

# Deep learning for image analysis

Linus Falk

March 28, 2023

## Contents

<b>1</b>	<b>A linear classifier</b>	<b>2</b>
<b>2</b>	<b>Feed forward neural networks: Backpropagation</b>	<b>5</b>

# 1 A linear classifier

This lecture will describe the linear classifier and ways to train it.

## Recap deep learning: "end to end"

Working with deep learning can reduce the number of steps that are involved in image analysis comparing it to traditional approaches that are more human driven (engineering and prior knowledge based). The advantages with the "end to end" approach in deep learning is that we remove the prior knowledge needed for the specific case, that also removes any bias when starting to solve the problem also, since there is no pre-formed idea of how the result should look like. This is in the case were the data set can be consider non biased and non skewed.

The disadvantage of this is that it is typically very data hungry and needs large data sets in order to train well. The deep learning method is data driven and suffers if the data is not satisfying. Traditional image analysis methods c.

## Problem formulation

Given a image we have an array of  $X \times Y \times 3$  values representing the pixels in the image (3 colors). Images poses many challenges since the image is a 2D representation of a 3D object in a environment that can change.

- **viewpoint variations** change of camera angle
- **Illumination variations**, environment
- **Deformations**, the object may deform/
- **Occlusions**
- **Background clutter**
- **Intraclass variations**, many different looking cats...

So in contrast to sorting numbers, there is no simple solution to solve this problem with hard-coding.

```
def predict(image):  
    # no clue....  
    return class_label
```

## Data-driven approach

One way to solve it is to collect a lot of images and corresponding labels. Use a machine learning method to train a classifier and then evaluate the result on a test set of images.

The task is to design a classifier:  $f(x, \mathbf{W})$  that can tell us which class  $y_i \in \{1, 2, \dots, N\}$  an input image  $x_i$  belongs to.

1. Select a classifier type,
  - a linear affine classifier for example:  $y = \mathbf{W}x + b$
2. select a performance measure, loss function
  - Hinge loss
  - Negative Log-likelihood
3. Learning: find the parameters  $\mathbf{W} = \{W, b\}$  which maximizes the performance, minimizes the overall loss

## Multiclass SVM loss

If given an example image with label  $(x_i, y_i)$ , where  $x_i$  is the image and  $y_i$  is the label. The shorthand for the score vector:  $s = f(x_i, W)$ . Then the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad (1)$$

The full training loss is the mean over all examples in the training data:

$$L = \frac{1}{N} \sum_{i=1}^N L_i \quad (2)$$

## Softmax classifier (Multinomial logistic regression)

The scores are unnormalized log probabilities of the classes:  $s = f(x_i, W)$ .

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad (3)$$

We want to maximize the log likelihood (or minimize the loss function, the negative log likelihood)

$$L_i = -\log P(Y = y_i | X = x_i) \quad (4)$$

and in summary we have:

$$L_i = -\log\left(\frac{e^{s_k}}{\sum_j e^{s_j}}\right) \quad (5)$$

## Learning

How do we minimize the loss over the training data then:

- $\arg \min$  loss(training data)
- Follow the slope, gradient descent

With the second alternative we want to follow the slope like a blind person finding its way down from the hills. In the 1D case this is done with the derivative, but in the multidimensional case it's done with gradients (partial derivatives). **Gradient descent** can be used to minimize the loss  $L$  by:

1. Initialize the weights  $\mathbf{W}_0$
2. Compute the gradient with respect to  $\mathbf{W}$ ,  $\nabla L / (\mathbf{W}_k; x) = (\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots)$
3. Take a small step in the negative direction of the gradient:  $\mathbf{W}_{k+1} = \mathbf{W}_k - \text{stepsize} \cdot \nabla L$
4. Iterate from (2) until convergence

## Linear regression and classification

We will distinguish between two types of problems: regression and classification. In regression, the output  $y$  is a continuous quantity, for example a temperature, currency or distance. In classification the output is a discrete class label, for example: true/false and cat/dog/horse etc. But what is a good model? A good model is the one which **minimizes** our **Loss function**, so to answer the question we must know what a good Loss function is to begin with.

## The statistical approach

One common approach is the statistical **maximum likelihood** where we make the assumption that each data points can be described by a linear model + some noise with a normal distribution. The probability density function or PDF for the scalar Normal/Gaussian distribution:

$$\mathcal{N}(z; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z-\mu)^2}{2\sigma^2}} \quad (6)$$

Where  $\mu$  is the mean or expected value of the distribution,  $\sigma$  is the standard deviation and  $\sigma^2$  is the variance.  $z \sim \mathcal{N}(z; \mu, \sigma^2)$  means that  $z$  is a Normal/Gaussian random variable with the mean  $\mu$  and variance  $\sigma^2$ ,  $\sim$  : "distributed according to"

## Maximum likelihood

A linear model with Gaussian noise can be modeled with:

$$y_i = wx_i + b + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n \quad (7)$$

We can also express this as a probability:

$$p(y_i | x_i, w, b) = \mathcal{N}(y_i; wx_i + b, \sigma^2) \quad (8)$$

By picking the weights  $w$  and bias  $b$  that makes the data as likely as possible.

$$\hat{w}, \hat{b} = \arg \max p(y_1, \dots, y_n | x_1, \dots, x_n, w, b) \quad (9)$$

We assume here that all  $\epsilon_i$  are independent.  $y$  and  $z$  are also independent so:  $p(y|z) \Rightarrow p(y)p(z)$

$$p(y_1, \dots, y_n | x_1, \dots, x_n, w, b) = \prod_{i=1}^n (y_i - (wx_i + b))^2 \quad (10)$$

And the loss function is given by

$$L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|^2 \quad (11)$$

## Data driven approach to image classification

The task at hand is to design a classifier  $f(x, \mathbf{w})$  that can tell us which class an image  $x_i$  belongs to. We formulate an approach:

1. Select a classifier type:
  - Start off with a linear (affine) classifier  $y = Wx + b$
2. We then select a performance measure
  - SVM/hinge loss or Softmax + NLL loss
3. We then need for our data set the parameters  $W$  which maximizes the performance, which means minimize the overall loss. We do this by
  - Solve it with gradient descent (the learning part)

Another task is to design a regression model  $f(x, W)$  that tells us the value of  $y_i \in \mathbb{R}$  given a sample  $x_i$

1. Select a regression model
  - We start with a linear (affine) model  $y = Wx + b$

2. Then select a performance measure
  - Sum of squared errors for example
3. Then we want for this data set find the parameters  $W$  which maximizes the performance aka minimize the loss
  - Solve the gradient descent

## The limitations of linear classifiers

Some problem are difficult to solve with linear classifiers, since they are linear. A good way to see if it's possible to use a linear classifier is to have a look in the geometric space between the input variable and the output. Is it possible to fit a straight line through the data?

## Neural networks - stacked non-linear classifiers

By stacking multiple linear classifiers and adding non-linear activation functions between them, we can do better. Now it's possible to fit non-linear data to the model. A stacked linear classifier can be expressed as:

$$f(x) = f_2(f_1(x)) = W_2 W_1 c = W^* x \quad (12)$$

This is however still linear so we need to have an activation function  $h(x)$ . We then have a generalized linear classifier  $f(x) = h(Wx)$ . The Logistic regression is of this type with a sigmoid activation function. A stacked non-linear classifier:

$$f(x) = f_2(h(f_1(x))) = W_2 h(W_1 x) \quad (13)$$

This gives us a lot more power and versatility model, known as a feed forward artificial neural network (ANN). Here are some other activation functions:

- $\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$
- $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\text{sigmoid}(2x) - 1$
- $\text{ReLU}(x) = \max(0, x)$

## 2 Feed forward neural networks: Backpropagation

This lecture will cover how backpropagation works, how to compute the derivatives and some implementation.

### Neural network - construction

An artificial neural network is a sequential construction of several generalized linear regression models. It consists of **Inputs**  $(1, x_1, x_2, \dots, x_p)$ , **Hidden units**  $(1, q_1, q_2, \dots, q_U)$  and an output  $\hat{y}$

$$\begin{aligned} q_1 &= h(b_1^{(1)} + \sum_{j=1}^p W_{1j}^{(1)} x_j) \\ q_2 &= h(b_2^{(1)} + \sum_{j=1}^p W_{2j}^{(1)} x_j) \\ &\vdots \\ q_U &= h(b_U^{(1)} + \sum_{j=1}^p W_{Uj}^{(1)} x_j) \end{aligned} \quad (14)$$

$$\hat{y} = b^{(2)} + \sum_{i=1}^U W_i^{(2)} q_i \quad (15)$$

In vector representation:

$$\begin{aligned}
W^1 &= \begin{bmatrix} W_{11}^{(1)} & \dots & W_{1p}^{(1)} \\ \vdots & & \vdots \\ W_{U1}^{(1)} & \dots & W_{Up}^{(1)} \end{bmatrix}, b^{(1)} = \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_U^{(1)} \end{bmatrix}, q = \begin{bmatrix} q_1 \\ \vdots \\ q_U \end{bmatrix} \\
b^{(2)} &= [b^2], W^{(2)} = [W_1^{(2)} \dots W_U^{(2)}] \\
\hat{y} &= W^{(2)}q + b^2,
\end{aligned} \tag{16}$$

## Deep neural network

A deep neural network with  $L$  can be written like this:

$$\begin{aligned}
q^{(0)} &= x \\
q^l
\end{aligned} \tag{17}$$