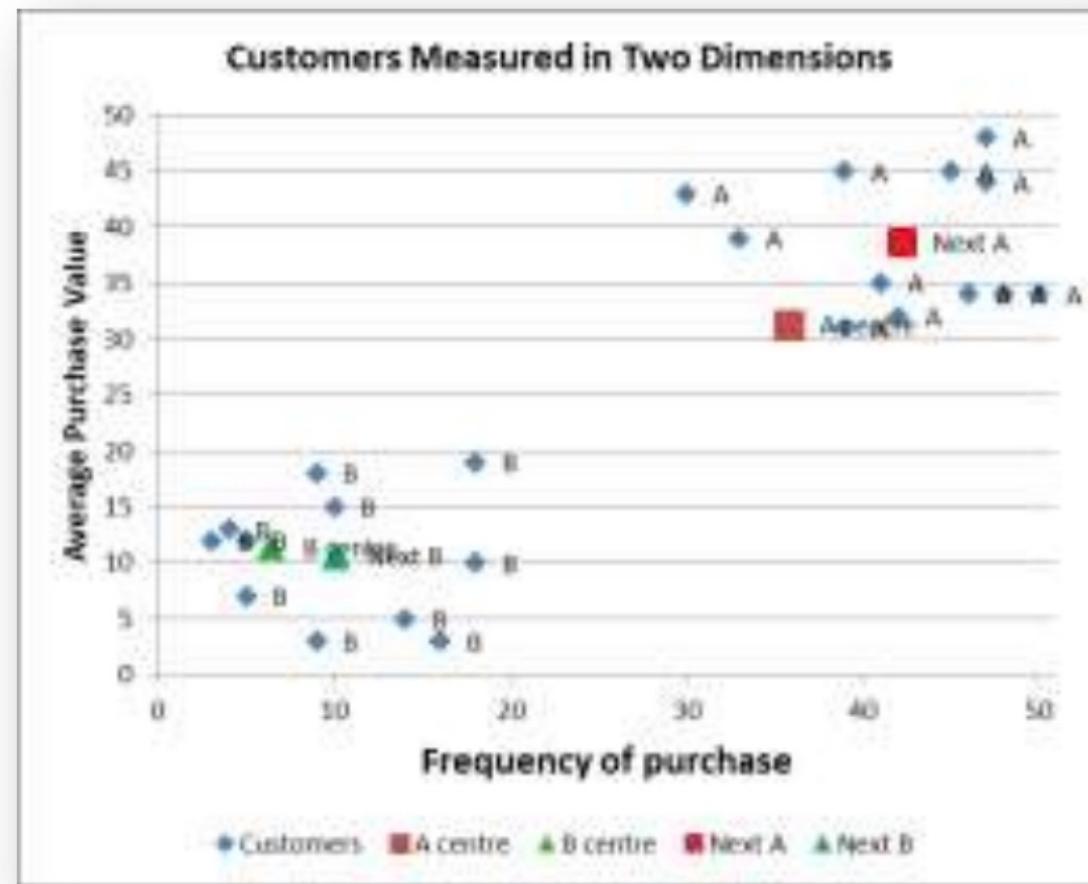
Step 3 : Find the Centre of Each Group

- Find the centre of all of the 'A's
 - Call this 'Next A'
- Find the centre of all of the 'B's
 - Call this 'Next B'



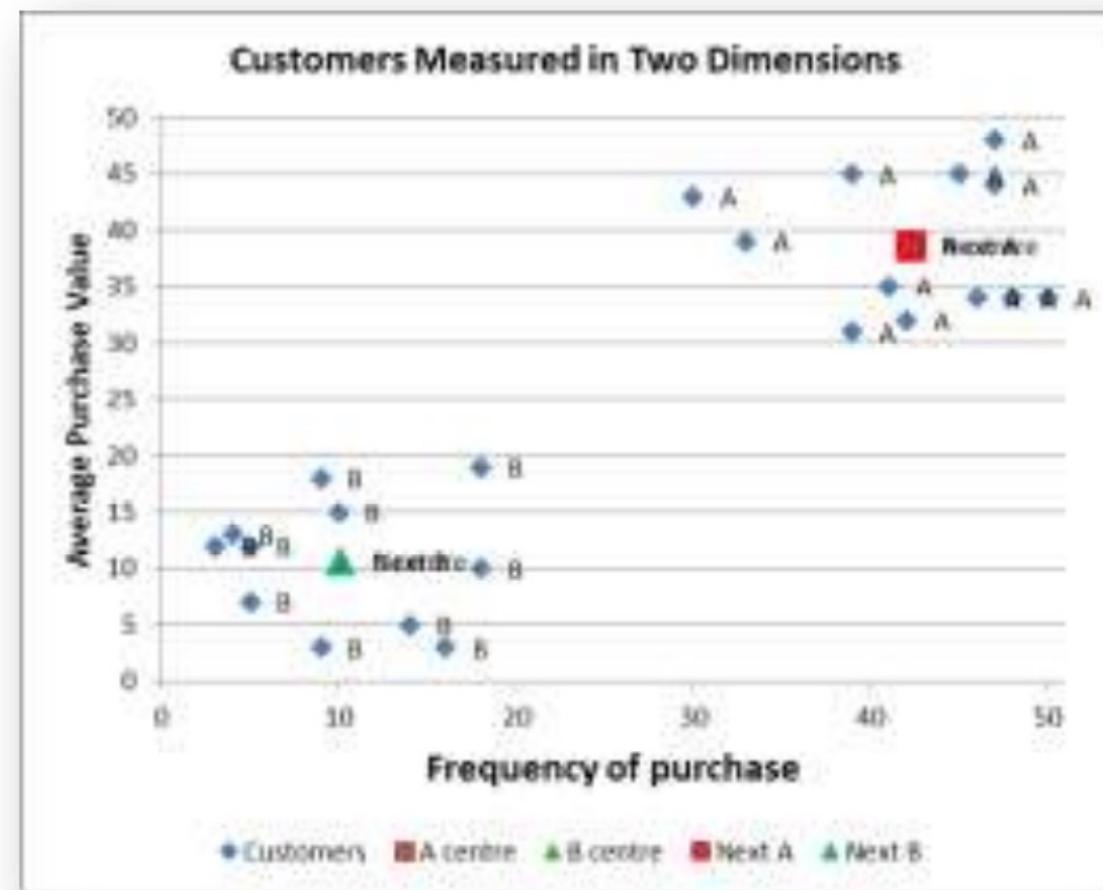
Step 4 : Use those Points as your New Best Guess

- Those centre points become the new 'best guesses'

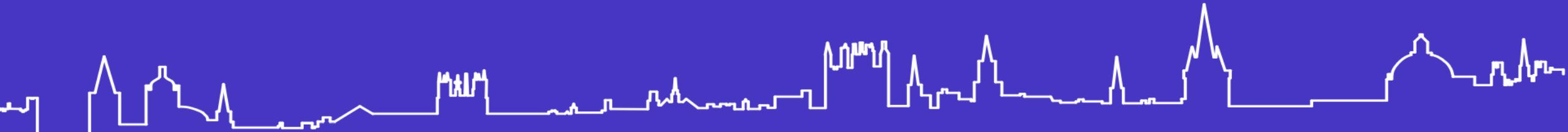


Repeat Steps 2 .. 4

- This simple procedure is then repeated
- Each time it is repeated the data falls into more distinct groups
- After a few repeats our ‘guesses’ are actually at the centre of the group and the set has been divided



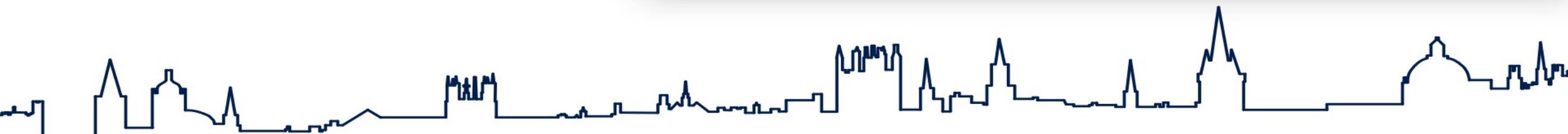
Clustering : Practical Example



Identifying Customer Groups using Clustering

- Our source data is over 10,000 customer records
- 5 factors for each customer
- We want to identify groups of customers based on the data

	A	B	C	D	E
1	Age	Gender	Married	Salary	Annual Spend
2	28.2	Male	Single	26908.95	331.56
3	43.5	Female	Married	39366.44	3071.18
4	27.7	Female	Single	28210.1	1357.19
5	18.9	Male	Single	26235.55	769.78
6	38	Male	Single	30822.14	100
7	32.9	Female	Single	18334.52	2854.59
8	43.5	Female	Married	36642.04	1926.37
10240	29.9	Female	Single	20955.77	3855.29
10241	43.7	Female	Married	36522.9	3801.28
10242	28.9	Female	Single	26408.18	2623.61
10243	28.8	Female	Single	25518.66	1958.84
10244	23.8	Female	Single	24218.53	1779.65
10245	26.6	Female	Single	26793.32	2352.41
10246	22.1	Male	Single	29096.49	410.25
10247	23	Female	Single	23437.99	1945.92
10248	29.8	Male	Single	28754.04	2673.34
10249					

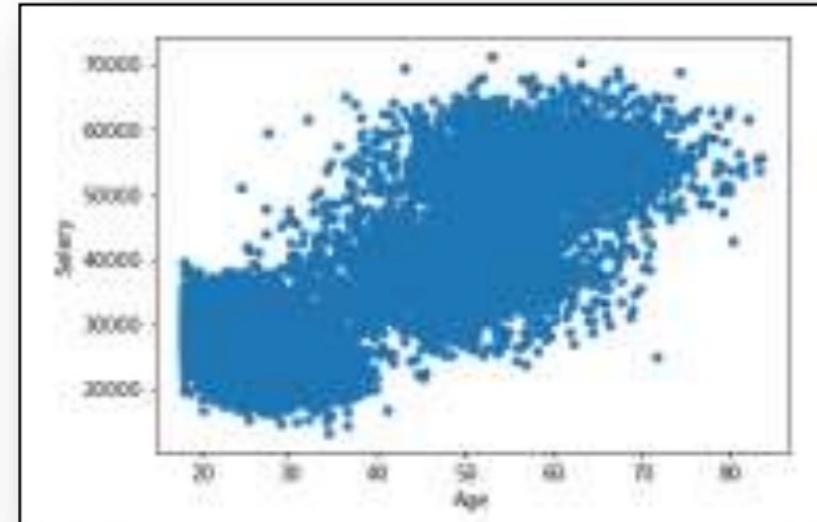


Import the Data and Take a look ...

```
customer_data = pd.read_csv("Retail Data v4.csv")
```

```
print(customer_data[0:5])
```

	Age	Gender	Married	Salary	Annual Spend
0	28.2	Male	Single	26908.95	331.56
1	43.5	Female	Married	39366.44	3071.18
2	27.7	Female	Single	28210.10	1357.19
3	18.9	Male	Single	26235.55	769.78
4	18.0	Male	Single	30822.14	100.00



```
customer_data.plot.scatter('Age', 'Salary')
```



Build the Model and Review Results

```
# create kmeans object
kmeans = KMeans(n_clusters=4)

# fit kmeans object to data
kmeans.fit(standardized_customer_data_df)

# save new clusters for chart
y_km = kmeans.fit_predict(standardized_customer_data_df)
```

```
print(customer_data[y_km==1][0:10])
```

	Married	Single	Female	Male	Age	Salary	Annual Spend
2	0	1	1	0	27.7	28210.10	1357.19
5	0	1	1	0	22.9	18334.52	2854.59
9	0	1	1	0	27.2	21901.06	566.21
11	0	1	1	0	31.0	24944.10	1661.97
12	0	1	1	0	34.9	23136.51	1343.36
14	0	1	1	0	32.8	21680.78	2600.72
15	0	1	1	0	28.7	27930.64	3717.44
17	0	1	1	0	28.0	25404.68	756.29
19	0	1	1	0	30.4	30677.96	1120.23
20	0	1	1	0	29.9	25182.12	1391.35

Chart the Results

