# Multi-Layer Perceptrons and Deep Learning

Christos Dimitrakakis

October 23, 2024

# Outline

Features and layers

Algorithms
    Random projection
    Back propagation
    Derivatives
    Cost functions

Python libraries
    sklearn
    PyTorch
    TensorFlow

# Features and layers

## Algorithms
- Random projection
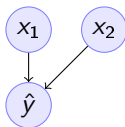- Back propagation
- Derivatives
- Cost functions

## Python libraries
- sklearn
- PyTorch
- TensorFlow

# Perceptron vs linear regression



- Network output
$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

- Chain rule
$$\nabla_\beta L = \nabla_{\hat{y}} L \nabla_\beta \hat{y}$$

- Network gradient
$$\nabla_\beta \hat{y} = (x_1, x_2)$$

## Cost functions
The only difference are the cost functions

- Perceptron
$$L = -\mathbb{I}\{y \neq \hat{y}\} \hat{y}$$
with
$$\nabla L = -\mathbb{I}\{y \neq \hat{y}\} yx$$

- Linear regression
$$L = (\hat{y} - y)^2,$$
with
$$\nabla_{\hat{y}} L = 2(\hat{y} - y).$$

# Layering and features

## Fixed layers

- Input to layer $x \in R^n$
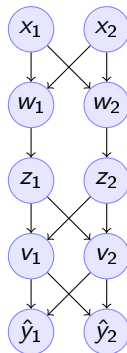- Output from layer $\hat{y} \in R^m$.

## Intermediate layers

- Linear layer
- Non-linear activation function.

## Linear layers types

- Dense
- Sparse
- Convolutional

## Activation funnction

- Sigmoid
- Softmax



| | |
|---|---|
| $x_1$ $x_2$ | Input layer |
| $w_1$ $w_2$ | Linear layer |
| $z_1$ $z_2$ | Sigmoid activation |
| $v_1$ $v_2$ | Linear layer |
| $\hat{y}_1$ $\hat{y}_2$ | Softmax activation |

# Linear layers

## Example: Linear regression with $n$ inputs, $m$ outputs.

- Input: Features $x \in \mathbb{R}^n$
- Dense linear layer with $B \in \mathbb{R}^{m \times n}$
- Output: $\hat{y} \in \mathbb{R}^m$

## Dense linear layer

- Parameters $B = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}$,
- $\beta_i = [\beta_{i,1}, \ldots, \beta_{i,n}]$, $\beta_i$ connects the $i$-th output $y_i$ to the features $x$:

$$y_i = \beta_i x$$

- In compact form:

$$y = Bx$$

# ReLU layers

▶ Typically used in the hidden layers of neural networks

$$f(x) = \max(0, x)$$

## Derivative

$$df/dxf(x) = \mathbb{I}\{x > 0\}$$
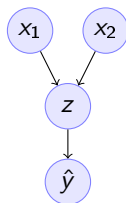
# Sigmoid activation

## Example: Logistic regression

- Input $x \in \mathbb{R}^n$
- Intermediate output: $z \in \mathbb{R}$,

$$z = \sum_{i=1}^{n} \beta_i x_i.$$

- Output: sigmoid activation $\hat{y} \in [0, 1]$.

$$f(z) = 1/[1 + \exp(-z)].$$

Now we can interpret $\hat{y} = P_{\boldsymbol{\beta}}(y = 1 | x)$.



Input layer

Linear layer

Sigmoid layer

## Loss function: negative log likelihood

$$\ell(\hat{y}, y) = -[\mathbb{I}\{y = 1\} \ln(\hat{y}) + \mathbb{I}\{y = -1\} \ln(1 - \hat{y})]$$

# Softmax layer

Example: Multivariate logistic regression with $m$ classes.

- Input: Features $\boldsymbol{x} \in \mathbb{R}^n$
- Fully-connected linear activation layer
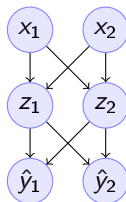$$z = Bx, \qquad B \in \mathbb{R}^{m \times n}.$$

- Softmax output
$$\hat{y}_i = \frac{\exp(z_i)}{\sum_{j=1^m} \exp(z_j)}$$

We can also interpret this as

$$\hat{y}_i \triangleq \mathbb{P}(y = i \mid \boldsymbol{x})$$

with usual loss $\ell(\hat{y}, y) = -\ln \hat{y}_y$



Input layer

Linear layer

Softmax layer

Features and layers

Algorithms
    Random projection
    Back propagation
    Derivatives
    Cost functions

Python libraries
    sklearn
    PyTorch
    TensorFlow

# Random projections

- ▶ Features $x$
- ▶ Hidden layer activation $z$
- ▶ Output $y$

## Hidden layer: Random projection

Here we project the input into a high-dimensional space

$$z_i = \text{sgn}(\beta_i^\top x) = y_i$$

where $\boldsymbol{B} = [\boldsymbol{\beta}_i]_{i=1}^m$, $\beta_{i,j} \sim \text{Normal}(0, 1)$

## The reason for random projections

- ▶ The high dimension makes it easier to learn.
- ▶ The randomness ensures we are not learning something spurious.

# Background on back-propagation

## The problem

- ▶ We need to minimise a loss function $L$
- ▶ We need to calculate

$$\nabla_{\boldsymbol{\beta}} \mathbb{E}_{\boldsymbol{\beta}}[L] \approx \frac{1}{T} \sum_{t=1}^{T} \nabla_{\boldsymbol{\beta}} \ell(x_t, y_t, \boldsymbol{\beta}).$$

- ▶ However $\ell(x_t, y_t, \boldsymbol{\beta})$ is a complex non-linear function of $\boldsymbol{\beta}$.
- ▶ We need many steps to calculate $\ell$. How can we then do it?

# The chain rule of differentiation



[1673] Liebniz

# Chain rule applied to the perceptron



[1976] Rosenblat

# Chain rule for deep neural netowrks



[1982] Werbos

# Backpropagation given a name
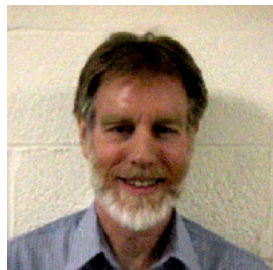
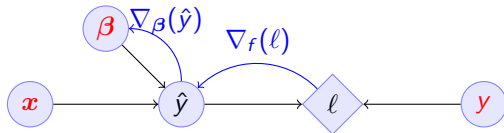1986: Learning representations by back-propagating errors.



Rumelhart



Hinton



Williams

# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}}\ell = \nabla_{\boldsymbol{\beta}}f\nabla_{\hat{y}}\ell$
- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$
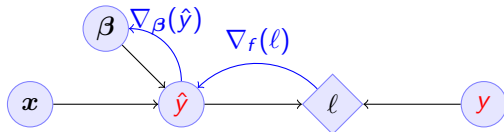
# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}} \ell = \nabla_{\boldsymbol{\beta}} f \nabla_{\hat{y}} \ell$
- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$

# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}}\ell = \nabla_{\boldsymbol{\beta}} f \nabla_{\hat{y}} \ell$
- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$
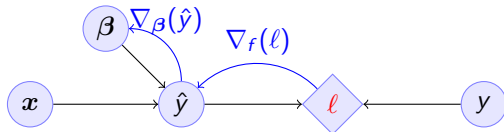
# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}}\ell = \nabla_{\boldsymbol{\beta}}f\nabla_{\hat{y}}\ell$
- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$

- Backward: return to calculate the gradients

$$\nabla_{\boldsymbol{\beta}}\ell(\hat{y}, y) = \nabla_{\boldsymbol{\beta}}f(\boldsymbol{\beta}, \boldsymbol{x}) \times \nabla_{\hat{y}}\ell(\hat{y}, y) \tag{1}$$

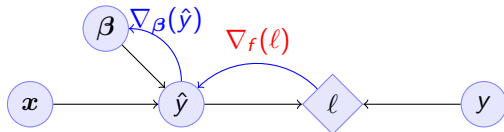$$= \nabla_{\boldsymbol{\beta}}f(\boldsymbol{\beta}, \boldsymbol{x}) \times 2[\hat{y} - y] \tag{2}$$

# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}} \ell = \nabla_{\boldsymbol{\beta}} f \nabla_{\hat{y}} \ell$
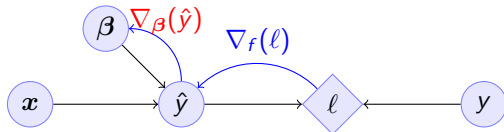- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$

- Backward: return to calculate the gradients

$$\nabla_{\boldsymbol{\beta}} \ell(\hat{y}, y) = \nabla_{\boldsymbol{\beta}} f(\boldsymbol{\beta}, \boldsymbol{x}) \times \nabla_{\hat{y}} \ell(\hat{y}, y) \tag{1}$$
$$= \nabla_{\boldsymbol{\beta}} f(\boldsymbol{\beta}, \boldsymbol{x}) \times 2[\hat{y} - y] \tag{2}$$

# Elementary back-propagation: linear regression



- $f : X \to Y$, $\ell : Y \times Y \to \mathbb{R}$, chain rule: $\nabla_{\boldsymbol{\beta}} \ell = \nabla_{\boldsymbol{\beta}} f \nabla_{\hat{y}} \ell$
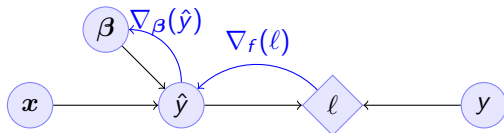- Forward: follow the arrows to calculate variables

$$\hat{y} \triangleq f(\boldsymbol{\beta}, x) = \sum_{i=1}^{n} \beta_i x_i, \qquad \ell(\hat{y}, y) = (\hat{y} - y)^2$$

- Backward: return to calculate the gradients

$$\nabla_{\boldsymbol{\beta}} \ell(\hat{y}, y) = \nabla_{\boldsymbol{\beta}} f(\boldsymbol{\beta}, \boldsymbol{x}) \times \nabla_{\hat{y}} \ell(\hat{y}, y) \tag{1}$$
$$= \nabla_{\boldsymbol{\beta}} f(\boldsymbol{\beta}, \boldsymbol{x}) \times 2[\hat{y} - y] \tag{2}$$

- Update:

$$\boldsymbol{\beta}_{t+1} = \boldsymbol{\beta}_t - \alpha_t \times \nabla_{\boldsymbol{\beta}} \ell(\hat{y}_t, y_t)$$

# Gradient descent with *back-propagation*

- Dataset $D$, cost function $L = \sum_t \ell_t$
- Parameters $\boldsymbol{B}_1, \ldots, \boldsymbol{B}_k$ with $k$ layers
- Intermediate variables: $\boldsymbol{z}_j = h_j(\boldsymbol{z}_{j-1}, \boldsymbol{B}_j)$, $\boldsymbol{z}_0 = \boldsymbol{x}$, $\boldsymbol{z}_k = \hat{\boldsymbol{y}}$.

# Gradient descent with *back-propagation*

- Dataset $D$, cost function $L = \sum_t \ell_t$
- Parameters $\boldsymbol{B}_1, \ldots, \boldsymbol{B}_k$ with $k$ layers
- Intermediate variables: $\boldsymbol{z}_j = h_j(\boldsymbol{z}_{j-1}, \boldsymbol{B}_j)$, $\boldsymbol{z}_0 = \boldsymbol{x}$, $\boldsymbol{z}_k = \hat{\boldsymbol{y}}$.
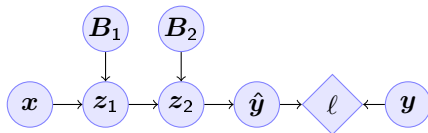
## Dependency graph

# Gradient descent with *back-propagation*

- Dataset $D$, cost function $L = \sum_t \ell_t$
- Parameters $\boldsymbol{B}_1, \ldots, \boldsymbol{B}_k$ with $k$ layers
- Intermediate variables: $\boldsymbol{z}_j = h_j(\boldsymbol{z}_{j-1}, \boldsymbol{B}_j)$, $\boldsymbol{z}_0 = \boldsymbol{x}$, $\boldsymbol{z}_k = \hat{\boldsymbol{y}}$.

## Dependency graph



## Backpropagation with steepest stochastic gradient descent

- Forward step: For $j = 1, \ldots, k$, calculate $\boldsymbol{z}_j = h_j(k)$ and $\ell(\hat{\boldsymbol{y}}, \boldsymbol{y})$
- Backward step: Calculate $\nabla_{\hat{\boldsymbol{y}}}\ell$ and $d_j \triangleq \nabla_{\boldsymbol{B}_j}\ell = \nabla_{\boldsymbol{B}_j}z_j d_{j+1}$ for $j = k \ldots, 1$
- Apply gradient: $\boldsymbol{B}_j -= \alpha d_j$.

# Other algorithms and gradients

## Natural gradient
Defined for probabilistic models

## ADAM
Exponential moving average of gradient and square gradients

## BFGS: Broyden–Fletcher–Goldfarb–Shanno algorithm
Newton-like method

# Linear layer

## Definition

This is a linear combination of inputs $x \in \mathbb{R}^n$ and parameter matrix $\boldsymbol{B} \in \mathbb{R}^{m \times n}$

where $\boldsymbol{B} = \begin{bmatrix} \boldsymbol{\beta_1} \\ \vdots \\ \boldsymbol{\beta_i} \\ \vdots \\ \boldsymbol{\beta_m} \end{bmatrix} = \begin{bmatrix} \beta_{1,1} & \cdots & \beta_{1,j} & \cdots & \beta_{1,m} \\ \vdots & \ddots & \vdots & \ddots & \cdots \\ \beta_{i,1} & \cdots & \beta_{i,j} & \cdots & \beta_{i,m} \\ \vdots & \ddots & \ddots & \ddots & \cdots \\ \beta_{n,1} & \cdots & \beta_{i,j} & \cdots & \beta_{n,m} \end{bmatrix}$

$$f(\boldsymbol{B}, \boldsymbol{x}) = \boldsymbol{B}\boldsymbol{x} \qquad f_i(\boldsymbol{B}, \boldsymbol{x}) = \boldsymbol{\beta_i} \cdot \boldsymbol{x} = \sum_{j=1}^{n} \beta_{i,j} x_j,$$

## Gradient

Each partial derivative is simple:

$$\frac{\partial}{\partial \beta_{i,j}} f_k(\boldsymbol{B}, \boldsymbol{x}) = \sum_{k=1}^{n} \frac{\partial}{\partial \beta_{i,j}} \beta_{i,k} x_k = x_j$$

# Sigmoid layer

▶ This layer is used for binary classification.
▶ It is used in the logistic regression model to obtain label probabilities.

$$f(z) = 1/(1 + \exp(-z))$$

## Derivative
So let us ignore the other inputs for simplicity:

$$\frac{d}{dz} f(z) = \exp(-z)/[1 + \exp(-z)]^2$$

# Softmax layer

▶ This layer is used for multi-class classification

$$y_i(z) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

## Derivative

$$\frac{\partial}{\partial z_i} y_i(z) = \frac{e^{z_i} e^{\sum_{j \neq i} z_j}}{\left(\sum_j e^{z_j}\right)^2}$$
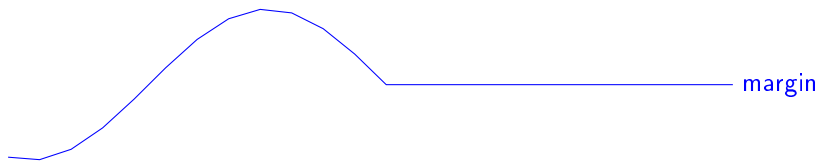
$$\frac{\partial}{\partial z_i} y_k(z) = \frac{e^{z_i + z_k}}{\left(\sum_j e^{z_j}\right)^2}$$

# Classification cost functions

## Error margin

If $z$ is a confidence level for the positive class then

$$\ell(z, y) = -\mathbb{I}\{zy < 0\}\, zy$$



margin

## Negative log likelihood (aka cross-entropy)

If $z$ are label probabilities, then

$$\ell(z, y) = -\ln z_y.$$

# Regression cost functions

## Squared error

If $z$ is a prediction for the dependent variable then

$$\ell(z, y) = (y - z)^2$$

This also corresponds to negative log likelihood under a Gaussianity assumption.

## Huber loss

If $z$ is a prediction, then

$$\ell(z, y) = \begin{cases} \frac{1}{2}(z - y)^2 & |z - y| \leq \delta \\ \delta(|z - y| - \frac{1}{2}\delta) & \text{otherwise.} \end{cases} \tag{3}$$

Features and layers

Algorithms
  Random projection
  Back propagation
  Derivatives
  Cost functions

Python libraries
  sklearn
  PyTorch
  TensorFlow

# sklearn neural networks

## Classification
Uses the cross entropy cost

```
from sklearn.neural_network import MLPClassifier
clf = MLPClassifier(hidden_layer_sizes=(5, 2))
clf.fit(X, y)
clf.predict(X_test)
```

▶ Main condition is layer sizes.

## Regression

```
from sklearn.neural_network import MLPRegressor
model = MLPRegressor(hidden_layer_sizes=(5, 2))
```

# Datasets

```
X_train = torch.tensor(X_train, dtype=torch.float32)
train_dataset = TensorDataset(X_train, y_train)
train_loader = DataLoader(train_dataset, batch_size=16, shuffle=True
mlp = nn.Sequential(nn.Linear(input_size , 50), nn.ReLU())
```

# TensorFlow

This is another library, no need to use this for this course