

# Differentiable Programming (and some Deep Learning)

Jonathon Hare

Vision, Learning and Control  
University of Southampton

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data

$$\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

Regularisation  $E_1 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2 + r(\|\boldsymbol{\theta}\|)$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

Regularisation  $E_1 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2 + r(\|\boldsymbol{\theta}\|)$

Modelling Uncertainty  $p(\boldsymbol{\theta} | \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

Regularisation  $E_1 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2 + r(\|\boldsymbol{\theta}\|)$

Modelling Uncertainty  $p(\boldsymbol{\theta} | \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$

Probabilistic Inference  $\mathbb{E}[g(\boldsymbol{\theta})] = \int g(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{N_s} \sum_{n=1}^{N_s} g(\boldsymbol{\theta}^{(n)})$

# Machine Learning - A Recap

All credit for this slide goes to Niranjan

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \nu$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

Regularisation  $E_1 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2 + r(\|\boldsymbol{\theta}\|)$

Modelling Uncertainty  $p(\boldsymbol{\theta} | \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$

Probabilistic Inference  $\mathbb{E}[g(\boldsymbol{\theta})] = \int g(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{N_s} \sum_{n=1}^{N_s} g(\boldsymbol{\theta}^{(n)})$

Sequence Modelling  $\mathbf{x}_n = f(\mathbf{x}_{n-1}, \boldsymbol{\theta})$

# What is Deep Learning?

Deep learning is primarily characterised by function compositions:

# What is Deep Learning?

Deep learning is primarily characterised by function compositions:

- Feedforward networks:  $\mathbf{y} = f(g(\mathbf{x}, \theta_g), \theta_f)$ 
  - Often with relatively simple functions (e.g.  $f(\mathbf{x}, \theta_f) = \sigma(\mathbf{x}^\top \theta_f)$ )

# What is Deep Learning?

Deep learning is primarily characterised by function compositions:

- Feedforward networks:  $\mathbf{y} = f(g(\mathbf{x}, \theta_g), \theta_f)$ 
  - Often with relatively simple functions (e.g.  $f(\mathbf{x}, \theta_f) = \sigma(\mathbf{x}^\top \theta_f)$ )
- Recurrent networks:  
$$\mathbf{y}_t = f(\mathbf{y}_{t-1}, \mathbf{x}_t, \theta) = f(f(\mathbf{y}_{t-2}, \mathbf{x}_{t-1}, \theta), \mathbf{x}_t, \theta) = \dots$$

# What is Deep Learning?

Deep learning is primarily characterised by function compositions:

- Feedforward networks:  $\mathbf{y} = f(g(\mathbf{x}, \theta_g), \theta_f)$ 
  - Often with relatively simple functions (e.g.  $f(\mathbf{x}, \theta_f) = \sigma(\mathbf{x}^\top \theta_f)$ )
- Recurrent networks:  
$$\mathbf{y}_t = f(\mathbf{y}_{t-1}, \mathbf{x}_t, \theta) = f(f(\mathbf{y}_{t-2}, \mathbf{x}_{t-1}, \theta), \mathbf{x}_t, \theta) = \dots$$

In the early days the focus of deep learning was on learning functions for classification. Nowadays the functions are much more general in their inputs and outputs.

# What is Differentiable Programming?

- Differentiable programming is a term coined by Yann Lecun<sup>1</sup> to describe a superset of Deep Learning.

---

<sup>1</sup><https://www.facebook.com/yann.lecun/posts/10155003011462143>

<sup>2</sup>See our ICLR 2019 paper: <https://arxiv.org/abs/1812.03928> and NeurIPS 2019 paper:  
<https://arxiv.org/abs/1906.06565>

# What is Differentiable Programming?

- Differentiable programming is a term coined by Yann Lecun<sup>1</sup> to describe a superset of Deep Learning.
- Captures the idea that computer programs can be constructed of parameterised functional blocks in which the parameters are learned using some form of gradient-based optimisation.

---

<sup>1</sup><https://www.facebook.com/yann.lecun/posts/10155003011462143>

<sup>2</sup>See our ICLR 2019 paper: <https://arxiv.org/abs/1812.03928> and NeurIPS 2019 paper: <https://arxiv.org/abs/1906.06565>

# What is Differentiable Programming?

- Differentiable programming is a term coined by Yann Lecun<sup>1</sup> to describe a superset of Deep Learning.
- Captures the idea that computer programs can be constructed of parameterised functional blocks in which the parameters are learned using some form of gradient-based optimisation.
  - The implication is that we need to be able to compute gradients with respect to the parameters of these functional blocks. We'll start explore this in detail next week...

---

<sup>1</sup><https://www.facebook.com/yann.lecun/posts/10155003011462143>

<sup>2</sup>See our ICLR 2019 paper: <https://arxiv.org/abs/1812.03928> and NeurIPS 2019 paper: <https://arxiv.org/abs/1906.06565>

# What is Differentiable Programming?

- Differentiable programming is a term coined by Yann Lecun<sup>1</sup> to describe a superset of Deep Learning.
- Captures the idea that computer programs can be constructed of parameterised functional blocks in which the parameters are learned using some form of gradient-based optimisation.
  - The implication is that we need to be able to compute gradients with respect to the parameters of these functional blocks. We'll start explore this in detail next week...
  - The idea of Differentiable Programming also opens up interesting possibilities:
    - The functional blocks don't need to be direct functions in a mathematical sense; more generally they can be *algorithms*.
    - What if the functional block we're learning parameters for is itself an algorithm that optimises the parameters of an internal algorithm using a gradient based optimiser?!<sup>2</sup>

---

<sup>1</sup><https://www.facebook.com/yann.lecun/posts/10155003011462143>

<sup>2</sup>See our ICLR 2019 paper: <https://arxiv.org/abs/1812.03928> and NeurIPS 2019 paper: <https://arxiv.org/abs/1906.06565>

# Is all Deep Learning Differentiable Programming?

- Not necessarily!
  - Most deep learning systems are trained using first order gradient-based optimisers, but there is an active body of research on gradient-free methods.

---

<sup>3</sup>including at least myself, my PhD students and Geoff Hinton!

# Is all Deep Learning Differentiable Programming?

- Not necessarily!
  - Most deep learning systems are trained using first order gradient-based optimisers, but there is an active body of research on gradient-free methods.
  - There is an increasing interest in methods that use different styles of learning, such as Hebbian learning, within deep networks. More broadly there are a number of us<sup>3</sup> who are interested in biologically motivated models and learning methods.

---

<sup>3</sup>including at least myself, my PhD students and Geoff Hinton!

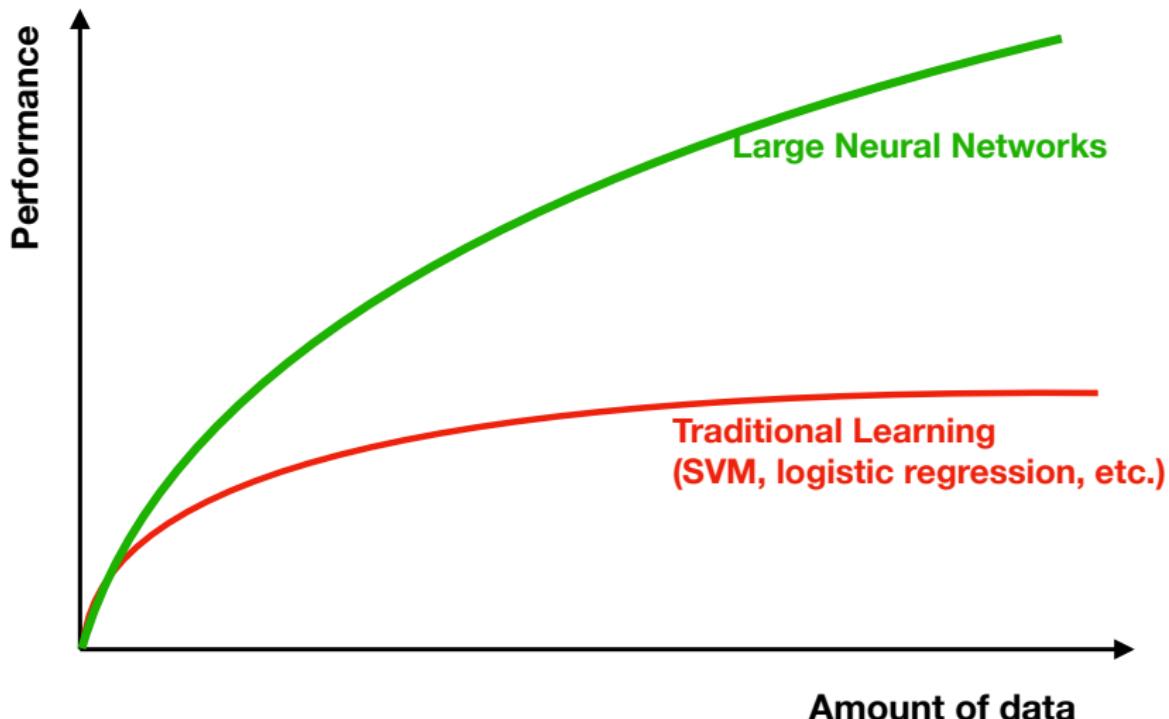
# Is all Deep Learning Differentiable Programming?

- Not necessarily!
  - Most deep learning systems are trained using first order gradient-based optimisers, but there is an active body of research on gradient-free methods.
  - There is an increasing interest in methods that use different styles of learning, such as Hebbian learning, within deep networks. More broadly there are a number of us<sup>3</sup> who are interested in biologically motivated models and learning methods.
  - This course will primarily focus on differentiable methods, but we'll look at how relaxations can be made to make non-differentiable operators learnable with gradient-based optimisers.

---

<sup>3</sup>including at least myself, my PhD students and Geoff Hinton!

# Why should we care about this?



---

Reference: Andrew Ng

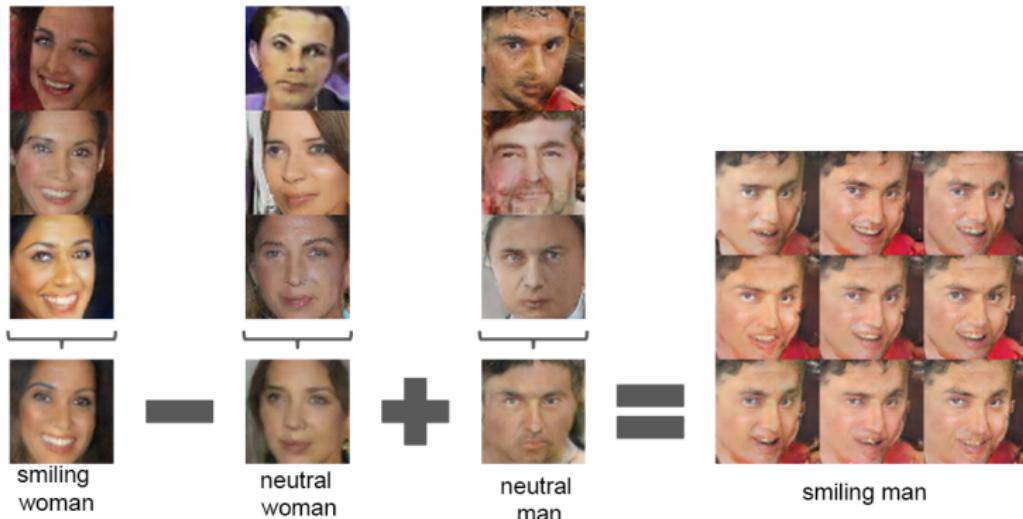
# Success stories - Object detection and segmentation



---

Pinheiro, Pedro O., et al. "Learning to refine object segments." European Conference on Computer Vision. Springer, 2016.

# Success stories - Image generation



Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv preprint arXiv:1511.06434 (2015).

# Success stories - Translation

## ENGLISH TEXT

The reason Boeing are doing this is to cram more seats in to make their plane more competitive with our products," said Kevin Keniston, head of passenger comfort at Europe's Airbus.

## TRANSLATED TO FRENCH

La raison pour laquelle Boeing fait cela est de creer plus de sieges pour rendre son avion plus competitif avec nos produits", a declare Kevin Keniston, chef du confort des passagers chez Airbus.

---

Wu, Yonghui, et al. "Google's neural machine translation system: Bridging the gap between human and machine translation." arXiv preprint arXiv:1609.08144 (2016).

# Types of Learning

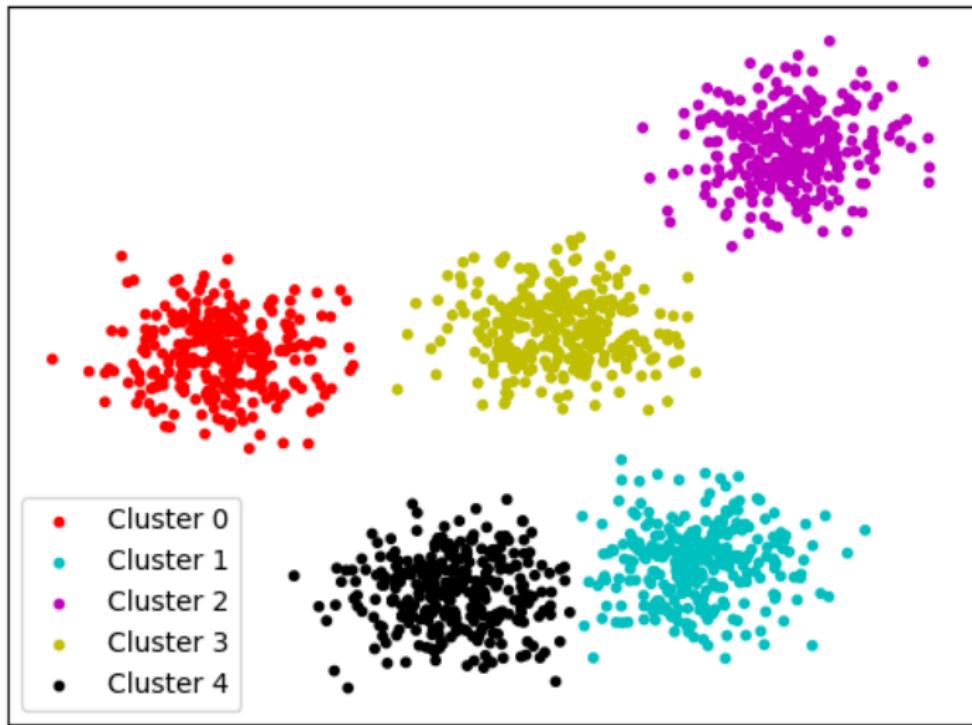
- Supervised Learning - learn to predict an output when given an input vector
- Unsupervised Learning - discover a good internal representation of the input
- Reinforcement Learning - learn to select an action to maximize the expectation of future rewards (payoff)
- Self-supervised Learning - learn with targets induced by a prior on the unlabelled training data
- Semi-supervised Learning - learn with few labelled examples and many unlabelled ones

# Supervised Learning

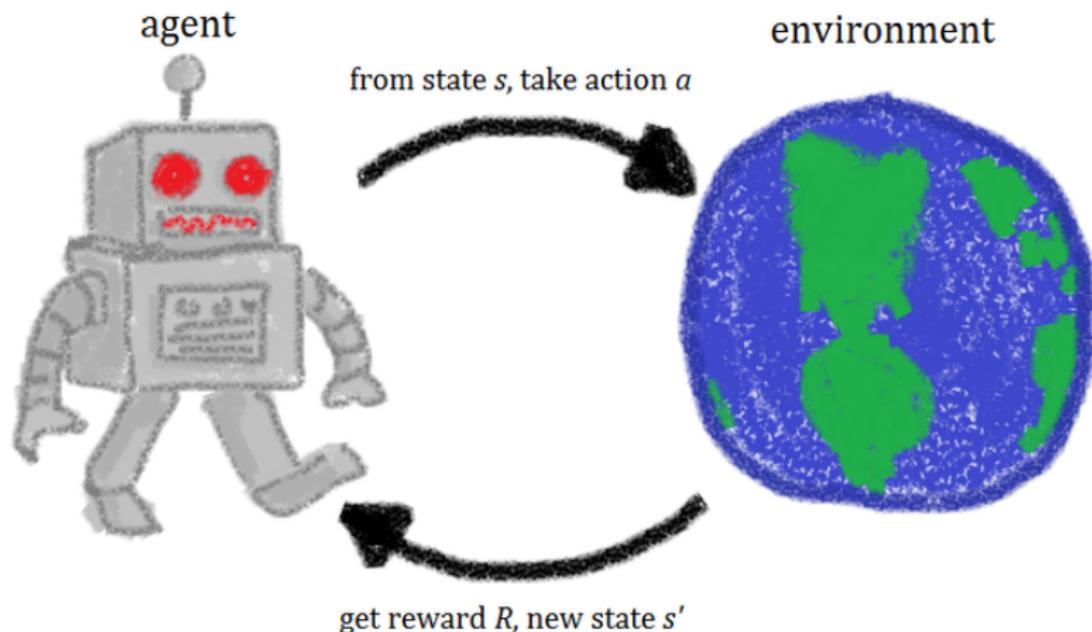


Newell, Alejandro, Kaiyu Yang, and Jia Deng. "Stacked hourglass networks for human pose estimation." ECCV'16. Springer, 2016.

# Unsupervised Learning



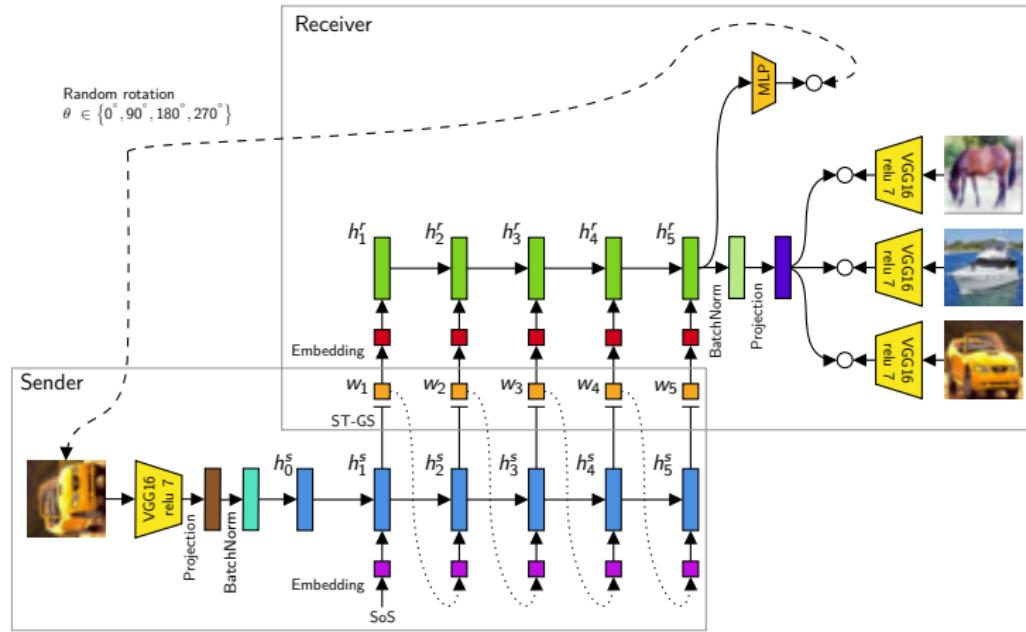
# Reinforcement Learning



Reference: Wikipedia

[https://simple.wikipedia.org/wiki/Reinforcement\\_learning](https://simple.wikipedia.org/wiki/Reinforcement_learning)

# Self-supervised Learning



Daniela Mihai and Jonathon Hare. Avoiding hashing and encouraging visual semantics in referential emergent language games. EmeCom @ NeurIPS 2019.  
<https://arxiv.org/abs/1911.05546>

# Semi-supervised Learning

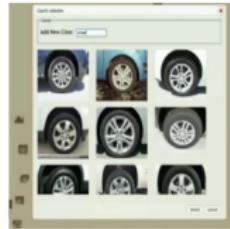
1. Start with unlabelled image data.



2. Use deep learning to automatically find structure in image.



3. Human adds labels to different clusters successfully classified, and distinguishes between images classified incorrectly.



4. Repeat stages 2&3 (re-train classifier using this additional info).

Jeremy Howard. The wonderful and terrifying implications of computers that can learn.  
TEDxBrussels. [http://www.ted.com/talks/jeremy\\_howard\\_the\\_wonderful\\_and\\_terrifying\\_im](http://www.ted.com/talks/jeremy_howard_the_wonderful_and_terrifying_im)

# Generative Models

- Many unsupervised and self-supervised models can be classed as 'Generative Models'.
- Given unlabelled data  $X$ , a unsupervised generative model learns  $P[X]$ .
  - Could be direct modelling of the data (e.g. Gaussian Mixture Models)
  - Could be indirect modelling by learning to map the data to a parametric distribution in a lower dimensional space (e.g. a VAEs Encoder) or by learning a mapping from a parameterised distribution to the real data space (e.g. a VAE Decoder or GAN)
- These are characterised by an ability to 'sample' the model to 'create' new data

## Generative vs. Discriminative Models (II)

Generative vs. discriminative approaches to classification use different statistical modelling.

- Discriminative models learn the boundary between classes. A discriminative model is a model of the conditional probability of the target  $Y$  given an observation  $X$ :  $P[Y|X]$ .
- Generative models of labelled data model the distribution of individual classes. Given an observable variable  $X$  and a target variable  $Y$ , a generative model is a statistical model that tries to model  $P[X|Y]$  and  $P[Y]$  in order to model the joint probability distribution  $P[X, Y]$ .<sup>4</sup>

---

<sup>4</sup>Some such models can be sampled conditionally based on a prior  $Y$  - e.g. a Conditional VAE: <https://papers.nips.cc/paper/5775-learning-structured-output-representation-using-deep-conditional-gene>

# Two Types of Supervised Learning

- Classification: The machine is asked to specify which of  $k$  categories some input belongs to.
  - Multiclass classification - target is one of the  $k$  classes
  - Multilabel classification - target is some number of the  $k$  classes
  - In both cases, the machine is a function  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  (although it is most common for the learning algorithm to actually learn  $\hat{f} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ ).

# Two Types of Supervised Learning

- Classification: The machine is asked to specify which of  $k$  categories some input belongs to.
  - Multiclass classification - target is one of the  $k$  classes
  - Multilabel classification - target is some number of the  $k$  classes
  - In both cases, the machine is a function  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  (although it is most common for the learning algorithm to actually learn  $\hat{f} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ ).
- Regression: The machine is asked predict  $k$  numerical values given some input. The machine is a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ .

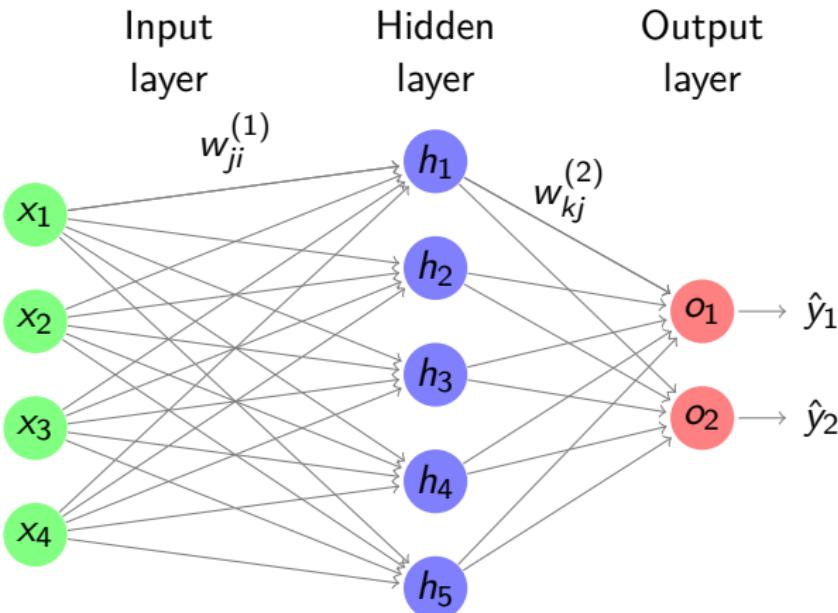
# Two Types of Supervised Learning

- Classification: The machine is asked to specify which of  $k$  categories some input belongs to.
  - Multiclass classification - target is one of the  $k$  classes
  - Multilabel classification - target is some number of the  $k$  classes
  - In both cases, the machine is a function  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  (although it is most common for the learning algorithm to actually learn  $\hat{f} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ ).
- Regression: The machine is asked predict  $k$  numerical values given some input. The machine is a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ .
- Note that there are lots of exceptions in the form the inputs (and outputs) can take though! We'll see lots of variations in the coming weeks.

# How Supervised Learning Typically Works

- Start by choosing a model-class:  $\hat{y} = f(\mathbf{x}; \mathbf{W})$  where the model-class  $f$  is a way of using some numerical parameters,  $\mathbf{W}$ , to map each input vector  $\mathbf{x}$  to a predicted output  $\hat{y}$ .
- Learning means adjusting the parameters to reduce the discrepancy between the true target output  $y$  on each training case and the output  $\hat{y}$ , predicted by the model.

# Let's look at an unbiased Multilayer Perceptron...



Without loss of generality, we can write the above as:

$$\hat{\mathbf{y}} = g(f(\mathbf{x}; \mathbf{W}^{(1)}); \mathbf{W}^{(2)}) = g(\mathbf{W}^{(2)} f(\mathbf{W}^{(1)} \mathbf{x}))$$

where  $f$  and  $g$  are activation functions.

# Common Activation Functions

- Identity
- Sigmoid (aka Logistic)
- Hyperbolic Tangent ( $\tanh$ )
- Rectified Linear Unit (ReLU) (aka Threshold Linear)

## Final layer activations

$$\hat{y} = g(\mathbf{W}^{(2)} f(\mathbf{W}^{(1)} \mathbf{x}))$$

- What form should the final layer function  $g$  take?

## Final layer activations

$$\hat{y} = g(\mathbf{W}^{(2)} f(\mathbf{W}^{(1)} \mathbf{x}))$$

- What form should the final layer function  $g$  take?
- It depends on the task (and on the chosen loss function)...
  - For regression it is typically linear (e.g. identity), but you might choose others if you say wanted to clamp the range of the network.
  - For binary classification (MLP has a single output), one would choose Sigmoid
  - For multilabel classification, typically one would choose Sigmoid
  - For multiclass classification, typically you would use the Softmax function

# Softmax

The softmax is an activation function used at the output layer of a neural network that forces the outputs to sum to 1 so that they can represent a probability distribution across a discrete mutually exclusive alternatives.

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \forall i = 1, 2, \dots, K$$

- Note that unlike the other activation functions you've seen, softmax makes reference to all the elements in the output.
- The output of a softmax layer is a set of positive numbers which sum up to 1 and can be thought of as a probability distribution.
- Note:

$$\frac{\partial \text{softmax}(\mathbf{z})_i}{\partial z_i} = \text{softmax}(z_i)(1 - \text{softmax}(z_i))$$

$$\begin{aligned}\frac{\partial \text{softmax}(\mathbf{z})_i}{\partial z_j} &= \text{softmax}(z_i)(1(i=j) - \text{softmax}(z_j)) \\ &= \text{softmax}(z_i)(\delta_{ij} - \text{softmax}(z_j))\end{aligned}$$

# Ok, so let's talk loss functions

- The choice of loss function depends on the task (e.g. classification/regression/something else)

# Ok, so let's talk loss functions

- The choice of loss function depends on the task (e.g. classification/regression/something else)
- The choice also depends on the activation function of the last layer

# Ok, so let's talk loss functions

- The choice of loss function depends on the task (e.g. classification/regression/something else)
- The choice also depends on the activation function of the last layer
  - For numerical reasons (see Log-Sum-Exp in a few slides) many times the activation is computed directly within the loss rather than being part of the model

# Ok, so let's talk loss functions

- The choice of loss function depends on the task (e.g. classification/regression/something else)
- The choice also depends on the activation function of the last layer
  - For numerical reasons (see Log-Sum-Exp in a few slides) many times the activation is computed directly within the loss rather than being part of the model
  - Some classification losses require *raw outputs* (e.g. a linear layer) of the network as their input
    - These are often called *unnormalised log probabilities* or *logits*
    - An example would be hinge-loss used to create a Support Vector Machine that maximises the margin — e.g.:
$$\ell_{\text{hinge}}(\hat{y}, y) = \max(0, 1 - y \cdot \hat{y})$$
with a true label,  $y \in \{-1, 1\}$ , for binary classification.

# Ok, so let's talk loss functions

- The choice of loss function depends on the task (e.g. classification/regression/something else)
- The choice also depends on the activation function of the last layer
  - For numerical reasons (see Log-Sum-Exp in a few slides) many times the activation is computed directly within the loss rather than being part of the model
  - Some classification losses require *raw outputs* (e.g. a linear layer) of the network as their input
    - These are often called *unnormalised log probabilities* or *logits*
    - An example would be hinge-loss used to create a Support Vector Machine that maximises the margin — e.g.:
$$\ell_{\text{hinge}}(\hat{y}, y) = \max(0, 1 - y \cdot \hat{y})$$
with a true label,  $y \in \{-1, 1\}$ , for binary classification.
- There are many different loss functions we might encounter (MSE, Cross-Entropy, KL-Divergence, huber, L1 (MAE), CTC, Triplet, ...) for different tasks.

# The Cost Function (measure of discrepancy)

Recall:

- Mean Squared Error (MSE) loss for a single data point (here assumed to be a vector, but equally applicable to a scalar) is given by
$$\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_i (\hat{y}_i - y_i)^2 = (\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y})$$
- We often multiply this by a constant factor of  $\frac{1}{2}$  — can anyone guess/remember why?

---

<sup>5</sup><http://neuralnetworksanddeeplearning.com/chap3.html>

# The Cost Function (measure of discrepancy)

Recall:

- Mean Squared Error (MSE) loss for a single data point (here assumed to be a vector, but equally applicable to a scalar) is given by
$$\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_i (\hat{y}_i - y_i)^2 = (\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y})$$
- We often multiply this by a constant factor of  $\frac{1}{2}$  — can anyone guess/remember why?
- $\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y})$  is the predominant choice for regression problems with linear activation in the last layer

---

<sup>5</sup><http://neuralnetworksanddeeplearning.com/chap3.html>

# The Cost Function (measure of discrepancy)

Recall:

- Mean Squared Error (MSE) loss for a single data point (here assumed to be a vector, but equally applicable to a scalar) is given by
$$\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_i (\hat{y}_i - y_i)^2 = (\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y})$$
- We often multiply this by a constant factor of  $\frac{1}{2}$  — can anyone guess/remember why?
- $\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y})$  is the predominant choice for regression problems with linear activation in the last layer
- For a classification problem with Softmax or Sigmoidal (or really anything non-linear) activations, MSE can cause slow learning, especially if the predictions are very far off the targets
  - Gradients of  $\ell_{MSE}$  are proportional to the difference in target and predicted multiplied by the gradient of the activation function<sup>5</sup>
  - The Cross-Entropy loss function is generally a better choice in this case

---

<sup>5</sup><http://neuralnetworksanddeeplearning.com/chap3.html>

# Binary Cross-Entropy

For the binary classification case:

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- The cross-entropy cost function is non-negative,  $\ell_{BCE} > 0$
- $\ell_{BCE} \approx 0$  when the prediction and targets are equal (i.e.  $y = 0$  and  $\hat{y} = 0$  or when  $y = 1$  and  $\hat{y} = 1$ )
- With Sigmoidal final layer,  $\frac{\partial \ell_{BCE}}{\partial \mathbf{W}_i^{(2)}}$  is proportional to just the error in the output ( $\hat{y} - y$ ) and therefore, the larger the error, the faster the network will learn!
- Note that the BCE is the negative log likelihood of the Bernoulli Distribution

## Binary Cross-Entropy — Intuition

- The cross-entropy can be thought of as a **measure of surprise**.
- Given some input  $x_i$ , we can think of  $\hat{y}_i$  as the estimated probability that  $x_i$  belongs to class 1, and  $1 - \hat{y}_i$  is the estimated probability that it belongs to class 0.
- Note the extreme case of infinite cross-entropy, if your model believes that a class has 0 probability of occurrence, and yet the class appears in the data, the 'surprise' of your model will be infinitely great.

## Binary Cross-Entropy for multiple labels

In the case of multi-label classification with a network with multiple sigmoidal outputs you just sum the BCE over the outputs:

$$\ell_{BCE} = - \sum_{k=1}^K [y_k \log(\hat{y}_k) + (1 - y_k) \log(1 - \hat{y}_k)]$$

where  $K$  is the number of classes of the classification problem,  $\hat{y} \in \mathbb{R}^K$ .

# Numerical Stability: The Log-Sum-Exp trick

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- Consider what might happen early in training when the model might confidently predict a positive example as negative
  - $\hat{y} = \sigma(z) \approx 0 \implies z \ll 0$

---

<sup>6</sup><https://www.xarg.org/2016/06/the-log-sum-exp-trick-in-machine-learning/>

# Numerical Stability: The Log-Sum-Exp trick

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- Consider what might happen early in training when the model might confidently predict a positive example as negative
  - $\hat{y} = \sigma(z) \approx 0 \implies z \ll 0$
  - if  $\hat{y}$  is small enough, it will become 0 due to limited precision of floating-point representations

---

<sup>6</sup><https://www.xarg.org/2016/06/the-log-sum-exp-trick-in-machine-learning/>

# Numerical Stability: The Log-Sum-Exp trick

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- Consider what might happen early in training when the model might confidently predict a positive example as negative
  - $\hat{y} = \sigma(z) \approx 0 \implies z \ll 0$
  - if  $\hat{y}$  is small enough, it will become 0 due to limited precision of floating-point representations
  - but then  $\log(\hat{y}) = -\infty$ , and everything will break!

---

<sup>6</sup><https://www.xarg.org/2016/06/the-log-sum-exp-trick-in-machine-learning/>

# Numerical Stability: The Log-Sum-Exp trick

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- Consider what might happen early in training when the model might confidently predict a positive example as negative
  - $\hat{y} = \sigma(z) \approx 0 \implies z \ll 0$
  - if  $\hat{y}$  is small enough, it will become 0 due to limited precision of floating-point representations
  - but then  $\log(\hat{y}) = -\infty$ , and everything will break!
- To tackle this problem implementations usually combine the sigmoid computation and BCE into a single loss function that you would apply to a network with linear outputs (e.g. `BCEWithLogitsLoss`).

---

<sup>6</sup><https://www.xarg.org/2016/06/the-log-sum-exp-trick-in-machine-learning/>

# Numerical Stability: The Log-Sum-Exp trick

$$\ell_{BCE}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- Consider what might happen early in training when the model might confidently predict a positive example as negative
  - $\hat{y} = \sigma(z) \approx 0 \implies z \ll 0$
  - if  $\hat{y}$  is small enough, it will become 0 due to limited precision of floating-point representations
  - but then  $\log(\hat{y}) = -\infty$ , and everything will break!
- To tackle this problem implementations usually combine the sigmoid computation and BCE into a single loss function that you would apply to a network with linear outputs (e.g. `BCEWithLogitsLoss`).
- Internally, a trick called ‘log-sum-exp’ is used to *shift* the centre of an exponential sum so that only numerical underflow can potentially happen, rather than overflow<sup>6</sup>.
  - Ultimately this means you’ll always get a numerically reasonable result (and will avoid NaNs and Infs originating from this point).

---

<sup>6</sup><https://www.xarg.org/2016/06/the-log-sum-exp-trick-in-machine-learning/>

# Multiclass classification with Softmax Outputs

- Softmax can be thought of making the  $K$  outputs of the network mimic a probability distribution.

---

<sup>7</sup>Note: Keras calls this function 'Categorical Cross-Entropy'; you would need to have a Softmax output layer to use this

# Multiclass classification with Softmax Outputs

- Softmax can be thought of making the  $K$  outputs of the network mimic a probability distribution.
- The target label  $y$  could also be represented as a distribution with a single 1 and zeros everywhere else.
  - e.g. they are “one-hot encoded”.

---

<sup>7</sup>Note: Keras calls this function ‘Categorical Cross-Entropy’; you would need to have a Softmax output layer to use this

# Multiclass classification with Softmax Outputs

- Softmax can be thought of making the  $K$  outputs of the network mimic a probability distribution.
- The target label  $y$  could also be represented as a distribution with a single 1 and zeros everywhere else.
  - e.g. they are “one-hot encoded”.
- In such a case, the obvious loss function is the *negative log likelihood* of the Categorical distribution (aka Multinoulli, Generalised Bernoulli, Multinomial with one sample)<sup>7</sup>:  $\ell_{NNL} = -\sum_{k=1}^K y_k \log \hat{y}_k$ 
  - Note that in practice as  $y_k$  is zero for all but one class you don't actually do this summation, and if  $y$  is an integer class index you can write  $\ell_{NNL} = -\log \hat{y}_y$ .

---

<sup>7</sup>Note: Keras calls this function ‘Categorical Cross-Entropy’; you would need to have a Softmax output layer to use this

# Multiclass classification with Softmax Outputs

- Softmax can be thought of making the  $K$  outputs of the network mimic a probability distribution.
- The target label  $y$  could also be represented as a distribution with a single 1 and zeros everywhere else.
  - e.g. they are “one-hot encoded”.
- In such a case, the obvious loss function is the *negative log likelihood* of the Categorical distribution (aka Multinoulli, Generalised Bernoulli, Multinomial with one sample)<sup>7</sup>:  $\ell_{NNL} = -\sum_{k=1}^K y_k \log \hat{y}_k$ 
  - Note that in practice as  $y_k$  is zero for all but one class you don't actually do this summation, and if  $y$  is an integer class index you can write  $\ell_{NNL} = -\log \hat{y}_y$ .
- Analogously to what we saw for BCE, Log-Sum-Exp can be used for better numerical stability.
  - PyTorch combines LogSoftmax with NNL in one loss and calls this “Categorical Cross-Entropy” (so you would use this with a *linear output layer*)

---

<sup>7</sup>Note: Keras calls this function ‘Categorical Cross-Entropy’; you would need to have a Softmax output layer to use this

## Reminder: Gradient Descent

- Define total loss as  $\mathcal{L} = -\sum_{(x,y) \in D} \ell(g(x, \theta), y)$  for some loss function  $\ell$ , dataset  $D$  and model  $g$  with learnable parameters  $\theta$ .
- Define how many passes over the data to make (each one known as an Epoch)
- Define a learning rate  $\eta$

Gradient Descent updates the parameters  $\theta$  by moving them in the direction of the negative gradient with respect to the **total loss**  $\mathcal{L}$  by the learning rate  $\eta$  multiplied by the gradient:

for each Epoch:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$$

## Reminder: Stochastic Gradient Descent

- Define loss function  $\ell$ , dataset  $D$  and model  $g$  with learnable parameters  $\theta$ .
- Define how many passes over the data to make (each one known as an Epoch)
- Define a learning rate  $\eta$

Stochastic Gradient Descent updates the parameters  $\theta$  by moving them in the direction of the negative gradient with respect to the loss of a **single item**  $\ell$  by the learning rate  $\eta$  multiplied by the gradient:

```
for each Epoch:  
    for each  $(x, y) \in D$ :  
         $\theta \leftarrow \theta - \eta \nabla_{\theta} \ell$ 
```