Faculty of Mathematics, Physics and Informatics Comenius University Bratislava



Neural Networks

Lecture 3

Supervised single-layer models



Linear NN models

Input vector:

Output vector:

Weight matrix:

$$\boldsymbol{x} = [x_{1}, x_{2, \dots}, x_{n}]^{T}$$
$$\boldsymbol{y} = [y_{1}, y_{2, \dots}, y_{m}]^{T}$$
$$\mathbf{W} \sim \text{type } [m \times n]$$



Linear transformation ϕ : $\Re^n \rightarrow \Re^m$, y = Wx

ignores saturation property of neurons
allows to find analytic solutions using linear algebra.

(Kohonen, 1970; Anderson, 1972; Cooper, 1973)

- Adding layers in a linear NN does not appear reasonable (since no complexity is added).
- But: It allows nonlinear learning dynamics in linear deep nets (Saxe, 2015).

Closed-form solution

Let's consider the data set: $A_{train} = \{(x^{(p)}, y^{(p)}), p = 1, 2, ..., N\}.$ We look for a matrix W that satisfies $y^{(p)} = Wx^{(p)}$, for each p. In matrix notation: Y = WX $[y^{(1)} y^{(2)} ... y^{(N)}] = W \times [x^{(1)} x^{(2)} ... x^{(N)}]$ $(m \times N)$ $(m \times n)$ $(n \times N)$

If *X* was regular (i.e., square matrix with N = n, with linearly independent rows) then X^{-1} would exist and $W = YX^{-1}$.

However, in general we cannot assume N = n, nor linear independence of input vectors ($\Rightarrow X^{-1}$ does not exist). Only generalized solutions exists:

 $W = YX^+$

X⁺ is called (Moore-Penrose) pseudoinverse matrix of *X*. Theorem: $\forall X, \exists X^+$ with properties: 1) $XX^+X = X$, 2) $X^+XX^+ = X^+$, 3) symmetric XX^+ and X^+X .

a) $X^{+} = X^{T} (XX^{T})^{-1}$, if $n \le N$ and rank(X) = n.

b) $X^+ = (X^T X)^{-1} X^T$, if n > N and rank(X) = N.

Auto-association in a linear network

Let us now look at autoassociation: $y^{(p)} \equiv x^{(p)}$, $\dim(y) = \dim(x) = n$

This can be useful for a memory, if n > N (i.e. few examples of large dimension).

Consider the task: to remember N prototypes $[x^{(1)} x^{(2)} \dots x^{(N)}]$.

Goal: train a linear model on prototypes and then submit a corrupted version of a prototype. Model should be able to reconstruct it.

Since $Y \equiv X$, then $W = XX^+$. How to interpret *W*?

In a special, restrictive case (N = n, linearly independent inputs), we would have a trivial solution W = I (identity).

How about a general (non-trivial) case?

Basics of linear vector spaces

Let's have a linear space \Re^n .

Linear manifold $\mathscr{L} = \{ \mathbf{x} \in \Re^n \mid \mathbf{x} = a_1 \mathbf{x}^{(1)} + a_2 \mathbf{x}^{(2)} + \dots + a_N \mathbf{x}^{(N)}, a_p \neq 0 \}$ $\mathscr{L} \subset \Re^n$

Orthogonal complement $\mathscr{L}^{\perp} = \{x \in \mathfrak{R}^n \mid x \perp \mathscr{L}\}$

Hence, $\mathscr{L} \oplus \mathscr{L}^{\perp} = \Re^n$

Each vector $x \in \Re^n$ can be uniquely decomposed: $x = x_P + x_C$ where $x_P \in \mathscr{L}$ and $x_C \in \mathscr{L}^{\perp}$.



What does an autoassociative NN do?

Training set $A_{train} = \{x^{(p)}, p = 1, 2, ..., N\}$ forms the linear manifold \mathscr{L} .

NN considers every departure x from \mathscr{L} as added noise that needs to be filtered out by projecting x to \mathscr{L} :

We need to show that output $Wx = XX^+x = x_p$ (filtered version of x), i.e. that operator $W = XX^+$ makes an orthogonal projection to \mathscr{L} .

Alternatively, the NN model with operator $W = I - XX^+$ is called novelty detector, where $Wx = x_C \in \mathscr{L}^{\perp}$.

Now assume: you learned N patterns, and you want to add (N+1)-st pattern.

How to change W efficiently?



Gram-Schmidt orthogonalization process

Let's have a base $u^{(1)}, u^{(2)}, ..., u^{(k)} \in \mathscr{L}$, for which we want to create an orthogonal base $v^{(1)}, v^{(2)}, ..., v^{(k)} \in \mathscr{L}$.

Procedure:

Let $v^{(1)} \equiv u^{(1)}$. In space with base formed by $v^{(1)}$, $u^{(2)}$ let's find vector $v^{(2)}$ such that $v^{(2)} \perp v^{(1)}$. Hence, $v^{(2)} = a_1 v^{(1)} + a_2 u^{(2)}$.

$$v^{(2)} = u^{(2)} - \frac{v^{(1)T} u^{(2)}}{|v^{(1)}|^2} v^{(1)}$$

Recursive formula: we have $v^{(1)}, v^{(2)}, ..., v^{(k-1)}$ and compute $v^{(k)}$ such that $v^{(k)} \perp v^{(i)}, i = 1, 2, ..., k-1$

$$\mathbf{v}^{(k)} = \mathbf{u}^{(k)} - \sum_{i=1}^{k-1} \frac{\mathbf{v}^{(i)T} \mathbf{u}^{(k)}}{|\mathbf{v}^{(i)}|^2} \mathbf{v}^{(i)}$$



How to use this recursion for a GI model?

General Inverse model

If we have patterns $x^{(1)}, x^{(2)}, ..., x^{(N)}$, we can use Gram-Schmidt process to calculate orthogonal base $z^{(1)}, z^{(2)}, ..., z^{(N)}$

Memory model *W* is computed recursively:

1. Initialize $W^{(0)} = 0$, i = 1

2. Repeat
$$z^{(i)} = x^{(i)} - W^{(i-1)} x^{(i)}$$

 $W^{(i)} = W^{(i-1)} + \frac{z^{(i)} z^{(i)T}}{|z^{(i)}|^2}$
 $i = i+1$
 $Hint: \sum_{i=1}^{k-1} \frac{v^{(i)T} u^{(k)}}{|v^{(i)}|^2} v^{(i)} = \sum_{i=1}^{k-1} \frac{v^{(i)} v^{(i)T}}{|v^{(i)}|^2} u^{(k)}$

Example: 8 faces from CMU image data repository

tace 1



tace 4



tace 7



tace 2



tace 5



tace 8



tace 3



tace 6



tace 1 - corrupt



Recall by GI and novelty detection

GI



Novelly detection







New imput



Supervised single-layer perceptrons

- i = 1, 2, ..., m neurons, that can work independently
- then each neuron independently splits input space into two subspaces



• For sigmoidal neurons we get learning rule: $(net_i \equiv o_i)$

 $w_{ij}(t+1) = w_{ij}(t) + \alpha . (d_i - y_i) . f'(o_i) . x$

- for linear neurons we get the least-means-square (LMS) learning rule $w_{ij}(t+1) = w_{ij}(t) + \alpha . (d_i - y_i) . x_j$
- closed-form solution may exist via pseudoinverse: for a linear model $W = D.X^+$, or $W = f^{-1}(D).X^+$, if a nonlinearity is used $[f^{-1}(d_i) = o_i]$

Softmax regression

- In classification tasks, where #output_neurons = #classes, it is convenient to interpret outputs as (conditional) probabilities
- hard assignment vs soft assignment (of classes)
- Softmax activ. function (Luce, 1959) enables that (by normalization): all nonnegative values sum up to one

$$y_{i} = \frac{\exp(o_{i})}{\sum_{k} \exp(o_{k})} \qquad \text{argmax}_{i} y_{i} = \operatorname{argmax}_{i} o_{i} \qquad o_{i} = \sum_{j} w_{ij} x_{j} + b_{i}$$

logit

- Important: the model remains differentiable
- Despite nonlinearity, the outputs of softmax regression are still determined by an affine transformation of input features (→ logits)
- Hence, softmax regression is a linear model.

(Zhang et al, 2020)

Log-likelihood estimation

- Assume a training set $\{X, D\}$, consisting of *N* pairs of input vectors $x^{(p)}$ and one-hot label vectors $d^{(p)}$. Let NN model predictions be $y^{(p)}$.
- We compare the predictions with ground truth by checking how probable the actual classes are according to our model, given the features:

$$P(\boldsymbol{D}|\boldsymbol{X}) = \prod_{p=1}^{N} P(\boldsymbol{d}^{(p)}|\boldsymbol{x}^{(p)})$$

• According to maximum likelihood estimation, we maximize P(D|X), which is equivalent to minimizing the negative log-likelihood, i.e. $-\log P(D|X)$

$$-\log P(\mathbf{D}|\mathbf{X}) = -\sum_{p=1}^{N} \log P(\mathbf{d}^{(p)}|\mathbf{x}^{(p)}) = \sum_{p=1}^{N} loss(\mathbf{d}^{(p)}, \mathbf{y}^{(p)})$$

• For each pattern *p* (with *m* classes assumed)

$$loss_{CE}^{(p)} = CE^{(p)} = -\sum_{i=1}^{m} d_i^{(p)} \log y_i^{(p)}$$

Cross-entropy loss calculation

- $CE = -\log P(\boldsymbol{d}|\boldsymbol{x}) = -\sum_i d_i \log y_i$,
- *i*, *k* =1,2,...,*m*

 $y_i = \frac{\exp(o_i)}{\sum_k \exp(o_k)}$

 $CE = \sum_{i} d_{i} \log \sum_{k} \exp(o_{k}) - \sum_{i} d_{i} o_{i} = \log \sum_{k} \exp(o_{k}) - \sum_{i} d_{i} o_{i}$

• The derivative:

$$\frac{d e_{CE}}{d o_i} = \frac{\exp(o_i)}{\sum_k \exp(o_k)} - d_i = P(y=i|\mathbf{x}) - d_i$$

- has the same effect as in case of the squared error
- allows easy gradient computation
- provides a useful link to information theory

Basics of information theory

- assume a random discrete event that can take values *j*
- Entropy: $H(p) = -\sum_{j} p_{j} \log p_{j}$, reaches max. for uniform distribution
- In order to encode data drawn randomly from the distribution *p*, we need at least *H*(*p*) "nats" to encode it (analogy to a "bit").
- $-\log q_j$ quantifies surprisal observing an event *j* having assigned it a (subjective) probability q_j
- Cross entropy *H*(*p*, *q*) is the **expected surprisal** of an observer with subjective probabilities *q* upon seeing data that was actually generated according to probabilities *p*.
- Kullback-Leibler divergence measures the distance b/w two distributions: $D_{KL}(p,q) = H(p,q) H(p) = -\sum_j p_j \log (p_j / q_j)$
- hence, minimizing $D_{KL}(p,q)$ corresponds to minimizing H(p,q)

Model learning: Sequential or batch mode

Sequential mode

- on line (example-by-example), stochastic
- able to track small changes in training data
- easier to implement, requires less local storage
- difficult to establish theoretical conditions for convergence

Batch mode

- adaptation performed at the end of each epoch, deterministic
- provides an accurate estimate of gradient vector
- parallelization possible

$$E_{av} = 1/(2N) \sum_{p=0}^{N} (d^{(p)} - y^{(p)})^2 \qquad \Delta w_i \propto -\partial E_{av}(t)/\partial w_i$$

Mini-batches – best of two "worlds" (typical minibatch size 50–256)

Online versus batch learning in param. space

Batch learning - does steepest descent on the error surface

Online learning - zig-zags around the direction of steepest descent



Example: training a single-layer perceptron

- MNIST data set, 10 classes of hand-written digits, images 28x28
- Accuracy 90+ %, check also http://yann.lecun.com/exdb/mnist/



https://www.oreilly.com/content/not-another-mnist-tutorial-with-tensorflow/

Summary

- Linear models studied in 1970s but concepts still useful
- Single-layer models are linear in parameters
- Hence, analytic solutions possible (via pseudoinverse)
- Tasks: (self-supervised) autoassociation, (supervised) regression, classification
- General Inverse model (projection to linear manifold)
- Data classification soft assignment (softmax)
- Cross-entropy error used for classification
- Link to information theory (via entropy)
- Batch vs online learning