

Understanding your data to use for Machine Learning

What is Data and how do we use it?

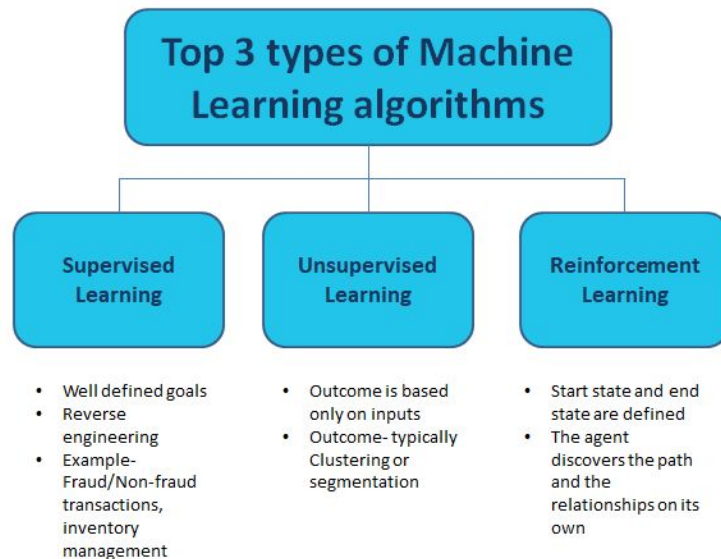
Data is a piece of raw information (facts; causes and consequences) that allows us to extract relevant knowledge for a particular event.

The 5 V's of data

1. *Volume* - The scale/quantity of Data.
2. *Variety* - Different forms of data – images, videos, audio, and so on.
3. *Velocity* - Rate of data streaming and generation. How often do we capture new data?
4. *Value* - Meaningfulness/relevance of data in terms of information that one might require to make a logical inference.
5. *Veracity* - Certainty and correctness in data we are working on.

Types of Machine Learning

1. **Supervised Learning** - Learning with existing information about the outcomes for an event (labeled data).
2. **Unsupervised Learning** - Finding similarities in the causal features and clusters information which may result in similar outcomes (unlabeled data).
3. **Reinforcement Learning** - Experience as you explore. Model uses trial-and-error to explore different techniques, and their outcomes define whether the model learns correctly. Favorable outputs are encouraged or 'reinforced', and non-favorable outputs are discouraged or 'punished'.



What kind of problems can we tackle using Supervised Learning?

Two types of learning objectives using labeled data:

1. **Classification** - Expected output is always from a discrete set of classes.
e.g.- whether an image contains a dog or cat, credit card fraud detection, etc.
2. **Regression** - Output is continuous in nature. There can be infinite possible answers.
e.g.- Price of a stock tomorrow, salary of an employee based on qualifications, etc.

Stages of Data Preprocessing

1. Data Exploration

2. Data Cleansing

3. Data Transformation

4. Feature Engineering

Data Exploration (Assessing your data)

1. What data is available and how much?
2. Do you have access to the ground truth, the values you're trying to predict?
3. What format will the data be in?
4. Which fields are most important?
5. Are the fields available in machine-readable form?
6. What important metrics are reported using this data?

Let's test it out!

Walmart Weekly Sales Dataset (<https://www.kaggle.com/yasserh/walmart-dataset>)

Data Cleansing

Detecting any incomplete or incorrect values.

- Dealing with missing values
- Dealing with outliers
- Correcting typos
- Grouping sparse classes
- Dropping duplicates

Data Transformation

Data you have available may not be in the right format or may require transformations to make it more useful.

- Categorical encoding
- Dealing with skewed data (Scaling)
- Bias mitigation (Correlation v/s Causation)

Feature Engineering

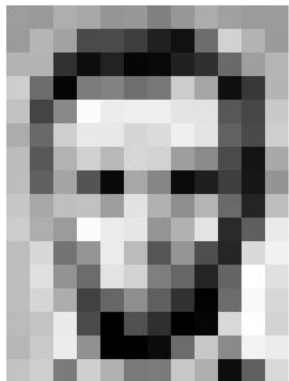
Feature engineering is the process of exploring new features based upon knowledge about current features and the required task.

- Feature Extraction
- Capturing Feature Relationships

Image Processing & Computer Vision

Representing the image data

- Array of pixels (represented as numbers) - colour channel.
- Images can be monochrome or coloured.
- Coloured images have 3 arrays - RGB channels



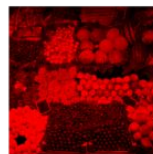
157	153	174	168	150	152	128	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	84	6	16	83	48	105	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	58	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	153	158	227	178	143	182	106	35	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	95	101	255	224
190	214	173	66	103	143	96	90	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	84	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	58	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	153	158	227	178	143	182	106	35	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	95	101	255	224
190	214	173	66	103	143	96	90	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

RGB



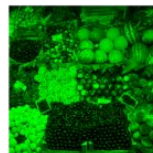
R



Red



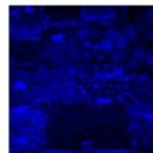
G



Green



B

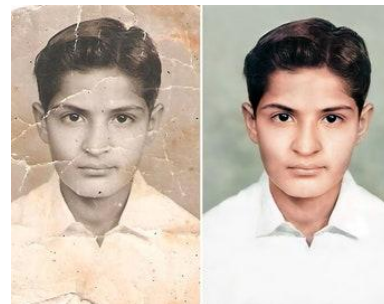
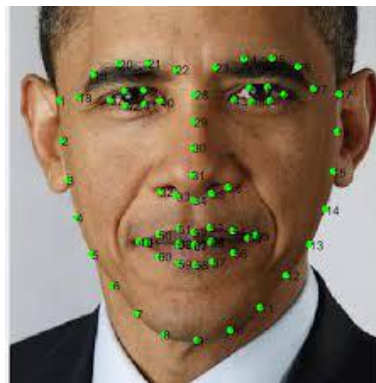
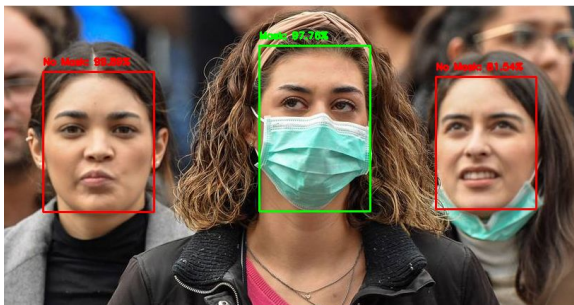


Blue



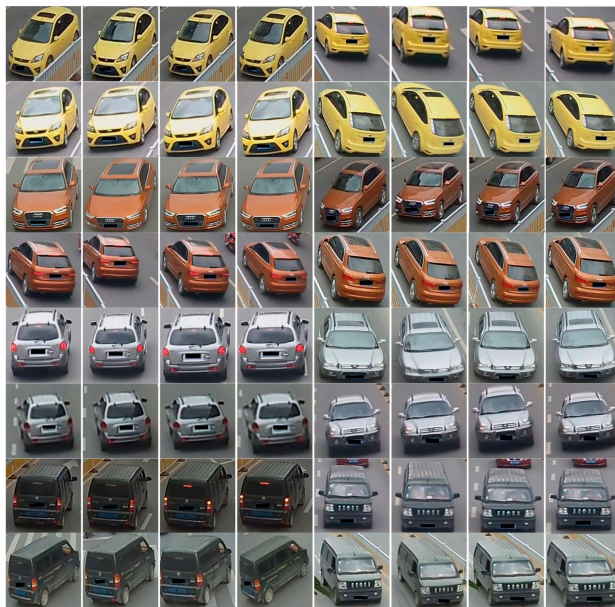
Applications of Computer Vision

- Object Classification
- Object Recognition
- Object Verification
- Object Detection
- Object Landmark Detection
- Object Segmentation
- Video motion analysis
- Scene reconstruction
- Image restoration



Difficulties while working with image data

- Size of the dataset
- Preprocessing overhead
- Quantity and quality of the data
- Eliminating redundant features



What other kinds of data?

More sources of data

- Simple Tabular Data - ML algorithms, DNN
- Image data - Convolution Networks (CNN)
- Time series data (Forecasting) - Recurrent Networks (RNN) / LSTMs
- Video data - CNN (Images) + RNN (Sequence)
- Textual data - Word embeddings, Vectorization (Word2Vec/Doc2Vec), LSTMs
- Audio data - CNN (Mel Spectrogram), LSTMs
- Reinforcement Learning ???

Thank you