

MSc Big Data for Business - *MAP 534*

Introduction to machine learning

Supervised classification

Logistic regression & feed forward neural networks

Introduction

Logistic regression

- Model

- Inference

- Intepretation

- prediction

Feed Forward Neural Networks

- Model

- Implementation

Coronary heart disease example

→ **Sample of males in a heart-disease high-risk region** (Western Cape, South Africa).

→ **Many covariates.**

sbp, systolic blood pressure.

tobacco, cumulative tobacco (kg).

ldl, low density lipoprotein, bad cholesterol.

famhist, family history of heart disease.

typea, type-A behavior.

alcohol, current alcohol consumption.

age, age.

chd, response, coronary heart disease.

Analysis

Understand which factors are linked to the disease - **prevention**.

- Importance of the effect.
- Positive or negative effect.
- Significativity.

Prediction

Predict, for a new patient, the risk to declare the disease - **monitoring**.

- Quality of the prediction.
- Parsimonious model.

Credit Default, Credit Score, Bank Risk, Market Risk Management

- **Data**: client profile, client credit history...
- **Input**: client profile.
- **Output**: credit risk.

Analysis

Understand which factors in this dataset are linked to the default - **customer analysis**.

- Importance of the effect.
- Positive or negative effect.
- Significativity.

Prediction

Predict, for a new client, the risk to default - **decision aid**.

- Quality of the prediction.
- Parsimonious model.

Introduction

Logistic regression

- Model

- Inference

- Intepretation

- prediction

Feed Forward Neural Networks

- Model

- Implementation

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

Logistic regression

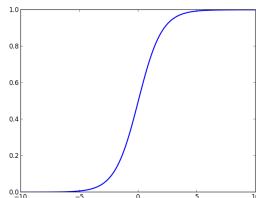
- The objective is to **predict the label** $Y \in \{0, 1\}$ based on $X \in \mathbb{R}^d$.
- Logistic regression **models the distribution of Y given X** .

The model

$$\mathbb{P}(Y = 1|X) = \sigma(\langle w, X \rangle + b),$$

where $w \in \mathbb{R}^d$ is a vector of model **weights** and $b \in \mathbb{R}$ is the **intercept**, and where σ is the **sigmoid** function.

$$\sigma : z \mapsto \frac{1}{1 + e^{-z}}.$$



- The sigmoid function is a **model choice to map \mathbb{R} into $(0, 1)$** .
- Another widespread solution for σ is $\sigma : z \mapsto \mathbb{P}(Z \leq z)$ where $Z \sim \mathcal{N}(0, 1)$, which leads to a **probit** regression model.

Log-odd ratio

$$\log \left(\mathbb{P}(Y = 1|X) \right) - \log \left(\mathbb{P}(Y = 0|X) \right) = \langle w, X \rangle + b.$$

Classification rule

Note that

$$\mathbb{P}(Y = 1|X) \geq \mathbb{P}(Y = 0|X)$$

if and only if

$$\langle w, x \rangle + b \geq 0.$$

→ This is a **linear classification** rule.

→ This classifier requires to **estimate w and b** .

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

→ $\{(X_i, Y_i)\}_{1 \leq i \leq n}$ are **i.i.d. with the same distribution as (X, Y)** .

Likelihood

$$\begin{aligned}\prod_{i=1}^n \mathbb{P}(Y_i | X_i) &= \prod_{i=1}^n \sigma(\langle w, X_i \rangle + b)^{Y_i} (1 - \sigma(\langle w, X_i \rangle + b))^{1-Y_i}, \\ &= \prod_{i=1}^n \sigma(\langle w, x_i \rangle + b)^{Y_i} \sigma(-\langle w, X_i \rangle - b)^{1-Y_i}\end{aligned}$$

and the **normalized negative loglikelihood** is

$$f(w, b) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-Y_i(\langle w, X_i \rangle + b)}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \langle w, X_i \rangle + b).$$

Logistic regression

Compute \hat{w}_n and \hat{b}_n as follows:

$$(\hat{w}_n, \hat{b}_n) \in \operatorname{argmin}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-Y_i(\langle w, X_i \rangle + b)}).$$

→ It is an **average of losses**, one for each sample point.

→ It is a **convex and smooth problem**.

Using the **logistic loss** function

$$\ell : (y, y') \mapsto \log(1 + e^{-yy'})$$

yields

$$(\hat{w}_n, \hat{b}_n) \in \operatorname{argmin}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \langle w, X_i \rangle + b).$$

Maximum likelihood estimate

Assume for now that the intercept is 0. Then, the likelihood is,

$$L_n(w) = \prod_{i=1}^n \left(\frac{e^{X_i^T w}}{1 + e^{X_i^T w}} \right)^{Y_i} \left(\frac{1}{1 + e^{X_i^T w}} \right)^{1-Y_i} = \prod_{i=1}^n \left(\frac{e^{X_i^T w Y_i}}{1 + e^{X_i^T w}} \right).$$

And the **negative log-likelihood** is

$$\ell_n(w) = -\log(L_n(w)) = \sum_{i=1}^n \left(-Y_i X_i^T w + \log(1 + e^{X_i^T w}) \right).$$

Derivatives

$$\begin{aligned} \frac{\partial (\log(L_n(w)))}{\partial w_j} &= \sum_{i=1}^n \left(Y_i X_{ij} - \frac{x_{ij} e^{X_i^T w}}{(1 + e^{X_i^T w})} \right) \\ &= \sum_{i=1}^n X_{ij} (Y_i - \sigma(\langle w, X_i \rangle)). \end{aligned}$$

→ **No explicit solution** for the maximizer of the loglikelihood... Parameter estimate obtained using **gradient based optimization** (see next lesson).

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

The **gradient of the negative loglikelihood** is,

$$\nabla \ell_n(w) = - \sum_{i=1}^n Y_i X_i + \sum_{i=1}^n \frac{\exp(\langle X_i, w \rangle)}{1 + \exp(\langle X_i, w \rangle)} X_i .$$

On the other hand, for all $1 \leq i \leq n$ and all $1 \leq j \leq d$,

$$\partial_j \left(\frac{\exp(\langle X_i, w \rangle)}{1 + \exp(\langle X_i, w \rangle)} X_i \right) = \frac{\exp(\langle X_i, w \rangle)}{(1 + \exp(\langle X_i, w \rangle))^2} X_{ij} X_i ,$$

where X_{ij} is the j th component of X_i .

Then, the **Hessian matrix** is

$$(H_n(w))_{\ell j} = \sum_{i=1}^n \frac{\exp(\langle X_i, w \rangle)}{(1 + \exp(\langle X_i, w \rangle))^2} X_{ij} X_{i\ell} ,$$

that is,

$$H_n(w) = \sum_{i=1}^n \frac{\exp(\langle X_i, w \rangle)}{(1 + \exp(\langle X_i, w \rangle))^2} X_i X_i^T .$$

$H_n(\beta)$ is a semi positive definite matrix, which implies that $\ell_n(\beta)$ is convex.

Asymptotic properties

Assumptions

→ $\hat{w}_n \rightarrow w^*$ **almost surely**.

→ There exists a continuous and nonsingular function H such that $n^{-1}H_n(w)$ **converges to $H(w)$** , uniformly in a ball around w^* .

For all $t \in \mathbb{R}^d$, using a Taylor expansion,

$$\mathbb{E} \left[\exp \left(-\frac{1}{\sqrt{n}} \langle t, \nabla \ell_n(w^*) \rangle \right) \right] \rightarrow_{n \rightarrow \infty} \exp \left(\frac{1}{2} t^T H(w^*) t \right).$$

Therefore,

$$-\nabla \ell_n(w^*) / \sqrt{n} \Rightarrow \mathcal{N}(0, H(w^*)).$$

On the other hand, by **Slutsky lemma**,

$$\sqrt{n}(\hat{w}_n - w^*) \Rightarrow \mathcal{N}(0, H(w^*)^{-1}).$$

Confidence interval

$\rightarrow \sqrt{n}(\hat{w}_j - w_j^*)$ converges in distribution to a centered Gaussian random variable with variance $(H(w^*)^{-1})_{jj}$.

Almost surely,

$$\hat{\sigma}_{n,j}^2 = (nH_n(\hat{w}_n)^{-1})_{jj} \rightarrow_{n \rightarrow \infty} (H(w^*)^{-1})_{jj}.$$

Then,

$$\sqrt{\frac{n}{\hat{\sigma}_{n,j}^2}}(\hat{w}_{n,j} - \beta_j^*) \rightarrow_{n \rightarrow \infty} \mathcal{N}(0, 1).$$

An asymptotic confidence interval $\mathcal{I}_{n,\alpha}$ of level $1 - \alpha$ is then

$$\mathcal{I}_{n,\alpha} = \left[\hat{w}_{n,j} - z_{1-\alpha/2} \sqrt{\frac{\hat{\sigma}_{n,j}^2}{n}}, \hat{\beta}_{n,j} + z_{1-\alpha/2} \sqrt{\frac{\hat{\sigma}_{n,j}^2}{n}} \right],$$

where $z_{1-\alpha/2}$ is the quantile of order $1 - \alpha/2$ of $\mathcal{N}(0, 1)$.

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

- The objective is to **predict the label** $Y \in \{1, \dots, M\}$ based on $X \in \mathbb{R}^d$.
- Softmax regression **models the distribution of Y given X** .

The model

For all $1 \leq m \leq M$,

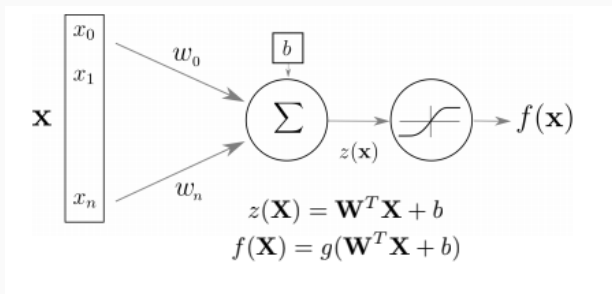
$$z_m = \langle w_m, X \rangle + b_m,$$

$$\mathbb{P}(Y = m|X) = \text{softmax}(z)_m,$$

where $z \in \mathbb{R}^M$, $w_m \in \mathbb{R}^d$ is a vector of model **weights** and $b_m \in \mathbb{R}$ is an **intercept**, and where softmax is the **softmax** function: for all $1 \leq m \leq M$,

$$\text{softmax}(z)_m = \frac{\exp(z_m)}{\sum_{j=1}^M \exp(z_j)}.$$

A neuron



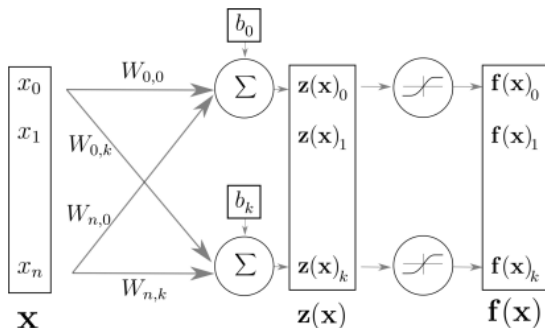
→ \mathbf{X} **input in \mathbb{R}^d** .

→ $z(\mathbf{X})$ **pre-activation in \mathbb{R}^M** , with **weight $\mathbf{W} \in \mathbb{R}^{d \times M}$** and **bias $b \in \mathbb{R}^M$** .

→ g **softmax function**.

On neuron is a multi-class extension of the logistic regression model.

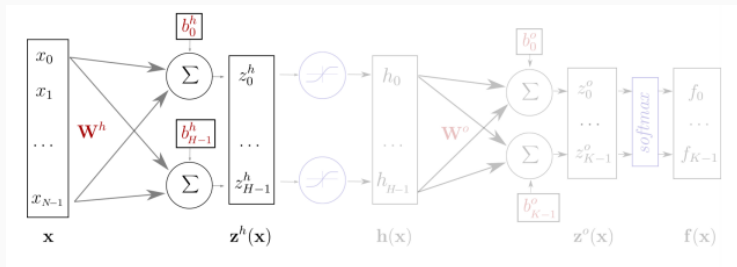
Layer of neurons and hidden states



$$f(\mathbf{X}) = g(\mathbf{z}(\mathbf{X})) = g(\mathbf{W}\mathbf{X} + \mathbf{b})$$

- \mathbf{X} **input** in \mathbb{R}^d .
- $\mathbf{z}(\mathbf{X})$ **pre-activation** in \mathbb{R}^k , with **weight** $\mathbf{W} \in \mathbb{R}^{d \times k}$ and **bias** $\mathbf{b} \in \mathbb{R}^k$.
- g **any activation function** (nonlinear & nondecreasing function).
- $\mathbf{f}(\mathbf{X})$ **hidden state** in \mathbb{R}^k which may be used as input of a new neuron...

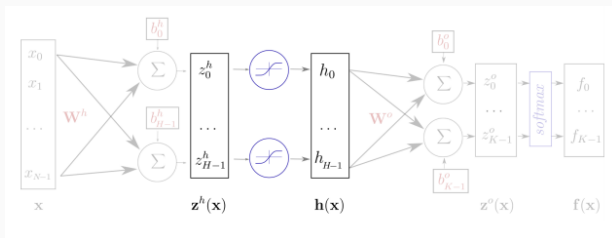
Feed Forward Network



→ \mathbf{X} **input** in \mathbb{R}^d .

→ $z^h(\mathbf{X})$ **pre-activation** in \mathbb{R}^H , with **weight** $\mathbf{W}^h \in \mathbb{R}^{d \times H}$ and **bias** $\mathbf{b}^h \in \mathbb{R}^H$.

Feed Forward Network

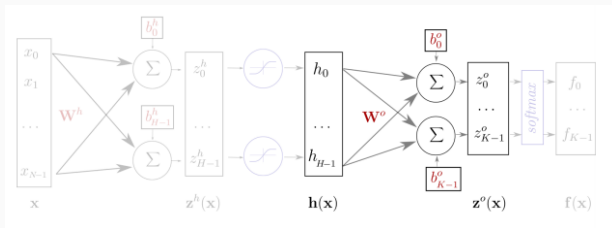


→ X input in \mathbb{R}^d .

→ $z^h(X)$ pre-activation in \mathbb{R}^H , with weight $W^h \in \mathbb{R}^{d \times H}$ and bias $b^h \in \mathbb{R}^H$.

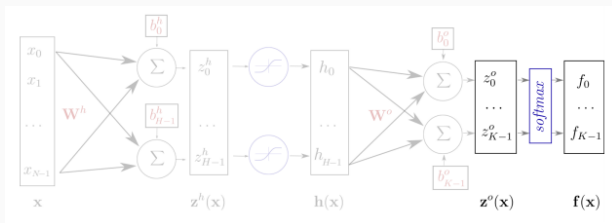
→ g any activation function to produce $h \in \mathbb{R}^H$.

Feed Forward Network



- \mathbf{X} **input** in \mathbb{R}^d .
- $z^h(\mathbf{X})$ **pre-activation** in \mathbb{R}^H , with **weight** $W^h \in \mathbb{R}^{d \times H}$ and **bias** $b^h \in \mathbb{R}^H$.
- g **any activation function** to produce $h \in \mathbb{R}^H$.
- $z^o(\mathbf{X})$ **pre-activation** in \mathbb{R}^M , with **weight** $W^o \in \mathbb{R}^{H \times M}$ and **bias** $b^o \in \mathbb{R}^M$.

Feed Forward Network



→ X **input** in \mathbb{R}^d .

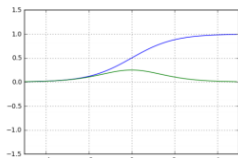
→ $z^h(X)$ **pre-activation** in \mathbb{R}^H , with **weight** $W^h \in \mathbb{R}^{d \times H}$ and **bias** $b^h \in \mathbb{R}^H$.

→ g **any activation function** to produce $h \in \mathbb{R}^H$.

→ $z^o(X)$ **pre-activation** in \mathbb{R}^M , with **weight** $W^o \in \mathbb{R}^{H \times M}$ and **bias** $b^o \in \mathbb{R}^M$.

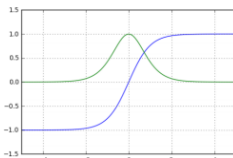
→ Apply the **softmax function** to produce the output, i.e. $\mathbb{P}(Y = m|X)$ for $1 \leq m \leq M$.

Activation functions



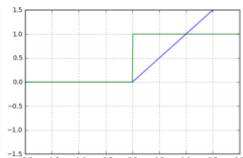
$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x))$$



$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\tanh'(x) = 1 - \tanh(x)^2$$



$$\text{relu}(x) = \max(0, x)$$

$$\text{relu}'(x) = 1_{x>0}$$

→ As there is no modelling assumptions anymore, **virtually any activation function** may be used.

→ The rectified linear unit (ReLU) activation function $\sigma(x) = \max(0, x)$ and its extensions are the default recommendation in modern implementations (Jarrett et al., 2009; Nair and Hinton, 2010; Glorot et al., 2011a), (Maas et al., 2013), (He et al., 2015). One of the major motivations arise from the **gradient based parameter optimization which is numerically more stable with this choice**.

Introduction

Logistic regression

Model

Inference

Intepretation

prediction

Feed Forward Neural Networks

Model

Implementation

→ This dataset contains images representing handwritten digits. Each image is made of **28 x 28 pixels**, and each pixel is represented by an integer (gray level). These arrays can be flattened into vectors in \mathbb{R}^{784} .

→ The **labels in $\{0, \dots, 9\}$ are represented using one-hot-encoding** and grayscale of each pixel in $\{0, \dots, 255\}$ are normalized to be in $(0, 1)$.

```
from keras.datasets import mnist
# Number of classes
num_classes = 10
# input image dimensions
img_rows, img_cols = 28, 28

# the data, shuffled and split between train and test sets
(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train = x_train.reshape(x_train.shape[0], img_rows, img_cols, 1)
x_test = x_test.reshape(x_test.shape[0], img_rows, img_cols, 1)
input_shape = (img_rows, img_cols, 1)

x_train = x_train.astype('float32')
x_test = x_test.astype('float32')

print('x_train shape:', x_train.shape)
print('x_test shape:', x_test.shape)
print('y_train shape:', y_train.shape)
print('y_test shape:', y_test.shape)
print(x_train.shape[0], 'train samples')
print(x_test.shape[0], 'test samples')
```

```
x_train shape: (60000, 28, 28, 1)
x_test shape: (10000, 28, 28, 1)
y_train shape: (60000,)
y_test shape: (10000,)
60000 train samples
10000 test samples
```


- This dataset contains images representing handwritten digits. Each image is made of **28 x 28 pixels**, and each pixel is represented by an integer (gray level). These arrays can be flattened into vectors in \mathbb{R}^{784} .
- The **labels in $\{0, \dots, 9\}$ are represented using one-hot-encoding** and grayscale of each pixel in $\{0, \dots, 255\}$ are normalized to be in $(0, 1)$.

```
plt.figure(figsize=(8, 2))
for i in range(4):
    plt.subplot(1, 4, i+1)
    plt.imshow(x_train[i].reshape(28, 28),
               interpolation="none", cmap="gray_r")
    plt.title('Label = %d' % y_train[i], fontsize=14)
    plt.axis("off")
plt.tight_layout()
```

Label = 5



Label = 0



Label = 4



Label = 1



The model with Keras

```
model_ffnn = Sequential()

model_ffnn.add(Flatten(input_shape=input_shape))

model_ffnn.add(Dense(128, activation='relu'))

model_ffnn.add(Dense(num_classes, activation='softmax'))
```

```
model_ffnn.compile(
    loss=keras.losses.categorical_crossentropy,
    optimizer=keras.optimizers.Adagrad(),
    metrics=['accuracy']
)

model_ffnn.summary()
```

Layer (type)	Output Shape	Param #
flatten_1 (Flatten)	(None, 784)	0
dense_1 (Dense)	(None, 128)	100480
dense_2 (Dense)	(None, 10)	1290
Total params: 101,770		
Trainable params: 101,770		
Non-trainable params: 0		

Figure 1: Feed Forward Neural network. h_1 is obtained with the **RELU activation function and is in \mathbb{R}^{128}** . The last layer is $h_2 \in \mathbb{R}^{10}$ and is **obtained with the softmax activation function** so that each component m models $\mathbb{P}(Y = m|X)$. This neural network with one hidden layer relies on 101.770 parameters.

→ This model relies on more than **100.000 unknown parameters** which should be estimated.

→ As for the logistic regression and the discriminant analysis, a common choice is to **minimize the negative loglikelihood of the data**:

$$\theta \mapsto -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{10} \mathbb{1}_{Y_i=k} \log \mathbb{P}_{\theta}(Y_i = k | X_i).$$

→ The negative loglikelihood is computed using $n = 60.000$ training samples and **minimized using gradient descent algorithms** - see next lesson.

→ Then, the performance of the model is assessed using 10.000 new (test) samples: the **accuracy is the frequency of labels which are well predicted by the model with the estimated parameters**.

```

batch_size = 32
epochs = 8

# Run the train
history = model_ffnn.fit(x_train, y_train,
                        batch_size=batch_size,
                        epochs=epochs,
                        verbose=1,
                        validation_data=(x_test, y_test))
score = model_ffnn.evaluate(x_test, y_test, verbose=0)
print('Test loss:', score[0])
print('Test accuracy:', score[1])

```

```

plt.figure(figsize=(5, 4))
plt.plot(history.epoch, history.history['acc'], lw = 1, label='Training')
plt.plot(history.epoch, history.history['val_acc'], lw = 1, label='Testing')
plt.legend()
plt.title('Accuracy of softmax regression', fontsize=16)
plt.xlabel('Epoch', fontsize=14)
plt.ylabel('Accuracy', fontsize=14)
plt.tick_params(labelright=True)
plt.grid('True')
plt.tight_layout()

```

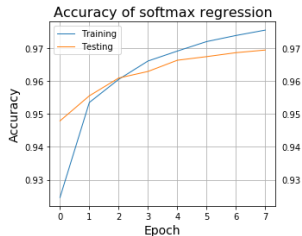


Figure 2: Minimization of the negative loglikelihood using a gradient descent algorithm (here AdaGrad). The gradient is computed using batches of 32 observations and the whole data set is used 8 times.