

Neuronale Netze

Zusammenfassung

Marius Bauer

31. Juli 2017

Inhaltsverzeichnis

1	Neuronale Netze - kurzer Überblick	4
2	Pattern Recognition	7
3	Recurrent Neural Networks	7
3.1	Sequence Learning	7
3.2	Elman vs. Jordan Networks - Simple RNN	7
3.3	Aufbau	8
3.4	Training	9
3.5	Vanishing / Exploding Gradient	9
4	Speech	10
4.1	Speech Recognition	10
4.2	Acoustic Model	11
4.2.1	Hidden Markov Model	11
4.2.2	Time Delay Neural Network (TDNN)	12
4.3	Word Model	13
4.3.1	Time Alignment	14
4.3.2	Multi-State-TDNN	14
4.3.3	NN-HMM Hybride	15
4.3.4	Viterbi-Algorithmus	15
4.3.5	Forward-Algorithmus	16
4.3.6	Backward-Algorithmus	16
5	Learning Vector Quantization	17
5.1	Vector Quantization	17
5.1.1	Applications	17
5.1.2	Training	17
5.2	Learning Vector Quantization	17
5.2.1	LVQ1	18
5.2.2	LVQ2	18
5.2.3	LVQ2.1	19
5.2.4	LVQ3	19
5.2.5	OLVQ	19
6	Self-Organizing-Maps	19
6.1	Principles Of Self-Organized-Learning	19
6.1.1	Principle 1: Self-Amplification	19
6.1.2	Principle 2: Competition	19
6.1.3	Principle 3: Cooperation	19

6.1.4	Principle 4: Structural Information	19
6.1.5	Hebbian Learning	19
6.2	Self-Organizing-Maps	19
7	Reinfocement Learning	19
8	Deep Learning in Computer Vision	20
9	Neural Network Applications in Machine Translation	20

1 Neuronale Netze - kurzer Überblick

Wieso werden neuronale Netze verwendet?

- Massive parallelism.
- Massive constraint satisfaction for ill-defined input.
- Simple computing units.
- Many processing units, many interconnections.
- Uniformity (-> sensor fusion)
- Non-linear classifiers/ mapping (-> good performance)
- Learning/ adapting
- Brain like ??

Wofür werden neuronale Netze verwendet?

- Classification
- Prediction
- Function Approximation
- Continuous Mapping
- Pattern Completion
- Coding

Welche Kriterien müssen beim Entwurf von neuronalen Netzen beachtet werden?

- Recognition Error Rate
- Training Time
- Recognition Time
- Memory Requirements
- Training Complexity
- Ease of Implementation
- Ease of Adaptation

Welche Parameter werden typischerweise vom Entwerfer bestimmt?

- Net Topology
- Node Characteristics
- Learning Rule
- Objective Function
- (Initial) Weights
- Learning Parameters

Wo werden neuronale Netze verwendet?

- Space Robot*
- Autonomous Navigation*
- Speech Recognition and Understanding*
- Natural Language
- Processing*
- Music*
- Gesture Recognition
- Lip Reading
- Face Recognition
- Household Robots
- Signal Processing
- Banking, Bond Rating
- Sonar

Advanced Neural Models

- Time-Delay Neural Networks (Waibel)
- Recurrent Nets (Elman, Jordan, Robinson,...)
- Higher Order Nets

- Modular System Construction
- Adaptive Architectures
- Hybrid Neural/Non-Neural Architectures
- Deep Nets

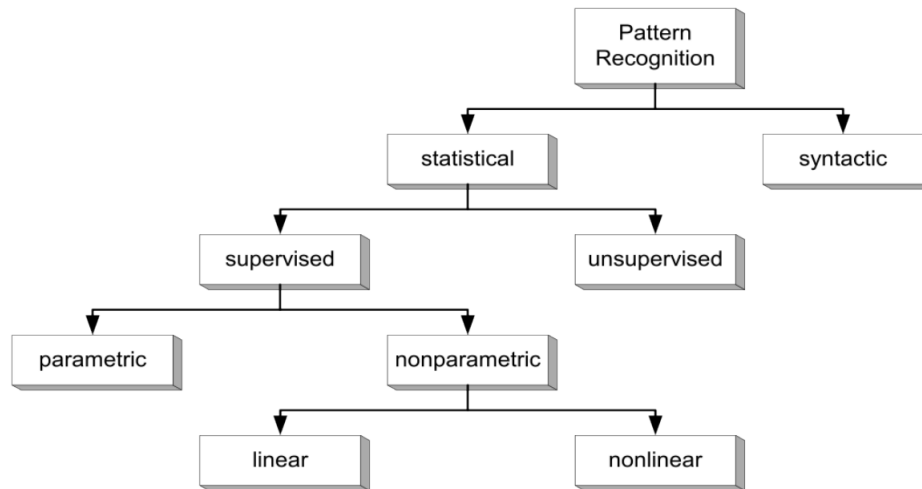
Welche Probleme können beim Entwurf von neuronalen Netzen auftreten?

- Local Minima
- Speed of Learning
- Architecture must be selected
- Choice of Feature Representation
- Scaling
- Systems, Modularity
- Treatment of Temporal Features and Sequences

Eigenschaften von neuronalen Netzen:

- non-linear classifier
- approximate posterior probability
- non-parametric training

2 Pattern Recognition



3 Recurrent Neural Networks

Recurrent Neural Networks unterscheiden sich von anderen neuronalen Netzen in dem Punkt, dass sie einen (oder mehrere) Zustände speichern und verwenden können. Dabei ist ein Zustand ein berechneter Wert des vorherigen Zeitschritts. Diese Zustände können wiederum als Input genutzt werden. Somit ist eine zeitliche Zustandsspeicherung möglich.

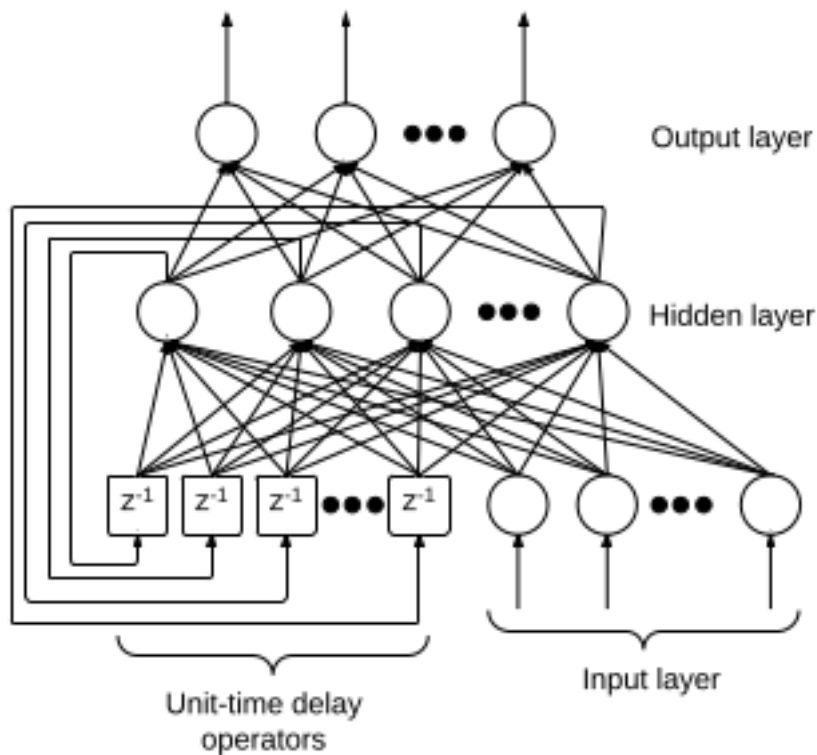
3.1 Sequence Learning

Sequenzen sind überall in unserem alltäglichem Leben vertreten: z.B. bei der Sequenzierung von Töne in Sprache. Bei der Spracherkennung wird der Fokus auf Wortsequenzen gelegt. Bei sequence learning gibt es vier generelle Probleme:

- sequence prediction: Was wird das nächste Element in der Sequenz sein?
- sequence generation: generieren einer Sequenz (eg. Wortsequenz)
- sequence recognition: Ist die Sequenz legitim?
- sequence decision making: ??

3.2 Elman vs. Jordan Networks - Simple RNN

Der Unterschied zwischen *Elman* und *Jordan* Netzen besteht darin, dass sie unterschiedliche Outputs zwischenspeichern und wiederverwenden. Bei Elman Netzen wird der Output des Hidden Layers als

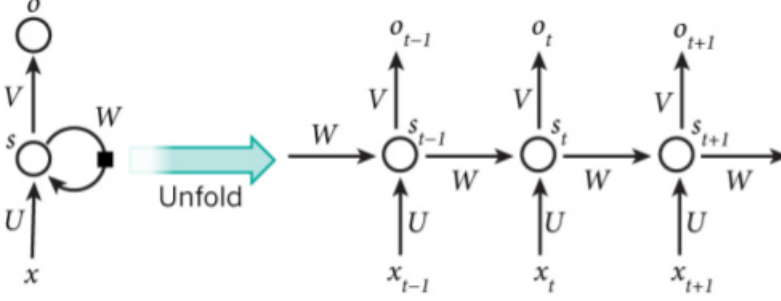


Input für den nächsten Zeitschritt verwendet. Jordan Netze auf der anderen Seite nutzen den Output des Netzes als Input für den nächsten Zeitschritt.

3.3 Aufbau

- Input Units $X = x_1, x_2, \dots, x_n$
- Output Units $Y = y_1, y_2, \dots, y_n$
- Hidden Units $H = h_1, h_2, \dots, h_n$
- Connections between Units

Im Prinzip sehen RNN genauso aus wie NN. Jedoch sind die Verbindungen anders. Diese können nämlich zurückführen und als Input dienen. Somit lassen sich vergangene States wiederverwenden. Alternative könnte man *recurrent hidden units* in den Aufbau mit aufnehmen. Diese speichern vergangene Zustände und leiten diese als Input in den darauffolgenden Zeitschritt in die entsprechenden Units.



3.4 Training

RNN werden mittels *backpropagation through time* trainiert. Dies ist eine Erweiterung des ursprünglichen Backpropagation, bei der die wiederkehrenden Zustände berücksichtigt werden. Dazu wird ein RNN *ausgerollt* um den zeitlichen Zusammenhang darzustellen: Zu jedem Zeitpunkt nutzen wir *forward pass* um die Aktivität in jedem Zeitschritt zu berechnen. Mittels *backward pass* wird die Fehlerableitung zu jedem Zeitschritt berechnet. In dem entrollten Netzwerk kommen die selben Gewichte öfters vor (shared weights), zu sehen in Bild XX. Jede Instanz eines Gewichtes muss den selben Gradienten erhalten (Einschränkung):

1. $w_j^{t_1} = w_j^{t_2} \Rightarrow \Delta w_j^{t_1} = \Delta w_j^{t_2}$
2. Berechne $\frac{\partial E}{\partial w_j^{t_1}}$ und $\frac{\partial E}{\partial w_j^{t_2}}$
3. $\Delta w_j^{t_1} = \Delta w_j^{t_2} = -\eta \left(\frac{\partial E}{\partial w_j^{t_1}} + \frac{\partial E}{\partial w_j^{t_2}} \right)$

Hierbei kann zum einen die Summe betrachtet werden, aber auch der Durchschnitt!

3.5 Vanishing / Exploding Gradient

4 Speech

Die Sprachproduktion beim menschlichen Körper besteht aus drei Teilen: 1. dem supraglottalen Sprachtrakt, 2. dem Kehlkopf und 3. dem subglottalen System (teil davon ist die Lunge). Die Lunge und der Kehlkopf werden dabei über motorische Bewegung angeregt. Genauer gesagt, die Lunge stellt die Luft bereit, welche vom Kehlkopf benötigt wird um angeregt zu werden. Die Lunge und der Kehlkopf zusammen sorgen somit für die *Anregung*. Je nach "Form" des Vokaltrakts (Rachen-, Mund- und Nasenraum) kommt es durch die Anregung zu einer anderen *Artikulation*. Dabei werden Töne erzeugt und "ausgegeben". Aneinander gereihte Töne ergeben dadurch Sprache. Die verschiedenen Töne produzieren verschiedene Spektre. Für bestimmte Laute bzw. Vokale lassen sich typische Muster im Spektralbereich finden und somit die Vokale unterscheiden.

4.1 Speech Recognition

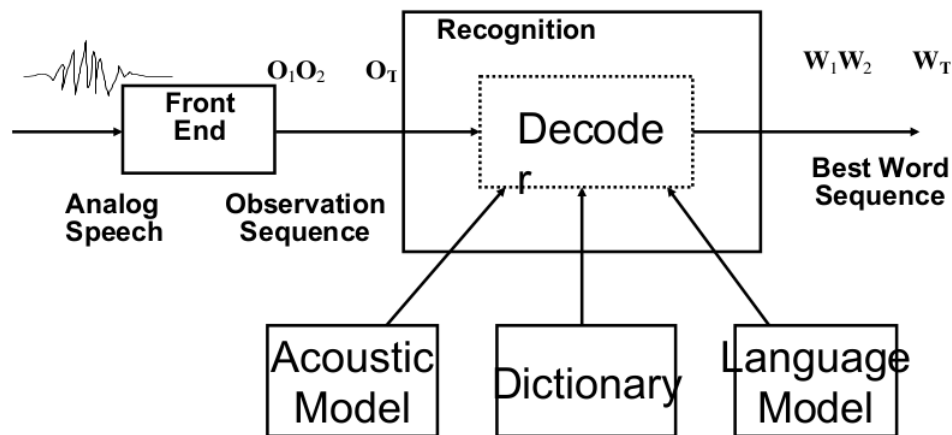


Abbildung 1: Komponenten eines Sprach Erkenners

Auf dem obigen Bild sind die einzelnen Komponenten eines Sprach Erkenners zu sehen. Links kommt das *analoge Signal* in das Front End rein und wird digitalisiert. Dazu wird es mittels Fouriertransformation in den Spektralbereich umgewandelt. Aus dem Frontend kommt die *observation sequence* heraus und wird an den *Decoder* weitergegeben. Im Decoder passiert die eigentliche Sprach Erkennung. Dazu wird ein *acoustic model*, ein *dictionary* und ein *language model* verwendet. Die einzelnen Komponenten werden wir uns später noch genauer anschauen. Der Decoder gibt am Ende die *best word sequence* aus. Das Ziel des Sprach Erkenners ist es aus einer gegebenen akkustischen Daten $A = a_1, a_2, \dots, a_k$ die Wortsequenz $W = w_1, w_2, \dots, w_n$ zu finden, welche die Wahrscheinlichkeit

$P(W|A)$ maximiert. Dazu brauchen wir die *Bayes Regel*:

$$P(W|A) = \frac{P(A|W) \cdot P(W)}{P(A)}$$

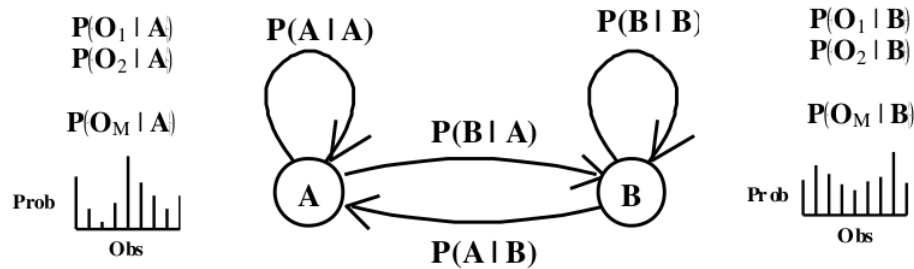
$P(W|A)$ ist das *acoustic model*, $P(W)$ ist das *language model* und $P(A)$ ist konstant für einen ganzen Satz.

4.2 Acoustic Model

4.2.1 Hidden Markov Model

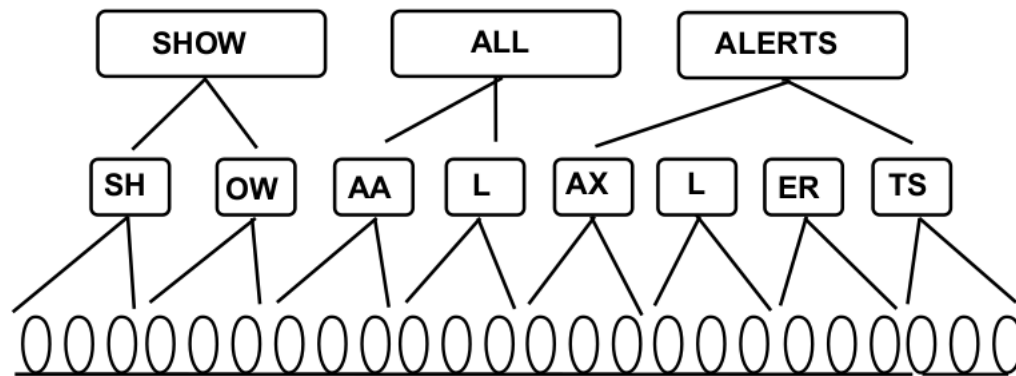
Elements:

- states: $S = \langle S_0, S_1, \dots, S_N \rangle$
- transition probabilities: $P(q_t = S_i | q_{t-1} = S_j) = a_{ji}$
- emission probabilities: $P(y_t = O_k | q_t = S_j) = b_j(k)$



In der Spracherkennung spiegeln die States die phonetischen Zustände wieder (genauer gesagt besteht jedes Phonem aus drei Zuständen: beginning, middle und end). Der Decoder sucht hierbei den besten Pfad zwischen den Modellen und der Sprache. Um die Emissionswahrscheinlichkeit zu approximieren können neben HMMs alternative Methoden verwendet werden:

- Mixture of Gaussians Networks
- Neural Networks
- Hierarchies of Neural Networks

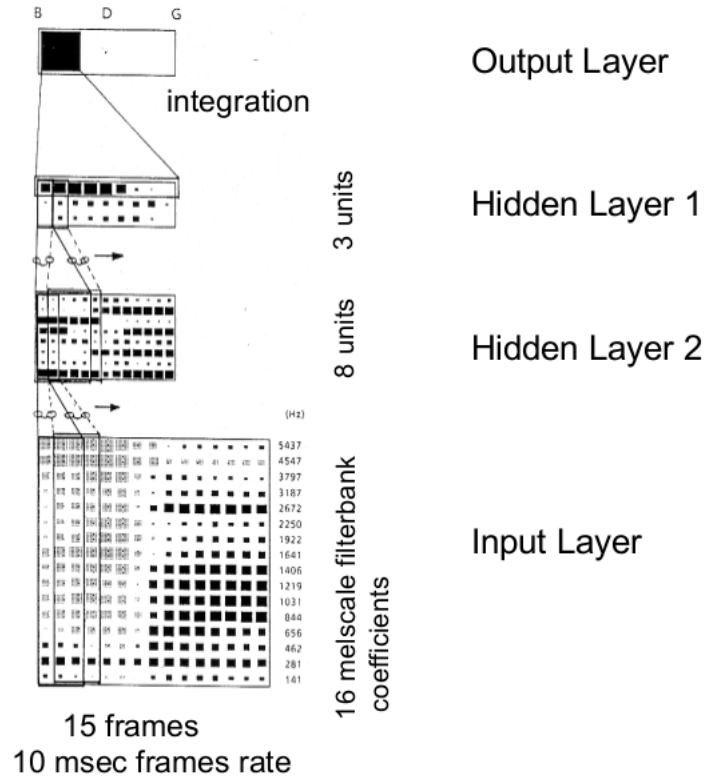


Die Wörter

werden zerlegt in Phoneme und die Phoneme in Zustände.

4.2.2 Time Delay Neural Network (TDNN)

Ein *time delay neural network* ist eine spezielle Art von neuronalem Netz, welches *time invariant* ist. Es besteht aus einem Input Layer, mehreren Hidden Layers und einem Output Layer. Als Eingabe dient eine Matrix welche *mel-scaled* ist. Die x-Achse spiegelt die Zeit wider, die y-Achse die Frequenz. Der erste Hidden Layer erhält als Eingabe ein 16 x 3 Ausschnitt der Eingabe, jeweils verschoben in der x-Richtung. Somit überlappen sich die Eingaben für den ersten Hidden Layer. Im Hidden Layer wird ebenfalls ein Fenster drüber gelegt und einem weiteren Hidden Layer als Input gegeben. Im Output Layer wird auf einzelne Buchstaben gemappt. Das bedeutet ein TDNN kann ein Muster in einem Signal finden egal wo es sich befindet. Alles was wir wissen ist, dass es darin vorkommt. Im letzten Hidden Layer werden horizontal die Aktivierungen der Neuronen angeschaut und daraus bestimmt, welches Phonem vorgekommen ist. Der letzte Hidden Layer wird auch als *phoneme layer bezeichnet*. Jedoch bezieht sich das TDNN in unserem Fall auf einzelne Phoneme, jedoch werden Wörter aus mehreren Phonemen gebildet. Siehe dazu 4.3.2. TDNN sind *convolutional networks* mit einer festen Größe. Dies ist beim Übergang zwischen den Schichten zu sehen. Convolutional neuronale Netze werden in der Bildverarbeitung verwendet. Hierbei ist es so, dass z.B. Straßenschilder aufgenommen, aber nicht immer an der selben Stelle im Bild sind.



Eine Erweiterung der TDNN stellen die *frequency shifting TDNN (FSTDNN)* dar. Diese Netze sind shift invariant in zwei Dimensionen. Das bedeutet nicht nur Zeitinvariant sondern auch Frequenzinvariant. Normalerweise wird ein Netz auf eine Person trainiert und die Fehlerrate steigt, sobald eine andere Person das Netz benutzt. Die Frequenzinvarianz sorgt dafür, dass Männer, Frauen oder Kinder das selbe Netz benutzen können und die Fehlerrate verändert sich nur gering. Dies nennt sich auch *Multi-Speaker Phoneme Recognition*. Darüber hinaus gibt es noch das Problem des Echo / der Verhallung. Wenn kein Mikrofon direkt am Mund getragen wird, wird das Signal verrauscht / verschmiert. Die Verhallung ist immer unterschiedlich und man kann kein Netz für jeden Raum bauen. Bei der Verhallung wird das ursprüngliche, gesprochene Signal mit dem Echo gefaltet. Neuronale Netze können diese Verhallung lernen und rausfiltern.

4.3 Word Model

Bisher haben wir nur die Erkennung von Phonemen behandelt. Diese müssen nun zu ganzen Wörtern zusammengebaut werden. Eine Idee wäre ein neuronales Netz für ein bestimmtes Vokabular zu trainieren. Jedoch ist das Unsinn, da ein neues Wort bedeutet, dass man Beispiele sammeln muss und das Netz neu trainieren muss.

Problems:

- Time Alignment: the same word might be spoken faster / slower
- Endpoint Detection: where does one word end?
- Large Dictionaries (Training Data, Time)

4.3.1 Time Alignment

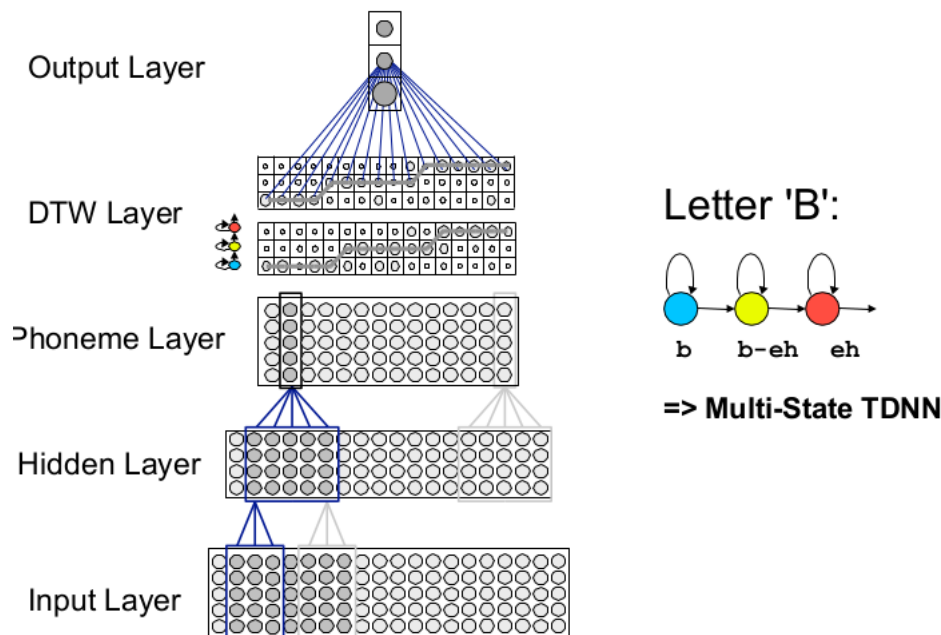
Hierbei geht es darum, dass verschiedene Sprecher unterschiedlich schnell/langsam sprechen. Irgendwie muss beim Lernen einen zeitlichen Zusammenhang herausgearbeitet werden.

Beim *Linear Sequence Alignment* wird das Referenzsignal und das unbekannte Signal linear aneinander angepasst (Stauen / Strecken). Problem hierbei ist, dass der Mensch in der Regel nicht linear schnell redet. Wir benötigen ein **non-linear alignment**.

Eine Lösung ist das *time warping* (auch dynamic programming genannt?!?). Hierbei wird eine Verbindung zwischen zwei Sequenzen x und y gesucht. Die Achsen repräsentieren die Sequenzen, welche Vektoren sind, deren Einträge gut oder schlecht aufeinander passen. Sind die Signale aufeinander abgestimmt, kann mittels des *Viterbi-Algorithmus* die passende Sequenz Q gefunden werden, welche $P(O, Q|\lambda)$ maximiert.

4.3.2 Multi-State-TDNN

Ein MS-TDNN ist eine Erweiterung des TDNN auf Wortebene. Um ein Wort zu erkennen muss eine bestimmte Sequenz von Phonemen auftreten. Im DTW-Layer(dynamic time warping) wird die Verbindung der Phoneme vorgenommen und im Output Layer lese ich das entsprechende Wort ab. Das bedeutet aber wiederum, dass ich für jedes einzelne Wort ein Output Neuron brauche. Das ist bei großen Wörterbüchern unpraktisch. Jedoch ist das nicht so schlimm, da wir auf dem Wortlevel Trainieren können.



Ich versuche die Distanz zwischen dem richtigen Wort und dem besten, falschem Wort zu maximieren. Das heißt wir optimieren die phonetische Sequenz und nicht die Wörter an sich. Im MS-TDNN befindet sich ein Decoder. Das Alinieren der Zustände und der Feature, welche das neuronale Netz über die Zeit gelernt hat, läuft wie ein HMM mit dynamik time warping oder mit Viterbi Decoding. Der Unterschied zu den Hybriden, welche danach dran kommen, ist, dass das es sich hierbei um hidden Features handelt.

4.3.3 NN-HMM Hybride

Die Idee ist, dass ein neuronales Netz gebaut wird, welches die Phoneme erkennt. Wenn der Output statisch interpretiert werden kann, wird ein HMM Decoder obendrüber gebaut, der für das Alignment und die Integration der Wörter zuständig ist. Die Wahrscheinlichkeit, welche wir als Output des NN bekommen, ist nicht die richtige. Um die richtige Wahrscheinlichkeit zu erhalten, müssen wir die Bayes Regel umformulieren. Wir wollen $P(A|W)$ haben aber $P(W|A)$.

Die (log)Wortwahrscheinlichkeit besteht aus der Summe der (log)Output Activations entlang des besten Pfades (bestimmt mit DTW oder Viterbi).

Kontextabhängige Phonemmodelle (a in Kontext von P etc.)

4.3.4 Viterbi-Algorithmus

Der Viterbi Algorithmus sucht das beste Alignment zwischen allen phonetischen Sequenzen, die vorkommen können und der gesprochenen Sprache.

4.3.5 Forward-Algorithmus

4.3.6 Backward-Algorithmus

5 Learning Vector Quantization

Der initiale Ansatz ist *vector quantization*, welcher als erstes besprochen wird. Die Weiterentwicklung des vector quantization ist *learning vector quantization*, welche danach behandelt wird.

5.1 Vector Quantization

Die Idee hinter vector quantization ist, dass der *data space* mit einer kleinen Anzahl an *prototype vectors* $U = u_1, u_2, \dots, u_k$ beschrieben wird. Jeder prototype vector repräsentiert eine Klasse (k Klassen). Dazu muss jeder Input Vektor $X = x_1, x_2, \dots, x_m$ einem der Vektoren in U zugewiesen werden. Die prototype vectors U sind in einem *codebook* zusammengefasst. Vector Quantization ist unsupervised.

5.1.1 Applications

- multimedia compression (storage and transmission)
- dimensionality reduction
- classification

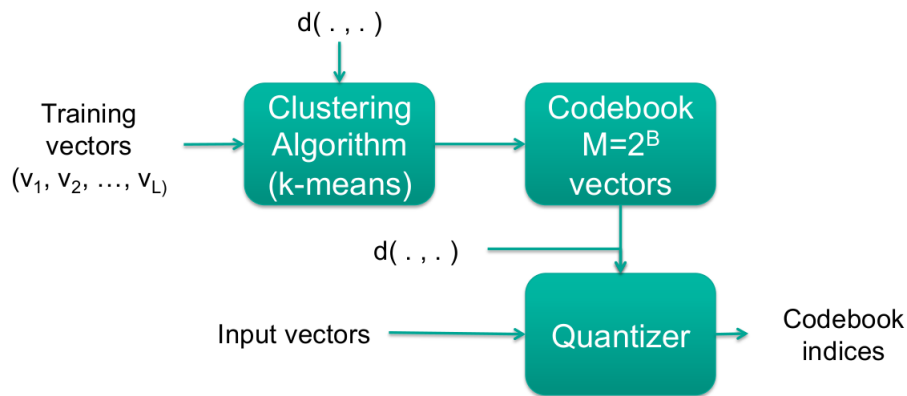
5.1.2 Training

Beim Trainieren werden die Input Vektoren mittels einem geeigneten Distanzmaß d und einem Clustering Algorithmus unterteilt. Daraus ergeben sich die Anzahl an Klassen und somit auch direkt die Anzahl an Codebuch Einträgen. Die Vektoren im Codebuch sind die *centroids* der einzelnen Regionen. Bei der eigentlichen Klassifikation (Encoding) wird für einen Input Vektor die Distanz zu den Codebuch Einträgen berechnet und bei der minimalen Entfernung der Index des prototype vectors gemerkt (bzw. übermittelt). Beim Decoding wird anhand des Index der jeweilige prototype vector ausgewählt und repräsentiert somit den Input Vektor. Der Fehler hierbei ist das Distanzmaß des Input Vektors und des prototype vectors.

5.2 Learning Vector Quantization

Da wir nun über Labels verfügen (also wissen in welche Klasse der Input Vektor gehört), lassen sich nun Wahrscheinlichkeitsverteilungen nutzen:

- $P(S_k)$: a prior Wahrscheinlichkeit der Klasse S_k
- $p(x|x \in S_k)$: bedingte Wahrscheinlichkeitsdichte
- discriminant functions: $\delta_k(x) = p(x|x \in S_k)P(S_k)$
- rate of misclassification minimized if: $\delta_c(x) = \max_k \{\delta_k(x)\}$



Prinzipiell weisen wir mehrere prototype vectors den einzelnen Klasse zu. Ein Input Vector erhält die selbe Klasse wie der prototype vector, an dem der Input Vektor am nächsten ist. Da wir nun aber mehrere prototype vectors je Klasse haben, müssen wir entscheiden, welcher der *winning codebook entry* (c ist der Index) ist. Die einzelnen Algorithmen unterscheiden sich in dieser Wahl:

5.2.1 LVQ1

$c = \operatorname{argmin}_i \{ \|x - m_i\| \}$: Index des prototype vectors

Learning rules:

- $m_c(t+1) = m_c(t) + \alpha(t)[x(t) - m_c(t)]$: x und m_c selbe Klasse
- $m_c(t+1) = m_c(t) - \alpha(t)[x(t) - m_c(t)]$: x und m_c unterschiedliche Klasse
- $m_c(t+1) = m_i(t)$: für $i \neq c$

5.2.2 LVQ2

Die Klassifizierung ist die selbe wie bei LVQ1. Das updaten ist anders:

- m_i und m_j sind die nächsten Nachbarn von x und werden simulaten geupdated.
- x muss in ein "window" um m_i und m_j fallen.
- d_i und d_j sind die Distanzen (z.B. euklidien) zwischen x und m_i und m_j
- $\min \left(\frac{d_i}{d_j}, \frac{d_j}{d_i} \right) > s$ where $s = \frac{1-w}{1+w}$ (recommended window: 0.2 to 0.3)

Ergänzungen: m_i ist der Gewinner (falsche klasse) m_j zweiter gewinner, richtige klasse -> dann updaten α die Lernrate wird kleiner und kleiner

5.2.3 LVQ2.1

5.2.4 LVQ3

Ergänzen: - die Frage ist welches k das korrekte ist. Zur Not kann man die Distortion betrachten und ein anderes k wählen +- - $\Delta_k = w_{\text{weit}}$ das x zu s_k gehört - stopping rule -> fehler kleine genug oder anzahl iteration erreicht

5.2.5 OLVQ

6 Self-Organizing-Maps

6.1 Principles Of Self-Organized-Learning

6.1.1 Principle 1: Self-Amplification

6.1.2 Principle 2: Competition

6.1.3 Principle 3: Cooperation

6.1.4 Principle 4: Structural Information

6.1.5 Hebbian Learning

6.2 Self-Organizing-Maps

SOM theory structure learning properties VQ kernel SOM applications

7 Reinforcement Learning

agent -> can act action -> influences state success -> scalar reward signal

in nutshell -> actions to maximize future reward

sequential decision making -> markov decision process policy $a = \pi(s)$ die policy ist eine wkeit-verteilung abhängig vom momentan zustand value function gegeben policy und state: wie gut ist action a in zustand s ? zwei zustände haben den gleichen reward, können sich aber später unterscheiden r = reward, γ = continuierungsfaktor (0.9) policy-based ein optimal policy: achieving maximum future reward value-based optimal value function: maximum value achievable under any policy

policy-based: actions sequence welche immer den besten reward gibt je action value-based: aktionssequenz welche am ENDE den maximalen reward liefert. (aktion innerhalb der sequenz muss nicht optimal sein) Bsp: roboter mit energieverbrauch und Höhe (laufen) Regression: lernen einer Funktion bellmann equation: q-learning: back: reward je action / for: chose action with highest reward F29: wie kommt die ableitung zustände? Q-Values: $Q^\pi(s, a)$

TD-Learning:

Beispiel für Q-Learning suchen

Fragen: - wie parametrisieren wir eine Policy? - Wie sieht eine Policy aus?

8 Deep Learning in Computer Vision

class: was sieht man im bild? local: was sieht man und wo? (bounding boxes) detection: local + segmentation: trivial

low/mid/high - level: wie viele pixel betrachtet werden

rezeptives Feld: beim pooling wird rausgezoomt, der betrachtete Bildbereich wird vergrößert.

- pooling fördert die betrachtung von lokalen zusammenhängen (die ohne nvl nicht zu erkennen sind) - überlappung fördert invarianz stationarity (translational invariance): V19F17: wenn das bild durchgeschoben wird, wird dennoch die selbe wkeit detektiert (aber in unterschiedlichen Neuronen)

CNN: convolutional layer, sub sampling layer, fully connected MLP multiple convolution

Image caption with attention

Calculate Parameters of NNs

9 Neural Network Applications in Machine Translation

NLP vs MT