

IMPERIAL

Machine Learning for Neuroscience

ML4NS

## Neuroscience-inspired Machine Learning and Applications in Neuroscience

Payam Barnaghi

Department of Brain Sciences &

School of Convergence Science in Human and Artificial Intelligence

Imperial College London

January 2025

I

1

### Neuroscience and AI

IMPERIAL

- Neuroscience provides a rich source of **inspiration** for new types of algorithms and architectures.
- This is often independent of and complementary to the mathematical and logic-based methods and ideas that have largely dominated traditional approaches to AI.
- Neuroscience could help to provide a systems neuroscience-level view of the brain.
- It can help view the network, architecture, functions, and representations the brain utilises.
- The precise mechanisms by which the processes/interactions are physically realised in a biological substrate are often less relevant at the implementation level (of AI models).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017

2

2

1

## Transferrable ideas from neuroscience

IMPERIAL

- By focusing on the computational and algorithmic levels, we can obtain transferrable insights into general mechanisms of brain function while leaving room to accommodate the distinctive opportunities and challenges that arise when building intelligent machines *in silico*.
- For example, biological considerations informed the development of successful regularisation schemes that support generalisation beyond training data.
- One such scheme, in which only a subset of units participate in the processing of a given training example (“dropout”), was motivated by the stochasticity that is inherent in biological systems populated by neurons that fire with Poisson-like statistics ([Hinton et al., 2012](#)).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017

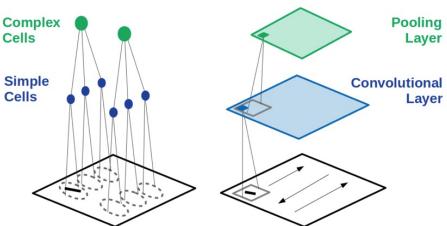
3

3

## Deep learning and biological systems

IMPERIAL

- In both biological and artificial systems, successive non-linear computations transform raw visual input into an increasingly complex set of features, permitting object recognition that is invariant to transformations of pose, illumination, or scale.
- Hubel and Wiesel discovered that simple cells (left, blue) have preferred locations in the image (dashed ovals) wherein they respond most strongly to bars of particular orientation. Complex cells (green) receive input from many simple cells and thus have more spatially invariant responses.



Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017  
Image source: Grace W. Lindsay, Gatsby Unit, UCL, <https://arxiv.org/pdf/2001.07092.pdf>

4

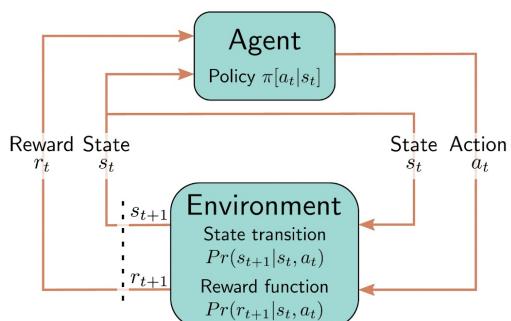
4

2

## Reinforcement learning (RL)

IMPERIAL

**Figure 19.6** Reinforcement learning loop. The agent takes an action  $a_t$  at time  $t$  based on the state  $s_t$ , according to the policy  $\pi[a_t|s_t]$ . This triggers the generation of a new state  $s_{t+1}$  (via the state transition function) and a reward  $r_{t+1}$  (via the reward function). Both are passed back to the agent, which then chooses a new action.



Source: Understanding Deep Learning, Simon J.D. Prince, MIT Press.

5

## Reinforcement learning

IMPERIAL

- Reinforcement learning methods address the problem of how to maximise future reward by mapping states in the environment to actions and are among the most widely used tools in AI research ([Sutton and Barto, 1998](#)).
- RL methods were originally inspired by research into animal learning. In particular, developing temporal-difference (TD) methods, a critical component of many RL models, was inextricably intertwined with research into animal behaviour in conditioning experiments.
- Examples of neuroscience-informed RL models include TD methods and related techniques that have gone on to supply the core technology for recent advances in AI, ranging from robotic control ([Hafner and Riedmiller, 2011](#)) to expert play in backgammon ([Tesauro, 1995](#)) and Go ([Silver et al., 2016](#)), and protein structures, AlphaFold ([Jumper et al., 2021](#)).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. Neuron. 2017

6

6

3

## Reinforcement learning in AlphaFold

- Several components of AlphaFold's machine learning pipeline rely on supervised learning, such as predicting residue-residue distance and orientation from large datasets.
- However, reinforcement learning enhances predictions by providing feedback and implementing an iterative optimisation process.

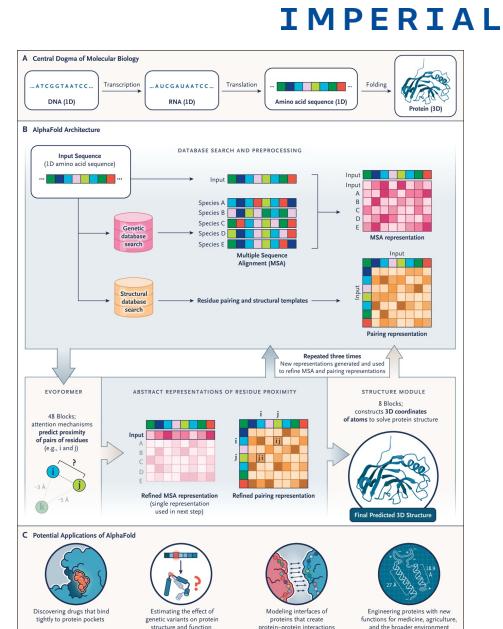


Image source: Altman RB. A Holy Grail - The Prediction of Protein Structure. N Engl J Med. 2023 Oct 12;389(15):1431-1434. doi: 10.1056/NEJMcb2307735.

7

## Attention-based models

- The brain does not learn by implementing a single, global optimisation principle within a uniform and undifferentiated neural network ([Marblestone et al., 2016](#)).
- Biological brains are modular, with distinct but interacting subsystems underpinning key functions such as memory, language, and cognitive control ([Anderson et al., 2004](#); [Shallice, 1988](#)).
- One illustrative example is recent AI work on attention. Until recently, most CNN models worked directly on entire images or video frames, with equal priority given to all image pixels at the earliest processing stage. The primate visual system works differently.

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. Neuron. 2017

8

## CNNs and attention

IMPERIAL

- Rather than processing all input in parallel, visual attention shifts strategically among locations and objects, centring processing resources and representational coordinates on a series of regions in turn ([Koch and Ullman, 1985](#); [Moore and Zirnsak, 2017](#); [Posner and Petersen, 1990](#)).
- Detailed neurocomputational models have shown how this piecemeal approach benefits behaviour, by prioritising and isolating the information that is relevant at any given moment ([Olshausen et al., 1993](#); [Salinas and Abbott, 1997](#)).
- One such network used this selective attentional mechanism to ignore irrelevant objects in a scene, allowing it to perform well in challenging object classification tasks in the presence of clutter ([Mnih et al., 2014](#)).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017

9

9

## Selective attention

IMPERIAL



Image source: TripAdvisor

10

10

# Attention-based models and Transformers

IMPERIAL

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

## Attention Is All You Need

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\*<sup>†</sup>**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaiser@google.com

**Illia Polosukhin\*<sup>‡</sup>**  
illia.polosukhin@gmail.com

### Abstract

II

11

# Standard neural networks

IMPERIAL

- A standard neural network layer  $f[x]$ , takes a  $D \times 1$  input  $x$ , and applies a linear transformation followed by an activation function  $a[\bullet]$ :

$$f[x] = a[\beta + \Omega x].$$

- Where,  $\beta$  contains the biases and  $\Omega$  contains the weights.

12

12

## Transformers\*\*

IMPERIAL

- Token: ["I", "like", "neuroscience"]

- Token embedding:

- "I" → [0.1, 0.2, 0.3]
- "like" → [0.1, 0.5, 0.6]
- "neuroscience" → [0.7, 0.9, 0.9]

- Input data:

0.1, 0.2, 0.3

0.1, 0.5, 0.6

0.7, 0.9, 0.9

13

13

## Transformers\*\*

IMPERIAL

- Input data:

$$\begin{bmatrix} 0.1, 0.2, 0.3 \\ 0.1, 0.5, 0.6 \\ 0.7, 0.9, 0.9 \end{bmatrix}$$

Final Input matrix:

$$\begin{bmatrix} (0.1+0.1), (0.2+0.2), (0.3+0.3) \\ (0.1+0.4), (0.5+0.5), (0.6+0.6) \\ (0.7+0.7), (0.9+0.8), (0.9+0.9) \end{bmatrix}$$

- Positional encoding:

Position 1: [0.1, 0.2, 0.3]

Position 2: [0.4, 0.5, 0.6]

Position 3: [0.7, 0.8, 0.9]

$$\begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}$$

14

14

## Transformers\*\*

**IMPERIAL**

- Query matrix (weights)
- Key matrix (weights)
- Value (weights)

$$X = \text{Input matrix: } \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}$$

$$\begin{array}{ccc} Q & K & V \\ \begin{bmatrix} 1, 0, 0 \\ 0, 1, 1 \\ 0, 0, 1 \end{bmatrix} & \begin{bmatrix} 1, 0, 0 \\ 0, 1, 0 \\ 0, 0, 1 \end{bmatrix} & \begin{bmatrix} 1, 0, 0 \\ 0, 1, 1 \\ 0, 0, 0 \end{bmatrix} \end{array}$$

$$X \cdot Q = \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}$$

$$X \cdot K = \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}$$

$$X \cdot V = \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}$$

15

15

## Transformers\*\*

**IMPERIAL**

$$\text{Attention}(Q_i, K_j) = \frac{Q_i \cdot K_j^T}{\sqrt{d_k}}$$

$$\begin{array}{l} X \cdot Q = \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix} \xrightarrow{\quad} Q \cdot K^T \\ \qquad \qquad \qquad \begin{bmatrix} 0.56, 1.22, 2.04 \\ 1.22, 2.69, 4.56 \\ 2.04, 4.56, 8.09 \end{bmatrix} \xrightarrow{\text{Divide by } \sqrt{3}} \begin{bmatrix} 0.3233, 0.7044, 1.1778 \\ 0.7044, 1.5531, 2.6327 \\ 1.1778, 2.6327, 4.6708 \end{bmatrix} \\ X \cdot K = \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix} \xrightarrow{\quad} \end{array}$$

Then apply softmax to these scores to normalise them into probabilities:

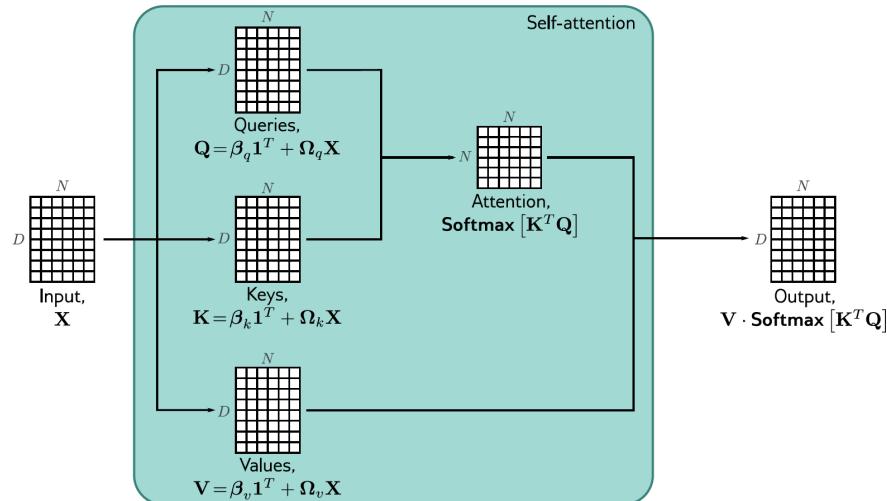
$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum e^{x_j}} \quad \begin{bmatrix} 0.2077 & 0.3041 & 0.4882 \\ 0.0979 & 0.2287 & 0.6734 \\ 0.0262 & 0.1122 & 0.8616 \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 0.2077 & 0.3041 & 0.4882 \\ 0.0979 & 0.2287 & 0.6734 \\ 0.0262 & 0.1122 & 0.8616 \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 0.2, 0.4, 0.6 \\ 0.5, 1.0, 1.2 \\ 1.4, 1.7, 1.8 \end{bmatrix}^T$$

16

16

## Self-attention in matrix form\*\*

IMPERIAL



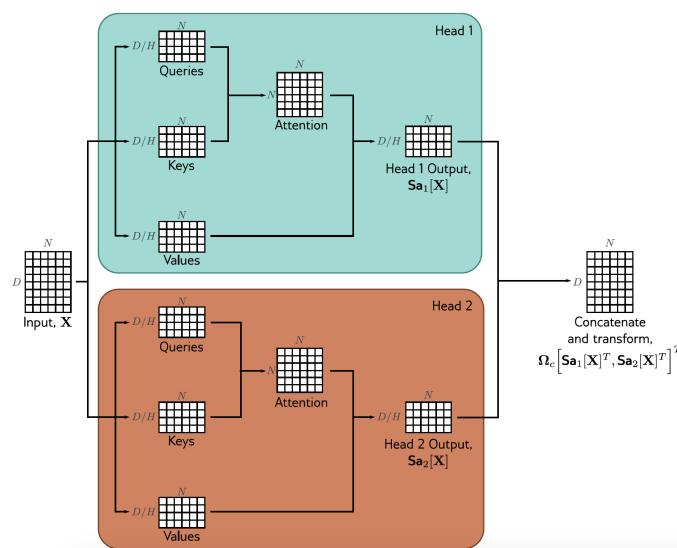
Understanding Deep Learning, Simon J.D. Prince, MIT Press.

17

17

## Multi-head attention\*\*

IMPERIAL



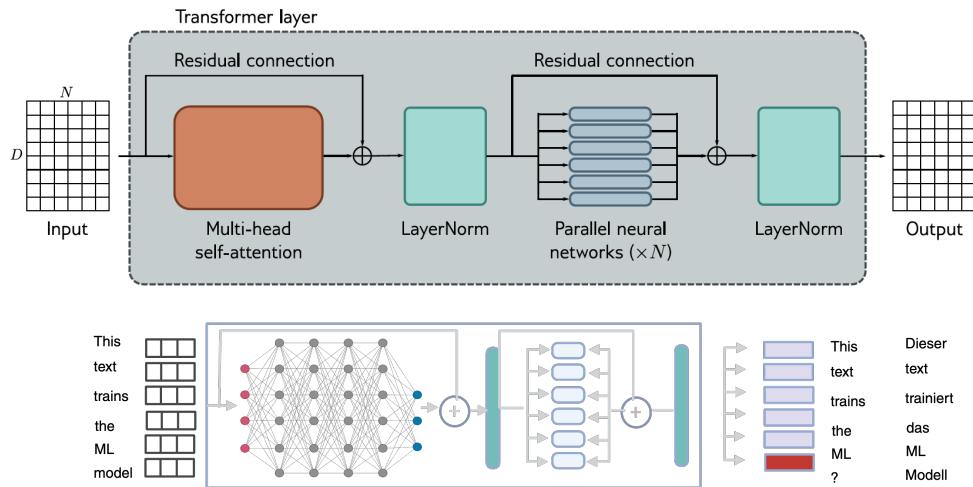
Understanding Deep Learning, Simon J.D. Prince, MIT Press.

18

18

## Transformer layers\*

IMPERIAL



Understanding Deep Learning, Simon J.D. Prince, MIT Press.

19

19

## Bidirectional Encoder Representations from Transformers (BERT)

IMPERIAL

- Encoder models like BERT exploit transfer learning.
- During pretraining, the parameters of the transformer architecture are learned using self-supervision from a large corpus of text.
- The goal here is for the model to learn general information about language statistics.
- In the fine-tuning stage, the resulting network is adapted to solve a particular task using a smaller body of supervised training data.

Understanding Deep Learning, Simon J.D. Prince, MIT Press.

20

20

## Pre-training

**IMPERIAL**

- In the pre-training stage, the network is trained using self-supervision. This allows the use of enormous amounts of data without the need for manual labels.
- For BERT, the self supervision task consists of predicting missing words from sentences from a large internet corpus.
- During training, the maximum input length is 512 tokens, and the batch size is 256.
- The system is trained for a million steps, corresponding to roughly 50 epochs of the 3.3-billion word corpus.

21

21

## Continual learning

**IMPERIAL**

- Intelligent agents must be able to learn and remember many different tasks that are encountered over multiple timescales.
- Both biological and artificial agents must thus have a capacity for continual learning, that is, an ability to master new tasks without forgetting how to perform prior tasks ([Thrun and Mitchell, 1995](#)). While animals appear relatively adept at continual learning, neural networks suffer from the problem of catastrophic forgetting ([French, 1999; McClelland et al., 1995](#)).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017

22

22

## Neuroimaging and continual learning

**IMPERIAL**

- In neuroscience, advanced neuroimaging techniques (e.g., two-photon imaging) now allow dynamic *in vivo* visualisation of the structure and function of dendritic spines during learning, at the spatial scale of single synapses ([Nishiyama and Yasuda, 2015](#)).
- This approach can be used to study neocortical plasticity during continual learning ([Cichon and Gan, 2015](#); [Hayashi-Takagi et al., 2015](#); [Yang et al., 2009](#)).

Source: Hassabis D, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*. 2017

23

23

## AI in neuroscience

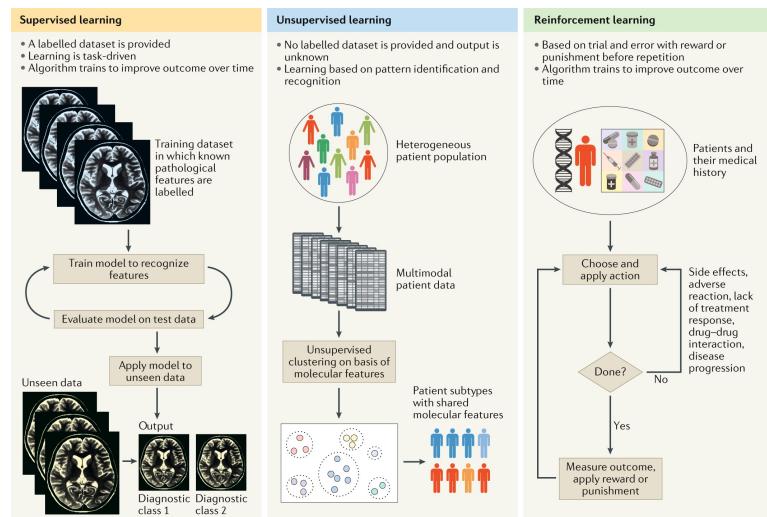
**IMPERIAL**

- Neuroimaging data analysis
- Behaviour and large population data analysis
- Computational models and simulations
- Graph and network analysis
- Fusion models and multimodal models for predictions/risk assessments
  
- An interesting read, How AI and neuroscience drive each other,  
<https://www.nature.com/articles/d41586-019-02212-4>

24

24

## Applications of ML to diagnosis and treatment of neurodegenerative diseases



Source : Myszczynska, M.A., Ojamies, P.N., Lacoste, A.M.B. et al. Applications of machine learning to diagnosis and treatment of neurodegenerative diseases. *Nat Rev Neuro* 16, 440–456 (2020).

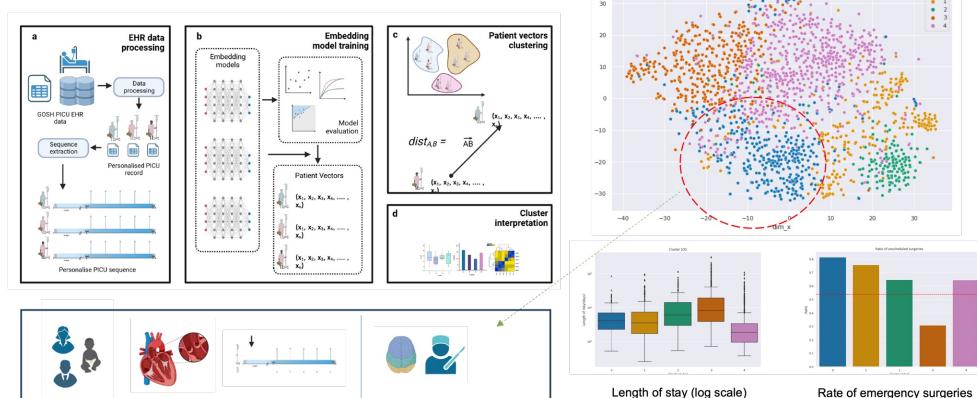
25

25

## ML and Electronic Healthcare Records Analysis

IMPERIAL

2 Years of GOSH Paediatric ICU Data



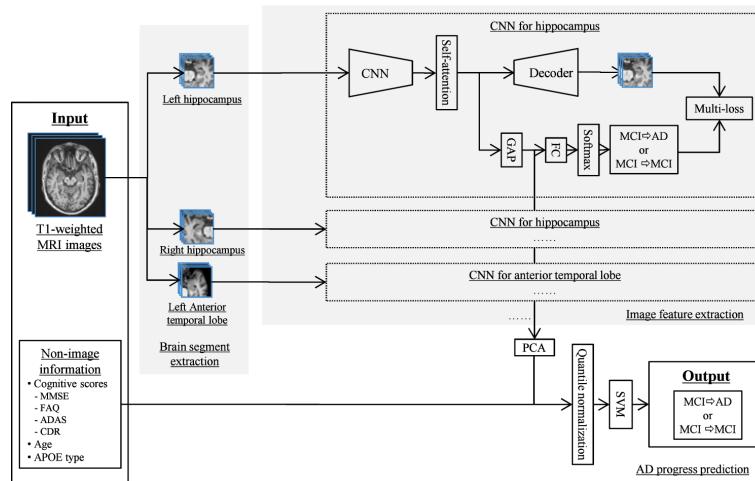
(J. Li et al., 2025).

26

26

## ML for predicting the progression of disease

IMPERIAL



Wang, C., Li, Y., Tsuboshita, Y. et al. A high-generalizability machine learning framework for predicting the progression of Alzheimer's disease using limited data. *npj Digit. Med.* 5, 43 (2022).

27

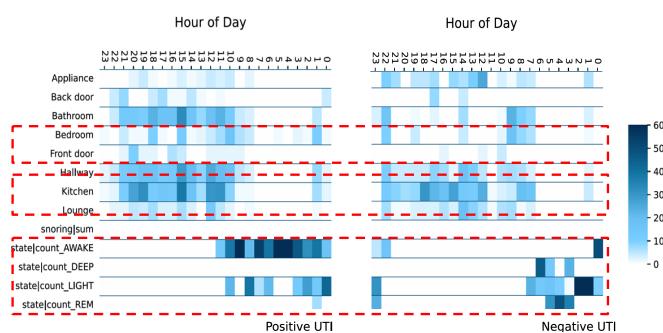
27

## ML in precision health and care

IMPERIAL

UTIs are one of the top 5 reasons on unplanned hospital admissions in People Living with Dementia. UTI symptoms include increase in frequency of going to bathroom, delirium, increased temperature.

The model uses 20,738 person-days of in-home monitoring data which consists sleep, movement, and physiology data collected within the UK DRI CR&T study cohort ( $n=108$ ) between 06/2021-07/2022.



(Alex Capstick et al., *npj Digital Medicine*, 2024).

28

28

**UTI risk analysis**

**IMPERIAL**

Date	Validation	Sensitivity	Specificity	Precision
Date	Test	86.6 (80.9 - 92.3)	94.5 (91.7 - 97.3)	87.3 (82.2 - 92.4)
Date-ID	Validation	69.0 (64.4 - 73.5)	94.1 (92.0 - 96.2)	81.9 (75.5 - 88.2)
Date-ID	Test	98.3 (95.5 - 101.1)	90.0 (85.5 - 94.5)	81.7 (74.4 - 89.1)
		74.7 (67.9 - 81.5)	87.9 (85.0 - 90.9)	77.0 (71.9 - 82.1)

Date	Female	Accuracy	No. of Participants	Positive : Negative	$\Pr(\hat{Y} = 1   \text{Sex})$
Date	Male	52.6 (31.8 - 73.4)	16	1 : 1.7	35.5 (32.4 - 38.5)
Date-ID	Female	88.5 (76.2 - 96.9)	25	1 : 5.5	22.3 (30.9 - 33.8)
Date-ID	Male	54.6 (26.8 - 82.4)	11	1 : 2.1	37.4 (33.6 - 41.1)
		85.3 (72.9 - 97.7)	20	1 : 4.1	38.0 (36.4 - 39.7)

(Alex Capstick et al., *npj Digital Medicine*, 2024).

29

29

**IMPERIAL**

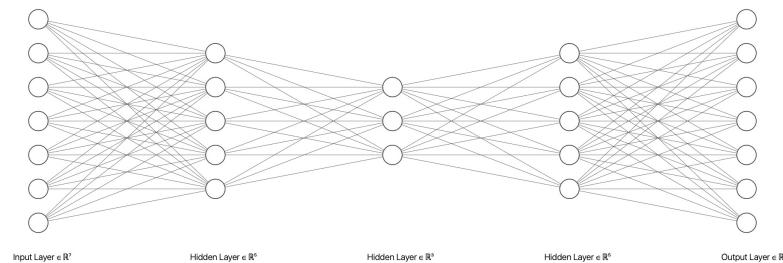
**Some other practical/useful methods and techniques**

30

## Deep Autoencoders

IMPERIAL

- A typical autoencoder consists of an encoder and a decoder.
- The encoder projects the input to hidden representations and the decoder maps the hidden layer to the reconstruction layer.

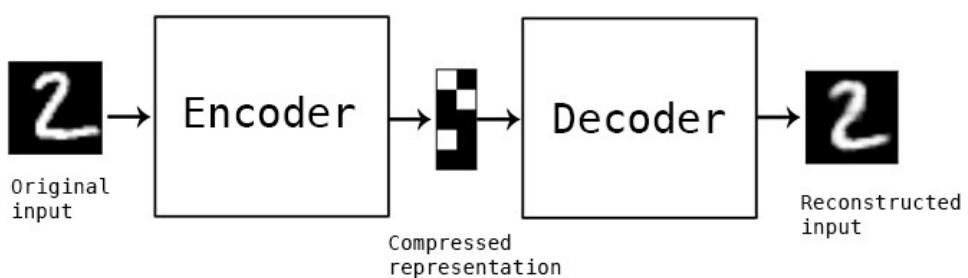


31

31

## Autoencoder example

IMPERIAL



Source: <https://blog.keras.io/building-autoencoders-in-keras.html>

32

32

## Autoencoders

IMPERIAL

- "Autoencoding" is a data representation algorithm in a lower dimension where the encoding and decoding functions:
  - data-specific
  - lossy,
  - *learned automatically from examples rather than engineered by a human.*
- Additionally, in almost all contexts where the term "autoencoder" is used, the encoder and decoder functions are implemented with neural networks.

Source: <https://blog.keras.io/building-autoencoders-in-keras.html>

33

33

## Autoencoders

IMPERIAL

- Practical applications of autoencoders include **data denoising** and **dimensionality reduction**.
- With appropriate dimensionality and sparsity constraints, autoencoders can learn data projections that are more efficient than PCA or other basic techniques.

Source: <https://blog.keras.io/building-autoencoders-in-keras.html>

34

34

## In other words

**IMPERIAL**

- An autoencoder is a neural network that is trained to attempt to copy its input to its output.
- Internally, it has a hidden layer  $h$  that describes a code (i.e. representation) used to represent the input.
- The network may be viewed as consisting of two parts: an encoder function  $h = f(x)$  and a decoder that produces a reconstruction  $r = g(h)$

Source: Ian Goodfellow et al., Deep Learning, MIT Press, 2016. <http://www.deeplearningbook.org>

35

35

## Important note on autoencoders

**IMPERIAL**

- If an autoencoder succeeds in simply learning to set  $g(f(x)) = x$  everywhere, then it is not especially useful.
- Instead, autoencoders are designed to be unable to learn to copy perfectly.
- Usually, autoencoders are restricted in ways that allow them to copy only approximately and to copy only input that resembles the training data. Because the model is forced to prioritise which aspects of the input should be copied, it often learns useful properties of the data.

Source: Ian Goodfellow et al., Deep Learning, MIT Press, 2016. <http://www.deeplearningbook.org>

36

36

## Dropouts in autoencoders

IMPERIAL

- You can add dropout to reduce overfitting.

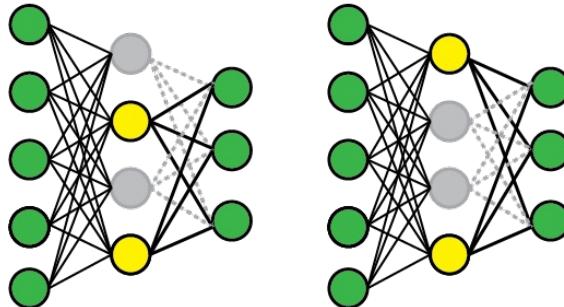


Image source: Matt Krause via <https://stats.stackexchange.com/questions/201569/what-is-the-difference-between-dropout-and-drop-connect>

37

37

## Undercomplete representation

IMPERIAL

- An autoencoder whose code dimension is less than the input dimension is called undercomplete.
- Learning an under-complete representation forces the autoencoder to capture the most salient features of the training data

38

38

19

## Learning process in autoencoders

IMPERIAL

- The learning process is described simply as minimising a loss function:

$$L(x, g(f(x)))$$

- where where  $L$  is a loss function penalising  $g(f(x))$  for being dissimilar from  $x$  such as the mean squared error.

Source: Ian Goodfellow et al., Deep Learning, MIT Press, 2016. <http://www.deeplearningbook.org>

39

39

## Linear decoders

IMPERIAL

- When the decoder is linear and  $L$  is the mean squared error, an undercomplete autoencoder learns to span the same subspace as PCA.
- In this case, an autoencoder trained to perform the copying task has learned the principal subspace of the training data as a side effect.
- Autoencoders with nonlinear encoder functions  $f$  and nonlinear decoder functions  $g$  can thus learn a more powerful nonlinear generalisation of PCA.

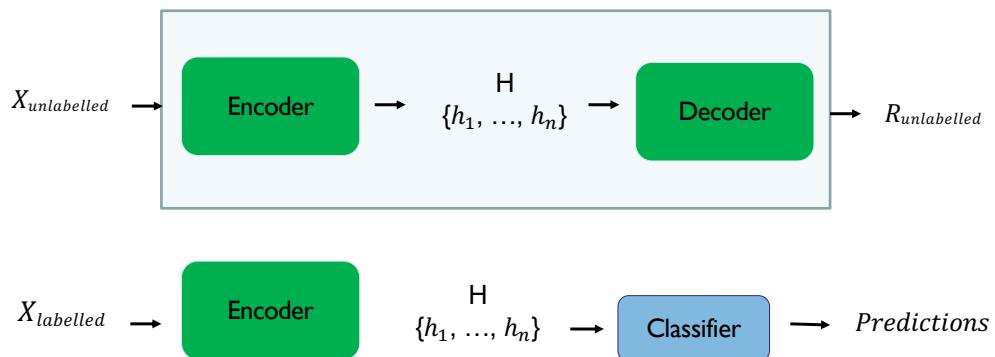
Source: Ian Goodfellow et al., Deep Learning, MIT Press, 2016. <http://www.deeplearningbook.org>

40

40

## A semi-supervised learning scenario

IMPERIAL



41

41

## Encoding categorical features

IMPERIAL

- Often, features are not given as continuous values but categorical.
- For example, a person could have features ["male", "female"], ["low", "mild", "high", "severe"]..
- Such features can be efficiently coded as integers; for instance, ["male", "has low symptoms"] could be expressed as [0, 1] while ["female", "has severe"] would be [1, 4].
- Such integer representation often cannot be used directly with machine learning models.

42

42

## Encoding categorical features

**IMPERIAL**

- Such integer representation cannot be used directly with *machine learning models* that expect continuous input and would interpret the categories as being ordered, which is often not desired (i.e. the set of browsers was ordered arbitrarily).
- One possibility to convert categorical features to features that can be used with *scikit-learn* estimators is to use a *one-of-K* or *one-hot* encoding, which is implemented in *OneHotEncoder*.
- This estimator transforms each categorical feature with  $m$  possible values into  $m$  binary features, with only one active.

43

43

## Encoding categorical features

**IMPERIAL**

- By default, the number of values each feature can take is inferred automatically from the dataset.
- It is possible to specify this explicitly using the parameter *n\_values*.
- For example, if there are two genders, three possible continents and four web browsers in our dataset. Then, we fit the estimator and transform a data point.
- In the result, the first two numbers encode the gender, the next set of three numbers is the continent and the last four are the web browser.

44

44

## OneHotEncoder

IMPERIAL

- The input to this transformer should be a matrix of integers, denoting the values taken on by categorical (discrete) features.
- The output will be a sparse matrix where each column corresponds to one possible value of one feature.
- It is assumed that input features take on values in the range [0, n\_values).
- This encoding is needed for feeding categorical data to many scikit-learn estimators, notably linear models and SVMs with the standard kernels.

45

45

## OneHotEncoder- example

IMPERIAL

46

46

## Imputation of missing values

**I M P E R I A L**

- For various reasons, many real-world datasets contain missing values, often encoded as blanks, NA or other placeholders.
- Such datasets, however, are incompatible with scikit-learn estimators, which assume that all values in an array are numerical and that all have and hold meaning.
- A basic strategy to use incomplete datasets is to discard entire rows and/or columns containing missing values.
- However, this comes at the price of losing data which may be valuable (even though incomplete).
- A better strategy is to impute the missing values, i.e., to infer them from the known part of the data.

47

47

## The *Imputer* class (in Python)

**I M P E R I A L**

- The *Imputer* class provides basic strategies for imputing missing values, either using the mean, the median or the most frequent value of the row or column in which the missing values are located.
- This class also allows for different missing values encodings.
- **strategy:** string, optional (default="mean")
- The imputation strategy.
  - If "mean", then replace missing values using the mean along the axis.
  - If "median", then replace missing values using the median along the axis.
  - If "most\_frequent", then replace missing using the most frequent value along the axis.

48

48

## Data imputation

IMPERIAL

- There are obviously other more sophisticated imputation techniques.
- However, the most important consideration is knowing the data, application context, and the techniques and solutions that would be more relevant to the given data/context and those that won't add more bias to the data.

49

49

## Feature selection

IMPERIAL

- The classes in the `sklearn.feature_selection` module can be used for feature selection/dimensionality reduction on sample sets, either to improve estimators' accuracy scores or to boost their performance on very high-dimensional datasets.

50

50

## Removing features with low variance

IMPERIAL

- *VarianceThreshold* is a simple baseline approach to feature selection.
- It removes all features whose variance doesn't meet some threshold. By default, it removes all zero-variance features, i.e. features that have the same value in all samples.
- There are several other methods/techniques and models for feature selection. We will/have studied some of those in this series.

51

51

IMPERIAL

## Some of the topics that we didn't cover in this module

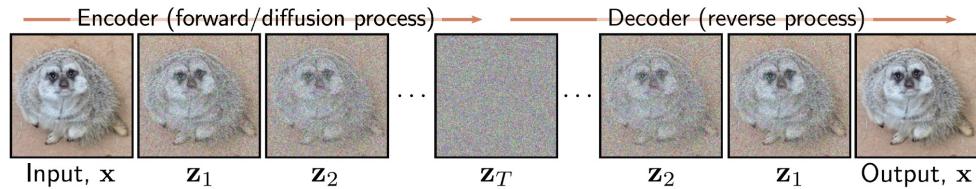
(some of these will be covered in the optional series in March/April)

52

52

## Diffusion models\*

IMPERIAL



**Figure 18.1** Diffusion models. The encoder (forward, or diffusion process) maps the input  $x$  through a series of latent variables  $z_1 \dots z_T$ . This process is pre-specified and gradually mixes the data with noise until only noise remains. The decoder (reverse process) is learned and passes the data back through the latent variables, removing noise at each stage. After training, new examples are generated by sampling noise vectors  $z_T$  and passing them through the decoder.

Source: Understanding Deep Learning, Simon J.D. Prince, MIT Press.

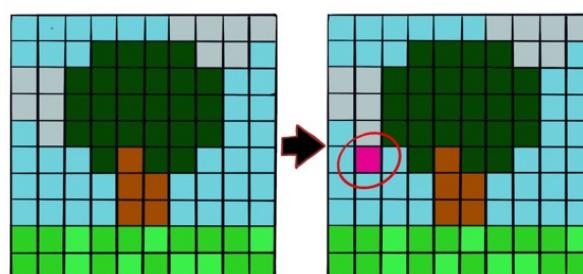
53

53

## Adversarial Examples

IMPERIAL

- Can we change one pixel in an image to deceive a machine learning model?



Original paper: Su, Jiawei, Danilo Vasconcelos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." *IEEE Transactions on Evolutionary Computation* (2019).

Image source: <https://christophm.github.io/interpretable-ml-book/adversarial.html>

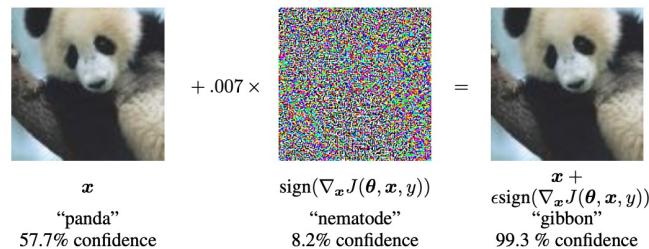
54

54

## A linear adversarial example\*

IMPERIAL

- A demonstration of fast adversarial example generation applied to GoogLeNet(Szegedy et al., 2014) on ImageNet.



Original paper: Explaining and Harnessing Adversarial Examples, Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy,  
<https://arxiv.org/abs/1412.6572>

55

## If you have any questions

IMPERIAL

- Please feel free to arrange a meeting or email ([p.barnaghi@imperial.ac.uk](mailto:p.barnaghi@imperial.ac.uk)).
- My office: 928, Sir Michael Uren Research Hub, White City Campus.

56

56