

# Multilayer Perceptrons: Expressiveness, overfitting, regularization

Marc Lelarge   Kevin Scaman   Jill-Jênn Vie

Oct 21, 2022

## Multilayer Perceptrons (MLP)

$$\mathbf{x}^{(0)} = \mathbf{x}$$

$$\mathbf{x}^{(\ell+1)} = \sigma(\mathbf{W}^{(\ell)} \mathbf{x}^{(\ell)} + \mathbf{b}^{(\ell)}) \quad \ell = 0, \dots, L-2$$

$$y = \mathbf{x}^{(L)} = \mathbf{W}^{(L-1)} \mathbf{x}^{(L-1)} + \mathbf{b}^{(L-1)}$$

# Multilayer Perceptrons (MLP)

$$\mathbf{x}^{(0)} = \mathbf{x} \in \mathbf{R}^{d_0}$$

$$\mathbf{x}^{(\ell+1)} = \sigma(\mathbf{W}^{(\ell)}\mathbf{x}^{(\ell)} + \mathbf{b}^{(\ell)}) \in \mathbf{R}^{d_\ell} \quad \ell = 0, \dots, L-2$$

$$y = \mathbf{x}^{(L)} = \mathbf{W}^{(L-1)}\mathbf{x}^{(L-1)} + \mathbf{b}^{(L-1)} \in \mathbf{R}^{d_L}$$

The  $\ell$ th layer has  $d_\ell$  neurons. Input layer  $\ell = 0$ , output layer  $\ell = L$ .  
 $\sigma$  is the link function. Usually,  $\sigma = \text{ReLU} = \max(\mathbf{0}, \mathbf{x})$ . #params?

# Multilayer Perceptrons (MLP)

$$\mathbf{x}^{(0)} = \mathbf{x} \in \mathbf{R}^{d_0}$$

$$\mathbf{x}^{(\ell+1)} = \sigma(\mathbf{W}^{(\ell)}\mathbf{x}^{(\ell)} + \mathbf{b}^{(\ell)}) \in \mathbf{R}^{d_\ell} \quad \ell = 0, \dots, L-2$$

$$y = \mathbf{x}^{(L)} = \mathbf{W}^{(L-1)}\mathbf{x}^{(L-1)} + \mathbf{b}^{(L-1)} \in \mathbf{R}^{d_L}$$

The  $\ell$ th layer has  $d_\ell$  neurons. Input layer  $\ell = 0$ , output layer  $\ell = L$ .  
 $\sigma$  is the link function. Usually,  $\sigma = \text{ReLU} = \max(\mathbf{0}, \mathbf{x})$ . #params?

```
1 from torch import nn
2
3 mlp = nn.Sequential(
4     nn.Linear(d0, d1),
5     nn.ReLU(),
6     nn.Linear(d1, d2),
7     ...
8     nn.Linear(dL-1, dL)
9 )
```

Example: logistic regression is a 1-layer perceptron

$$y = \sigma(\mathbf{w}^T \mathbf{x} + b) \text{ where } \sigma = \text{sigmoid} = 1/(1 + \exp(-x))$$

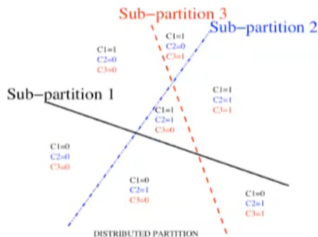
# Distributed Representations: The Power of Compositionality - Part 1



- Distributed (possibly sparse) representations, learned from data, can capture the **meaning** of the data and state
- Parallel composition of features: can be exponentially advantageous



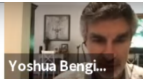
Not Distributed



Distributed

zoom

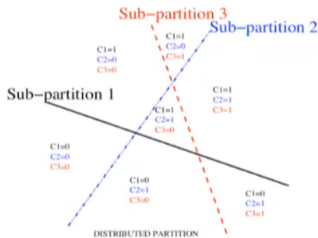
## Distributed Representations: The Power of Compositionality - Part 1



- Distributed (possibly sparse) representations, learned from data, can capture the **meaning** of the data and state
- Parallel composition of features: can be exponentially advantageous



Not Distributed



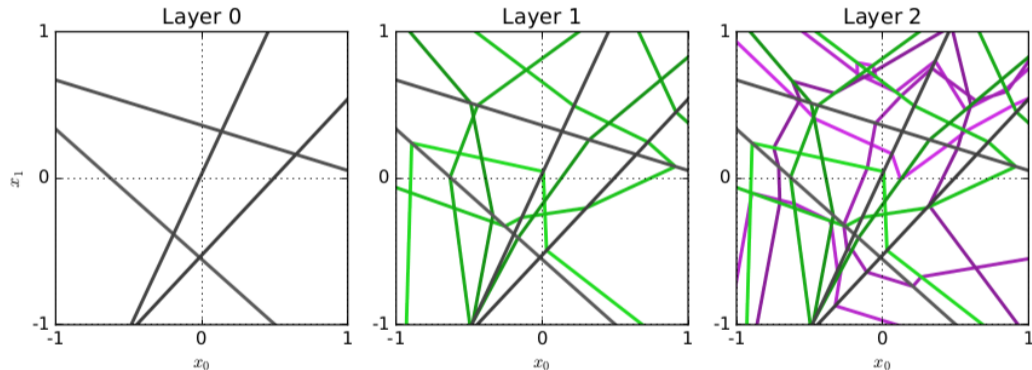
Distributed

zoom

A ReLU-based MLP with inputs in  $\mathbf{R}^n$ ,  $L$  layers of width  $k \geq n$ , can compute functions that have  $\Omega((k/n)^{n(L-1)} n^k)$  linear regions.

## Expressiveness

The number of activation patterns ( $\sim$  regions) of a ReLU-based MLP with  $L$  layers of width  $k$ , inputs in  $\mathbf{R}^n$  is upper bounded (tightly) by  $O(k^{nL})$  as  $L \rightarrow \infty$ .



Maithra Raghu et al. "On the expressive power of deep neural networks". In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2847–2854

# Universal approximation theorems

Fixed depth 2 arbitrary width  $k$  (Pinkus, 1999) (Cybenko, 1989)

Let  $\sigma \in \mathcal{C}(\mathbf{R})$  a continuous function from  $\mathbf{R}$  to  $\mathbf{R}$ .

Then:  $\sigma$  is not polynomial  $\iff$

For all  $\varepsilon > 0$ ,  $n, m \in \mathbf{N}$ , compact  $K \subseteq \mathbf{R}^n$ , function  $f \in \mathcal{C}(K, \mathbf{R}^m)$ ,  
there exist latent dimension  $k$  and weights  $\mathbf{W}, \mathbf{b}, \mathbf{C}$  such that

$$\sup_{\mathbf{x} \in K} \|f(\mathbf{x}) - \text{MLP}(\mathbf{x})\| < \varepsilon \quad \text{MLP}(\mathbf{x}) = \mathbf{C}\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}).$$

# Universal approximation theorems

Fixed depth 2 arbitrary width  $k$  (Pinkus, 1999) (Cybenko, 1989)

Let  $\sigma \in \mathcal{C}(\mathbf{R})$  a continuous function from  $\mathbf{R}$  to  $\mathbf{R}$ .

Then:  $\sigma$  is not polynomial  $\iff$

For all  $\varepsilon > 0$ ,  $n, m \in \mathbf{N}$ , compact  $K \subseteq \mathbf{R}^n$ , function  $f \in \mathcal{C}(K, \mathbf{R}^m)$ ,  
there exist latent dimension  $k$  and weights  $\mathbf{W}, \mathbf{b}, \mathbf{C}$  such that

$$\sup_{\mathbf{x} \in K} \|f(\mathbf{x}) - \text{MLP}(\mathbf{x})\| < \varepsilon \quad \text{MLP}(\mathbf{x}) = \mathbf{C}\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}).$$

In other words, 2-layer MLPs are **dense** in  $\mathcal{C}(K, \mathbf{R}^m)$ .

They are **universal approximators** of continuous functions.

# Universal approximation theorems

## Arbitrary depth, minimal width (Park, 2020)

For any function  $f \in L^p(\mathbf{R}^n, \mathbf{R}^m)$  and any  $\varepsilon > 0$   
there exists a MLP with ReLU of width  $\max(n+1, m)$  such that

$$\|f - \text{MLP}\|_p = \left( \int_{\mathbf{R}^n} \|f(x) - \text{MLP}(x)\|^p dx \right)^{1/p} < \varepsilon.$$

# Universal approximation theorems

## Arbitrary depth, minimal width (Park, 2020)

For any function  $f \in L^p(\mathbf{R}^n, \mathbf{R}^m)$  and any  $\varepsilon > 0$   
there exists a MLP with ReLU of width  $\max(n + 1, m)$  such that

$$\|f - \text{MLP}\|_p = \left( \int_{\mathbf{R}^n} \|f(x) - \text{MLP}(x)\|^p dx \right)^{1/p} < \varepsilon.$$

Moreover:

## Arbitrary depth, constrained width (Kidger and Lyons, 2020)

Let  $\mathcal{N}$  be the space of MLP :  $\mathbf{R}^n \rightarrow \mathbf{R}^m$  with any layers having  $n + m + 2$  neurons.  
Then:  $\mathcal{N}$  is **dense** in  $\mathcal{C}(K, \mathbf{R}^m)$  where compact  $K \subseteq \mathbf{R}^n$ .

## Some other theoretical results

### Infinite-depth limit

- ▶ Untrained MLP with random weights (Karakida, Akaho & Amari, 2018)

The Fisher information matrix i.e.  $\frac{\partial^2 \mathcal{L}}{\partial \theta^2}$  has eigenvalues having mean  $O(1/M)$ , variance  $O(1)$  and max  $O(M)$ .

- ▶ Neural Tangent Kernels (Jacot, Gabriel & Hongler, 2018)

### Turing-completeness

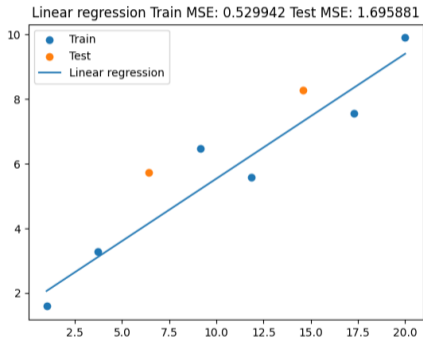
- ▶ RNNs are Turing-complete (Siegelmann & Sontag, 1995)
- ▶ LSTMs can perform unbounded counting while GRUs cannot (Weiss, Goldberg & Yahav, 2018)
- ▶ Neural Turing Machines with external memory (Graves, Wayne & Danihelka, 2014)

Training MLP with random labels?! (Maennel et al., 2020)

# Overfitting

## Fitting

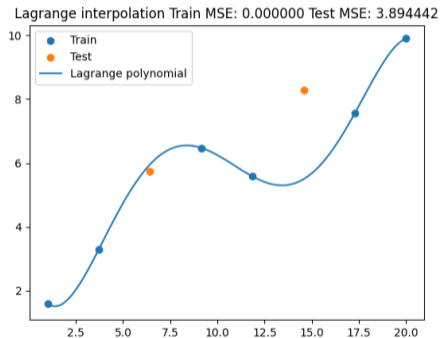
Polynomial degree 1  $Y = wX + b$



Linear regression  
`scipy.stats.linreg`

## Overfitting

Polynomial degree 6  $Y = \sum_{k=0}^6 w_k X^k$



Lagrange interpolation  
`scipy.interpolation.lagrange`

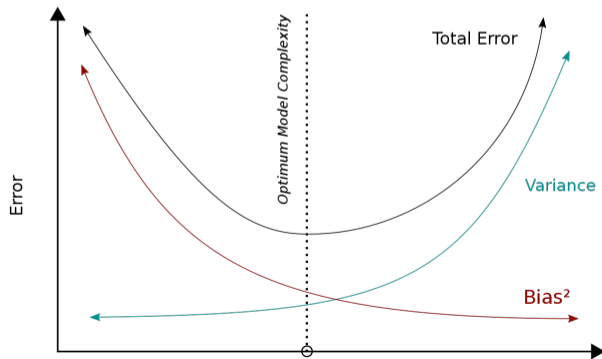
## Bias-variance decomposition

Samples  $x_i, y_i \in \mathbf{R}$ . We train a model  $g_\theta$ .

Let us assume that  $y_i = f(x_i) + \varepsilon_i$  where  $\varepsilon_i \in \mathcal{N}(0, \sigma^2)$ .

Then: the generalization error for the squared loss verifies

$$\begin{aligned}\mathbb{E}[(Y - g_\theta(X))^2] &= \text{Bias}(g_\theta)^2 + \text{Var}(g_\theta) + \sigma^2 \\ &= \mathbb{E}[g_\theta - f]^2 + \mathbb{E}[(g_\theta - \mathbb{E}g_\theta)^2] + \sigma^2.\end{aligned}$$



## Bias-variance decomposition

Samples  $x_i, y_i \in \mathbf{R}$ . We train a model  $g_\theta$ .

Let us assume that  $y_i = f(x_i) + \varepsilon_i$  where  $\varepsilon_i \in \mathcal{N}(0, \sigma^2)$ .

Then: the generalization error for the squared loss verifies

$$\begin{aligned}\mathbb{E}[(Y - g_\theta(X))^2] &= \text{Bias}(g_\theta)^2 + \text{Var}(g_\theta) + \sigma^2 \\ &= \mathbb{E}[g_\theta - f]^2 + \mathbb{E}[(g_\theta - \mathbb{E}g_\theta)^2] + \sigma^2.\end{aligned}$$

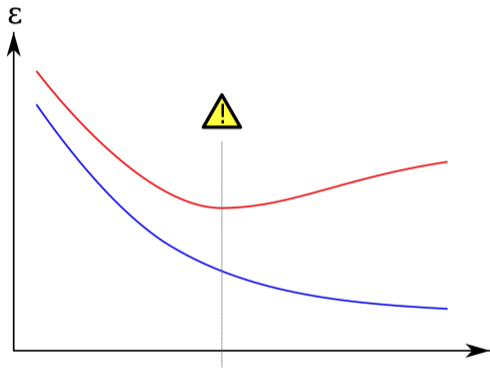
Proof.

$$\mathbb{E}f = f \quad \mathbb{E}Y = f \quad \text{Var}(Y) = \sigma^2$$

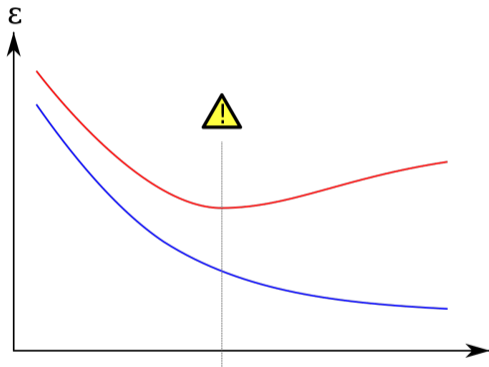
As  $\varepsilon$  and  $g_\theta$  are independent:

$$\begin{aligned}\mathbb{E}[(Y - g_\theta)^2] &= \mathbb{E}[Y^2 + g_\theta^2 - 2Yg_\theta] \\ &= \mathbb{E}[Y^2] + \mathbb{E}[g_\theta^2] - \mathbb{E}[2Yg_\theta] \\ &= \text{Var}(Y) + \mathbb{E}[Y]^2 + \text{Var}(g_\theta) + \mathbb{E}[g_\theta]^2 - 2f\mathbb{E}[g_\theta] \\ &= \text{Var}(Y) + \text{Var}(g_\theta) + (f - \mathbb{E}[g_\theta])^2 \\ &= \text{Var}(Y) + \text{Var}(g_\theta) + \mathbb{E}[f - g_\theta]^2 \\ &= \sigma^2 + \text{Var}(g_\theta) + \text{Bias}(g_\theta)^2.\end{aligned}$$

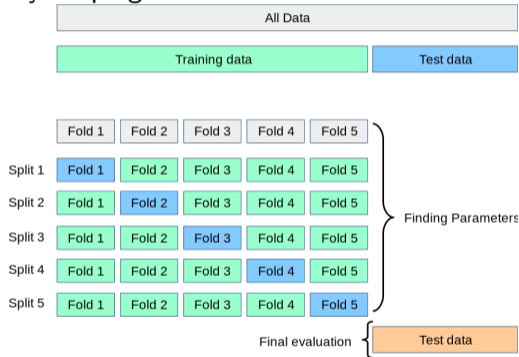
## How to detect overfitting?



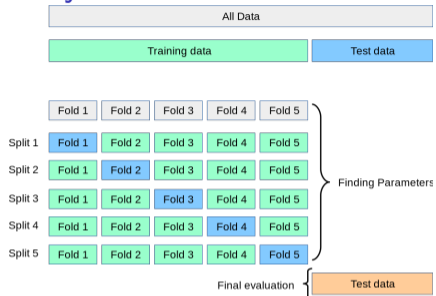
# How to detect overfitting?



By keeping a validation set

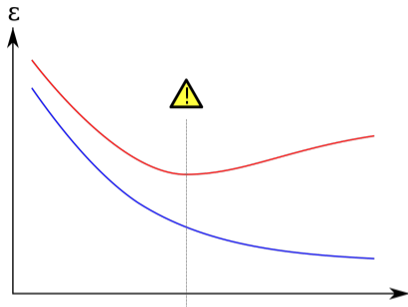


# Hyperparameter selection by cross-validation



```
1 trainval, test = split data into 80:20
2 train, valid = split trainval into 80:20
3
4 for each hyperparameter  $\lambda$ :
5     minimize error on train using  $\lambda$ 
6     valid_score $_{\lambda}$   $\leftarrow$  evaluate metric on valid
7  $\lambda^* \leftarrow \lambda$  achieving best valid_score $_{\lambda}$ 
8 minimize error on trainval using  $\lambda^*$  (= refit)
```

## Early stopping



```
1 def training_loop(train, valid):  
2     for each epoch:  
3         for x, y in train:  
4             do one step of gradient  
5             valid_score ← evaluate metric on valid  
6             if valid_score is worse than before:  
7                 return
```

## Example

Let us consider logistic regression i.e. 1-layer MLP:

$$f(\mathbf{x}_i) = \sigma(\mathbf{W}\mathbf{x}_i + b)$$

Logistic loss:  $\mathcal{L} = \sum_i -(1 - y_i) \log(1 - f(\mathbf{x}_i)) - y_i \log f(\mathbf{x}_i)$

If all samples have same target  $y_i = 1$  (or if there's only 1 sample), what will happen?

## Example

Let us consider logistic regression i.e. 1-layer MLP:

$$f(\mathbf{x}_i) = \sigma(\mathbf{W}\mathbf{x}_i + b)$$

Logistic loss:  $\mathcal{L} = \sum_i -(1 - y_i) \log(1 - f(\mathbf{x}_i)) - y_i \log f(\mathbf{x}_i)$

If all samples have same target  $y_i = 1$  (or if there's only 1 sample), what will happen?

MLP believes everything is a cat.

- ▶ Minimize  $-\log \sigma(\mathbf{W}\mathbf{x}_1 + b)$
- ▶  $\sigma(\mathbf{W}\mathbf{x}_1 + b) \rightarrow 1$
- ▶  $\mathbf{W}\mathbf{x}_1 + b \rightarrow +\infty$
- ▶ Parameters  $|\mathbf{W}|$  and  $b$  diverge to  $+\infty$

Add penalty to loss  $\|\mathbf{W}\|_2^2 + \|b\|_2^2$  (= assuming a Gaussian prior centered in  $\mathbf{0}$ ), called  $L_2$  regularization

# Regularize to generalize

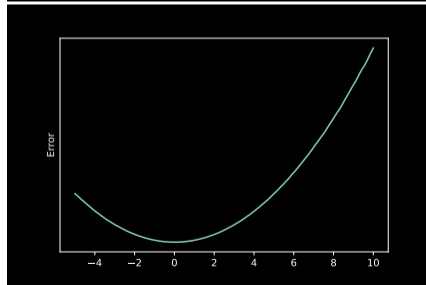
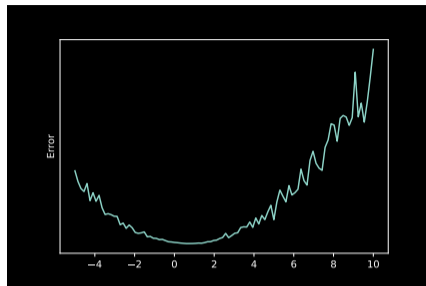
Minimizing loss:

May fall in local minima or diverge to  $\infty$

Minimizing loss + regularization:

Easier to optimize

We will see an example next week.



# Take home message

## Expressiveness

- ▶ 2-layer MLPs are universal approximators of continuous functions
- ▶ Try to overfit a single batch; otherwise your model cannot express the data.

## Bias-variance trade-off

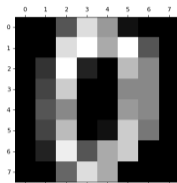
There is incompressible error due to inherent noise

## Overfitting

- ▶ Don't look at test data it's forbidden
- ▶ Implement early stopping
- ▶ And  $L_2$  regularization

# Today's practical: datasets

## Digits



$1797 \times 8 \times 8$  images representing numbers between 0 and 9.

Red Wine Quality<sup>1</sup> (Cortez et al., 2009)

1599 wines  $\times$  11 features<sup>2</sup>, have to predict quality which is an integer between 0 and 10 (in practice between 3 and 8).

Also: Faces, Cats and dogs

---

<sup>1</sup><https://kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>

<sup>2</sup>fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, alcohol

- [1] Guido F Montufar et al. “On the number of linear regions of deep neural networks”. In: *Advances in Neural Information Processing Systems*. Vol. 27. 2014, pp. 2924–2932.
- [2] Maithra Raghu et al. “On the expressive power of deep neural networks”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2847–2854.