

Neuronale Netze

Zusammenfassung

Marius Bauer

2. August 2017

Inhaltsverzeichnis

1	Neuronale Netze - kurzer Überblick	5
2	Pattern Recognition	8
2.1	Data Feature	9
2.1.1	Preprocessing:	9
2.1.2	Comparison	9
2.2	Parametric Methods	9
2.2.1	Bayes Decision Theory	9
2.2.2	Mixture Densities	10
2.3	Classifier Discrimination Function	10
2.4	Curse of Dimensionality	11
2.4.1	Principle Component Analysis (PCA)	11
2.5	Non-Parametric Methods	11
2.5.1	Parzen Window	11
2.5.2	k-nearest Neighbours	11
2.5.3	Clustering	12
2.6	Fisher-Linear Discriminant	12
2.7	Perceptron	12
3	Recurrent Neural Networks	13
3.1	Sequence Learning	13
3.2	Elman vs. Jordan Networks - Simple RNN	13
3.3	Aufbau	13
3.4	Training	14
3.5	Vanishing / Exploding Gradient	15
4	Backpropagation	16
4.1	Regularization	18
4.1.1	Weight Elimination	18
4.1.2	Optimal Brain Damage	18
4.1.3	Cascade Correlation	19
4.1.4	Meiosis Network	19
5	Speech	20
5.1	Speech Recognition	20
5.2	Acoustic Model	21
5.2.1	Hidden Markov Model	21
5.2.2	Time Delay Neural Network (TDNN)	22
5.3	Word Model	23
5.3.1	Time Alignment	24

5.3.2	Multi-State-TDNN	24
5.3.3	NN-HMM Hybride	25
5.3.4	Viterbi-Algorithmus	25
5.3.5	Forward-Algorithmus	26
5.3.6	Backward-Algorithmus	26
6	Learning Vector Quantization	27
6.1	Vector Quantization	27
6.1.1	Applications	27
6.1.2	Training	27
6.2	Learning Vector Quantization	27
6.2.1	LVQ1	28
6.2.2	LVQ2	28
6.2.3	LVQ2.1	29
6.2.4	LVQ3	29
6.2.5	OLVQ	29
7	Self-Organizing-Maps	30
7.1	Principles Of Self-Organized-Learning	30
7.1.1	Principle 1: Self-Amplification	30
7.1.2	Principle 2: Competition	30
7.1.3	Principle 3: Cooperation	30
7.1.4	Principle 4: Structural Information	30
7.1.5	Hebbian Learning	30
7.2	Self-Organizing-Maps	30
8	Reinforcement Learning	31
8.1	Bellman Equation	33
8.2	Q-Learning	33
8.3	Temporal Difference Learning	34
8.3.1	SARSA	34
8.4	Challenges	34
9	Deep Learning in Computer Vision	35
10	Neural Network Applications in Machine Translation	36
10.1	Conventional Statistical MT	36
10.2	Neural MT	37

11 Speaker Independence	40
11.1 Frequency Invariance	41
11.2 Multi-Speaker Reference Model	41
11.3 Cross-Language DNNs	42
12 Hand Writing	44
13 Natural Language Processing	45
13.1 Language Model	45
13.1.1 Word Embeddings	46
13.1.2 Word2Vec	46
13.2 Sentence Modeling	46
13.2.1 Dynamic k-max Convolutional NN	47
14 Gradient Optimizations and 2nd order Methods	49
14.1 Logistic Regression	49
14.1.1 Gradient Descent - General Approach	49
14.1.2 Batch Gradient Descent	50
14.1.3 Stochastic Gradient Descent	50
14.1.4 Mini-Batch Gradient Descent	50
14.2 Learning Rate Scheduling	50
14.2.1 Adagrad - ADaptive GRADient method	51
14.2.2 Adadelta	51
14.2.3 RMSprop	51
14.2.4 Adam - ADaptive Moment estimation	52
15 Error Functions	53
15.1 Binary Cross Entropy or Negative Log Likelihood	53
16 Activation Functions	54
16.1 Step Function	54
16.2 Linear Function	54
16.3 Logistic / Sigmoid Function	54
16.4 Hyperbolic Tangent function	55
16.5 Softmax Function	55

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

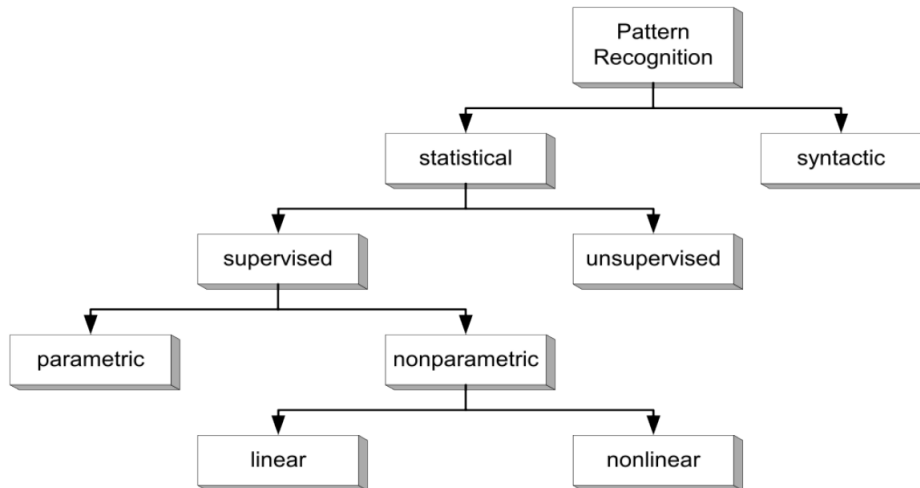
Welche Probleme können beim Entwurf von neuronalen Netzen auftreten?

- 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



- supervised:
 - labels / classes of training data are known
 - train a classifier
- unsupervised:
 - labels / classes are not known
 - find the structure of the data
 - methods: clustering, autoassociative nets
- parametric:
 - assume underlying probability distribution
 - estimate the parameters
 - eg.: Gaussian classifier
- unparametric:
 - doesn't assume underlying probability distribution
 - estimate probability of error from training data
 - nearest neighbours, parzen window, perceptron

	×	÷	+	−	<	>	=	≠	
Nominal							✓		Qualitative
Ordinal					✓		✓		Qualitative
Interval			✓		✓		✓		Quantitative
Ratio	✓		✓		✓		✓		Quantitative

2.1 Data Feature

2.1.1 Preprocessing:

- Normalization
 - Max-Min (Rescaling): $x' = \frac{x - \min(x)}{\max(x) - \min(x)}$
 - Standardization: $x' = \frac{x - \bar{x}}{\sigma}$
 - unit length scaling: $x' = \frac{x}{\|x\|}$
- incomplete data: null filling - smoothing

2.1.2 Comparison

- Edit Distance
- L_p Norm
- Probabilistic Distribution comparison
 - Cross Entropy
 - Kullback-Leibler Divergence

2.2 Parametric Methods

2.2.1 Bayes Decision Theory

$$P(\omega_j|x) = \frac{p(x|\omega_j) \cdot P(\omega_j)}{p(x)}$$

$$p(x) = \sum_j p(x|\omega_j) \cdot P(\omega_j)$$

Since we have overlapping probability distribution false classifications occur. Simply speaking to reduce false classification always decide for the class with the highest probability (errors still can occur!).

2.2.2 Mixture Densities

needed???

2.3 Classifier Discrimination Function

The *classifier discrimination function* iterates over all classes and chooses those which has the highest probability. In the case of the Gaussian classifier it looks as follows:

$$\begin{aligned} g_i &= P(\omega_i|x) \\ &= p(x|\omega_i) \cdot P(\omega_i) \\ &= \log(p(x|\omega_i)) + \log(P(\omega_i)) \end{aligned}$$

Problems:

- limited training data
- limited computation
- class-labeling potentially costly and errorful
- classes or features might not be known

Parametric solution: assume that $p(x|\omega_i)$ has a parametric form. The most common representative: multivariate normal density.

Univariate normal density:

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi \sigma}} \cdot e^{(-\frac{1}{2} \cdot \frac{x-\mu}{\sigma})^2}$$

Multivariate normal density:

$$p(\underline{x}) = \frac{1}{\sqrt{2\pi}^d \cdot \sqrt{|\Sigma|}} \cdot e^{-\frac{1}{2} \cdot \frac{(\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})}{\Sigma}}$$

The following form takes effect for the classifier discrimination function:

$$g_i(\underline{x}) = -\frac{1}{2} \cdot \frac{(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i)}{\Sigma_i} - \frac{d}{2} \cdot \log(2\pi) - \frac{1}{2} \cdot \log(|\Sigma_i|) + \log(P(\omega_i))$$

For the Gaussian classifier we have to estimate the *covariance matrix* Σ_i and the *mean vector* μ_i . To estimate the parameters we can use the *MLE: maximum likelihood estimation* for the univariate case. For the multivariate case we can use:

$$\begin{aligned} \underline{\mu} &= \frac{1}{N} \cdot \sum_{k=1}^N \underline{x}_k \\ \Sigma &= \frac{1}{N} \cdot \sum_{k=1}^N (\underline{x}_k - \underline{\mu}) \cdot (\underline{x}_k - \underline{\mu})^T \end{aligned}$$

2.4 Curse of Dimensionality

The problem with the classifier design in general is that we can not say which features are the most valuable ones, are more features better and is features are useful will they be ignored?

However in general we can say that more features perform worse because we have limited training data and we still have to estimate the parameters (Eg. 1000 sample training data and 1000 parameters to estimate). We have to select the best feature and might have to reduce the dimensions (for example with the *Principle Component Analysis (PCA)*).

2.4.1 Principle Component Analysis (PCA)

The idea is that single dimension are correlated and we want to remove the dimension with the least informations:

1. find the axis with the highest variance
2. rotate the space along the axis
3. not the dimensions are uncorrelated
4. remove the dimension with the lower variance

2.5 Non-Parametric Methods

Non-Parametric Methods do not assume any distributions and try to find structures in the data itself.

2.5.1 Parzen Window

1. chose a window with volume V
2. count the numbers of samples in that window $p(x) = \frac{k}{N}$ where k are the samples in the window and N is the total numbers of samples.

The hard part is to chose the window size. the thumb rule is: $V_n = \frac{1}{\sqrt{n}}$

2.5.2 k-nearest Neighbours

Idea: Let the volume be a function of the data. Include k-nearest neighbors in estimate. Set $k = \sqrt{n}$ k-nearest neighbor rule for classification.

1. Find k-nearest neighbors of x
2. Determine the class most frequently represented among those k samples (take a vote)
3. Assign x to that class.

The problem with k-nearest neighbour is k. k should be large for a reliable estimate but also small to guarantee that all k neighbours are reasonably close.

There is no data like more data.

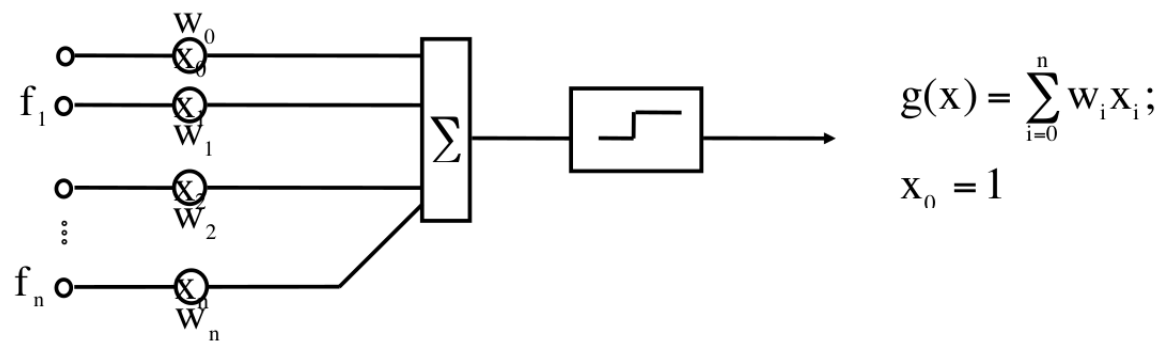
2.5.3 Clustering

needed?

2.6 Fisher-Linear Discriminant

needed??

2.7 Perceptron



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: ??

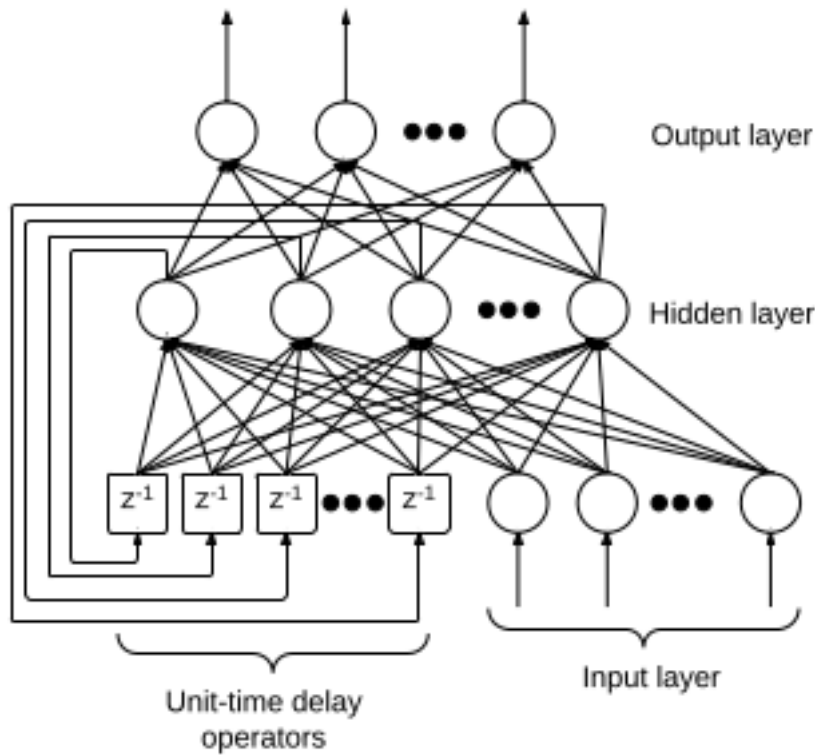
3.2 Elman vs. Jordan Networks - Simple RNN

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.

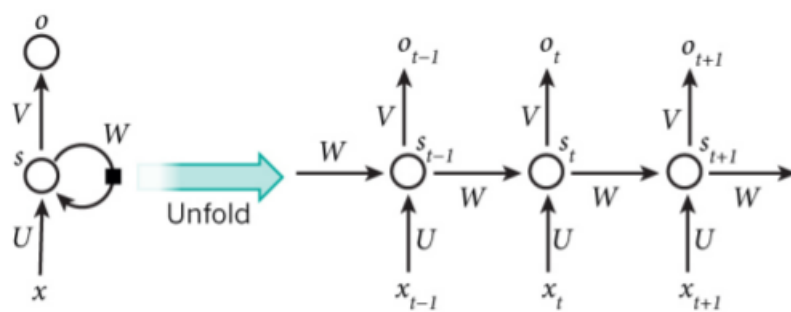


3.4 Training

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 Backpropagation

Improve BP

- Parallel Hardware
- Efficient Implementation
- Faster Gradient Descent Search
- Selective Choice of Patterns
- Efficient Architectures

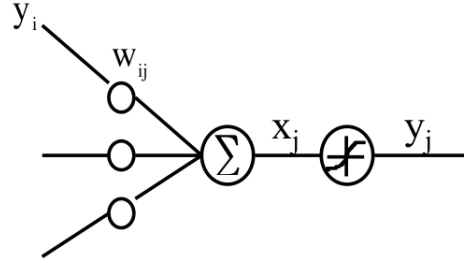
$$E = \frac{1}{2} \sum_j (y_j - d_j)^2 \quad y_j = \frac{1}{1 + e^{-x_j}} \quad x_j = \sum_i y_i w_{ij}$$

$$1). \quad \frac{\partial E}{\partial y_j} = y_j - d_j$$

$$2). \quad \frac{\partial E}{\partial x_j} = \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} y_j [1 - y_j]$$

$$3). \quad \frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial x_j} \cdot \frac{\partial x_j}{\partial w_{ij}} = \frac{\partial E}{\partial x_j} \cdot y_i$$

$$4). \quad \frac{\partial E}{\partial y_i} = \sum_j \frac{\partial E}{\partial x_j} \cdot \frac{\partial x_j}{\partial y_i} = \sum_j \frac{\partial E}{\partial x_j} \cdot w_{ij}$$



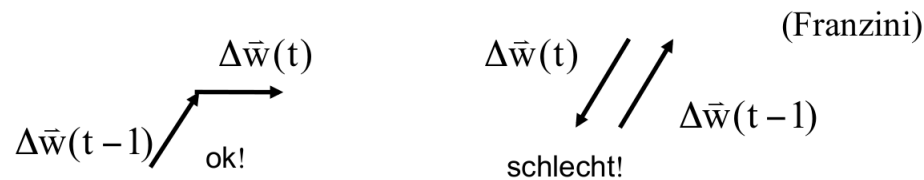
The weight update faces different problems: vanishing / exploding gradient, small weights update (slow learning), etc. One way to fix some of the problems is to set a *step size* and *momentum* to the weight update:

$$\Delta w_{ij}(t) = -\epsilon \frac{\delta E}{\delta w_{ij}(t)} + \alpha \Delta w_{ij}(t-1)$$

ϵ is the step size (or learning rate) and α is the momentum. The momentum says how much of the last weight update should be added to the current weight update. This will prevent slow training a little.

Eine weitere Beschleunigung des Trainings erhält man indem für bestimmte Beispiel kein BP angewende wird. Hierbei muss der Fehler groß genug sein, ansonsten wird der Fehler ignoriert.

Ebenso können wir die Lernrate dynamisch anpassen, um zu verhindern, dass weight updates hin und her springen.



$$\cos \Theta = \frac{\sum_{i,j} (\Delta w_{ij}(t-1) \cdot \Delta w_{ij}(t))}{\sqrt{\sum_{i,j} (\Delta w_{ij}(t-1))^2 \cdot \sum_{i,j} (\Delta w_{ij}(t))^2}}$$

$$\cos \Theta = 1 \quad (\Theta = 0^\circ, 360^\circ, \dots)$$

$$\cos \Theta = -1 \quad (\Theta = 180^\circ)$$

=> ϵ größer machen

=> ϵ kleiner machen



$$\epsilon(t) = \epsilon(t-1) \cdot \text{const} \cdot \frac{\cos \Theta + 1}{2}$$

Um ein MLP schnell trainieren zu können, kann sich der *Quickprop Algorithmus* eignen:

$$\Delta w(t) = \frac{s(t)}{s(t-1) - s(t)} \cdot \Delta w(t-1)$$

wobei $s(t) = \frac{\delta E}{\delta w_{ij}}(t)$. Jedoch hängt das schnelle Training nicht nur vom Lernalgorithmus ab. The initial weight should be initialst by looking at the activation function and look for the intervall with the biggest derivant. **Generalisierung:** The generalization states how well the system will perform on the validation set.

Reasons:

- Overfitting
- too much parameters or too less trainings data

- wrong topology

Generalization for linear Systems:

$$\langle \epsilon_{test} \rangle = \langle \epsilon_{train} \rangle + 2\sigma^2 \frac{p}{n}$$

Generalization for linear Systems:

$$\langle \epsilon_{test} \rangle = \langle \epsilon_{train}(\lambda) \rangle + 2\sigma_{eff}^2 \frac{p_{eff}(\lambda)}{n}$$

Methods to better generalization:

- reduce network complexity (weight decay, weight elimination, optimal brain damage, optimal brain surgeon): prevent overfitting
- stepwis increase of size of network (cascade correlation, meiosis network, automativ structure optimization): increasae size of network so it can learn

4.1 Regularization

Regularization is used to prevent overfitting.

4.1.1 Weight Elimination

Weight Elimination is added to the error to punish the network for large weights (λ : 0.0001...0.001):

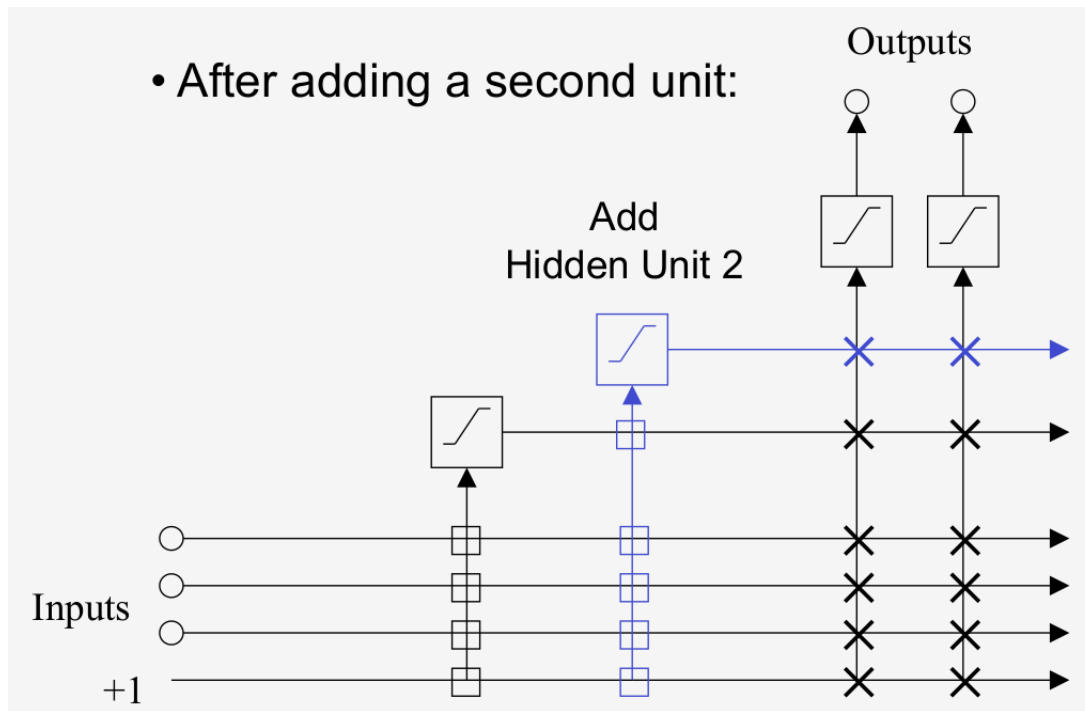
$$E = MSE + \lambda \sum_{i,j} \frac{w_{i,j}^2}{1 + w_{i,j}^2}$$

- - slow learnin
- - worse performance on trainingsdata
- + better generalization

4.1.2 Optimal Brain Damage

Idea: remove certrain connection of the network to reduce the complexity of the network and prevent overfitting. Easy: remove connections with very small $|w_{i,j}|$. Better: remove thos connections which influence the error the smallest (unimportant weights). Therefore calculate:

$$\frac{\delta E}{\delta w_{ij}^2}$$



4.1.3 Cascade Correlation

- can create deep networks without a drastic slowdown
- amount hidden units does not need to be estimated empirisch
- each time step only one layer of connections is trained
- learns very fast
- incremental learning

4.1.4 Meiosis Network

Idea: adding of hidden units depends on the "uncertainty" of the network. The mean and variance is learned.

$$w_{ij}^* = \mu(w_{ij}) + \sigma(w_{ij})\phi(0, 1)$$

Start with one hidden unit and split unit if

$$\frac{\sum_i \sigma_{ij}}{\sum_i \mu_{ij}} > 1.0 \text{ and } \frac{\sum_k \sigma_{ik}}{\sum_k \mu_{ik}} > 1.0$$

5 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 gereihe 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.

5.1 Speech Recognition

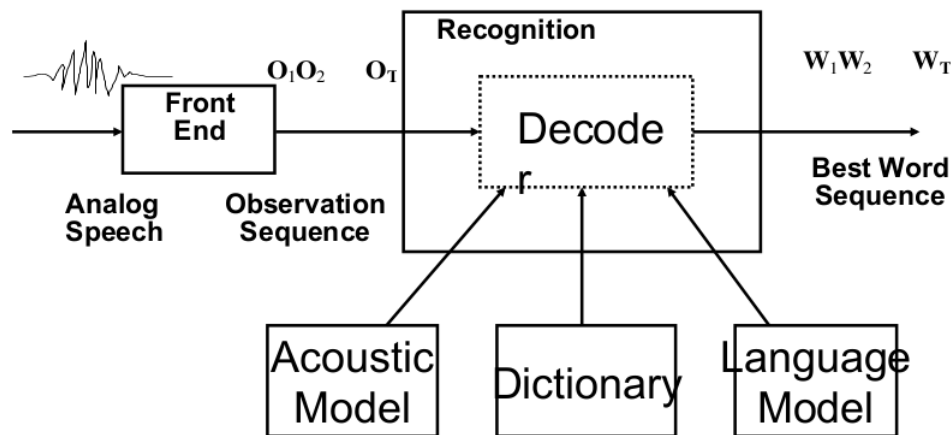


Abbildung 1: Komponenten eines Sprach Erkenners

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*:

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

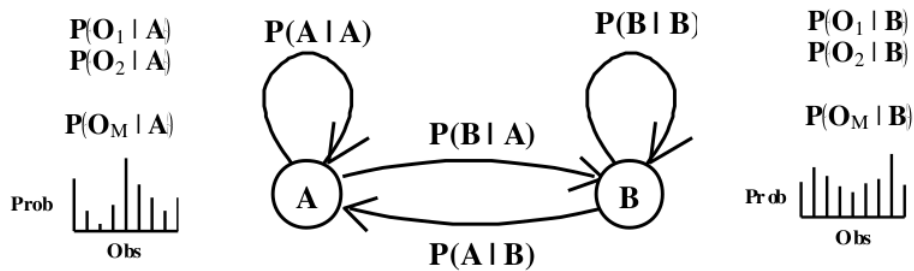
$P(W|A)$ ist das *acoustic model*, $P(W)$ ist das *language model* und $P(A)$ ist konstant für einen ganzen Satz.

5.2 Acoustic Model

5.2.1 Hidden Markov Model

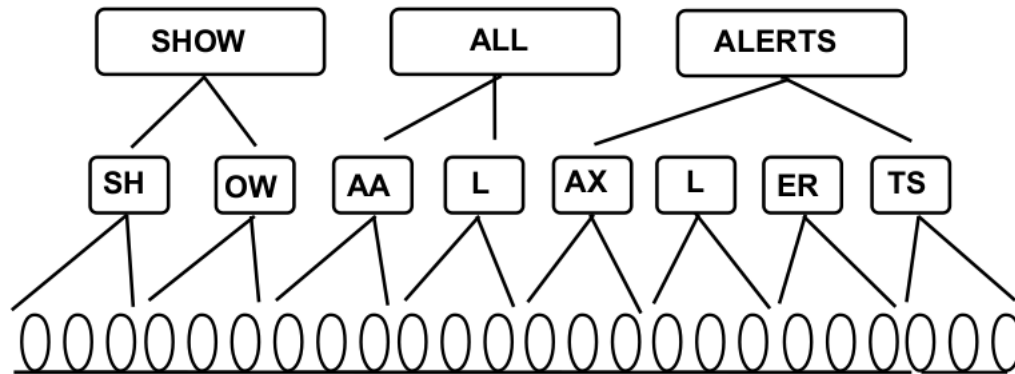
Elements:

- 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

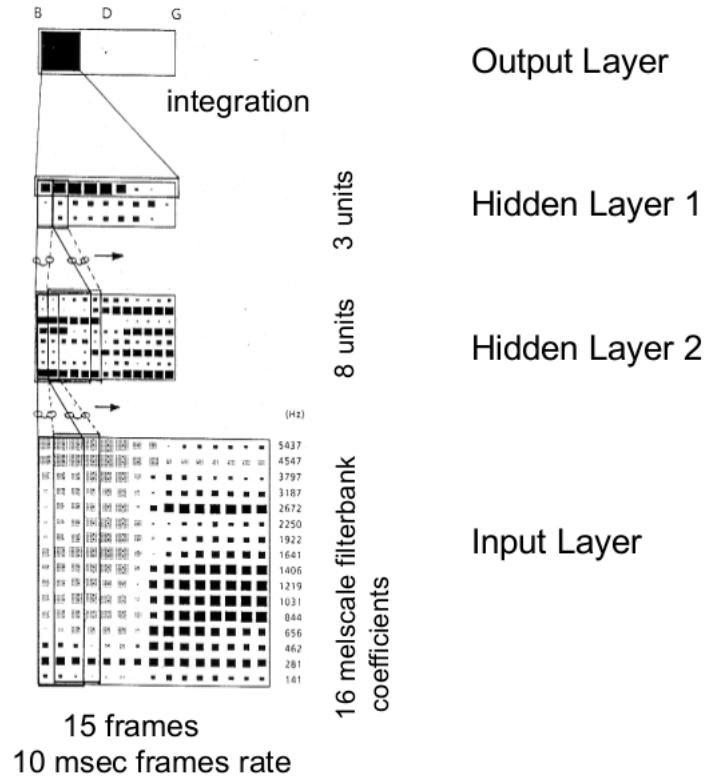


Die Wörter

werden zerlegt in Phoneme und die Phoneme in Zustände.

5.2.2 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 5.3.2. 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.

5.3 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)

5.3.1 Time Alignment

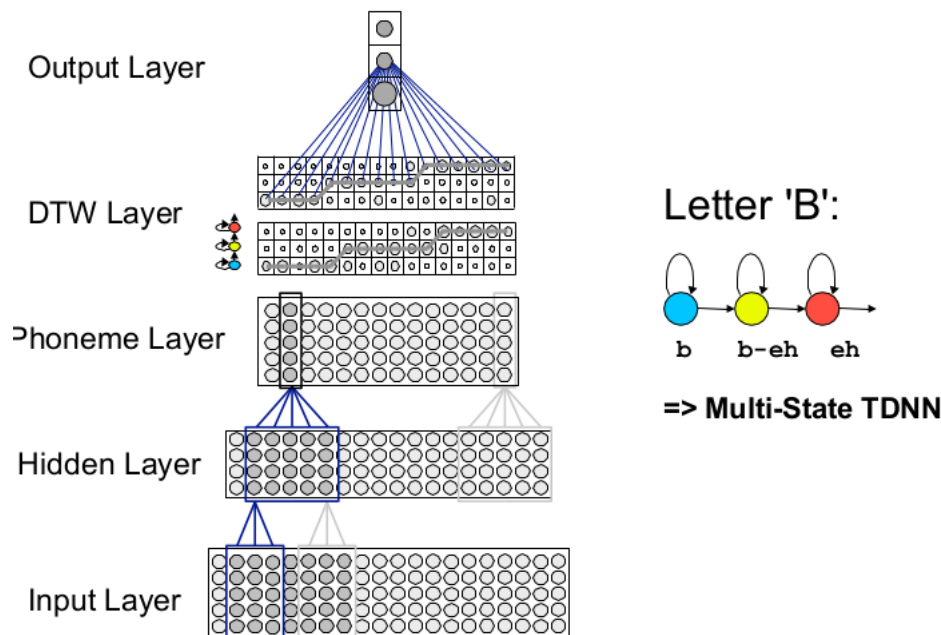
Hierbei geht es darum, dass verschiedene Sprecher unterschiedlich schnell/langsam sprechen. Irgendwie muss beim Lernen einen zeitlichen Zusammenhang herausgearbeitet werden.

Beim *Linear Sequence Alignment* wird das Referenzsignal und das unbekannte Signal linear aneinander angepasst (Stauchen / 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.

5.3.2 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 dazu den Hybriden, welche danach dran kommen, ist, dass das es sich hierbei um hidden Features handelt.

5.3.3 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.)

5.3.4 Viterbi-Algorithmus

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

5.3.5 Forward-Algorithmus

5.3.6 Backward-Algorithmus

6 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.

6.1 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.

6.1.1 Applications

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

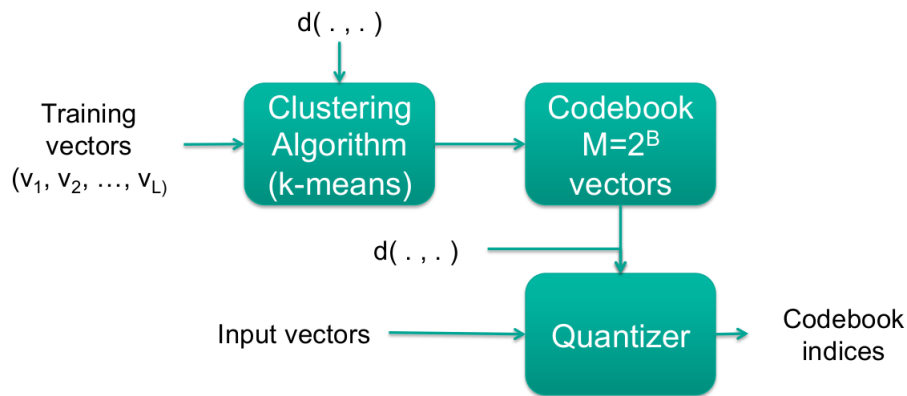
6.1.2 Training

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.

6.2 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:

6.2.1 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 $i \neq c$

6.2.2 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

6.2.3 LVQ2.1

6.2.4 LVQ3

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

6.2.5 OLVQ

7 Self-Organizing-Maps

7.1 Principles Of Self-Organized-Learning

7.1.1 Principle 1: Self-Amplification

7.1.2 Principle 2: Competition

7.1.3 Principle 3: Cooperation

7.1.4 Principle 4: Structural Information

7.1.5 Hebbian Learning

7.2 Self-Organizing-Maps

SOM theory structure learning properties VQ kernel SOM applications

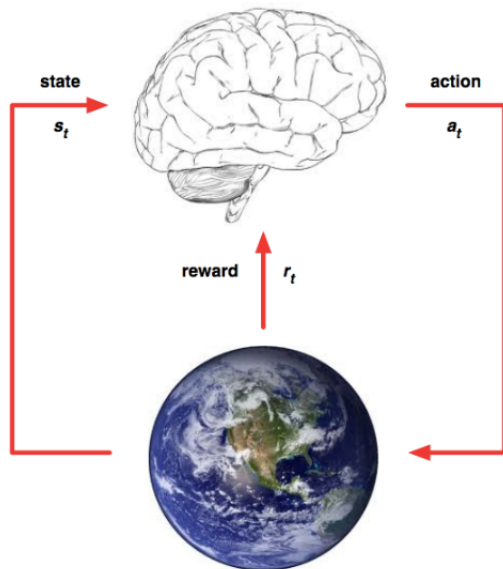
8 Reinforcement Learning

- agent can act
- actions change the agent's future state
- scalar rewards for success
- (sequential decision making

Select actions to maximize the future reward.

Examples of Reinforcement Learning:

- Fly stunt maneuver in a helicopter
- Defeat the world champion at backgammon
- Manage an investment portfolio
- Control a power station
- Make a humanoid robot walk
- Play many different Atari games better than humans



- At each step t the agent:
 - Receives state s_t
 - Receives scalar reward r_t
 - Executes action a_t
- The environment:
 - Receives action a_t
 - Emits state s_t
 - Emits scalar reward r_t
- The evolution of this process is called a Markov Decision Process (MDP)

Policy $a = \pi(s)$: probability distribution of actions given a state

Value function $Q^\pi(s, a)$: expected total reward from state s and action a under policy π

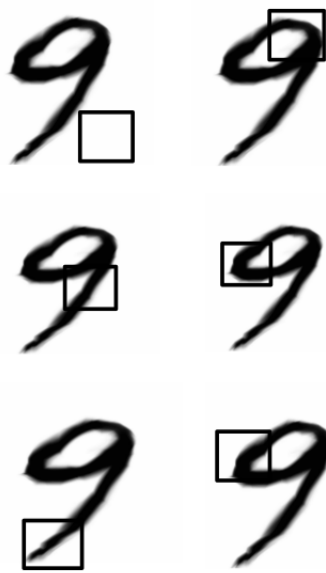
$$Q^\pi(s, a) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots | s, a]$$

Policy-based RL: search directly for the optimal policy π (achieving maximum future reward)

Value-based RL: estimate optimal value function $Q^*(s, a)$ (maximum value achievable under any policy)

To train these networks we use *Monte Carlo methods* to learn from experience. **Policy-based Reinforcement Learning:** parameterize the policy (for example with a neural network) and based on the cumulative rewards, update the policy. We need to get the gradient of the rewards with respect to the policy:

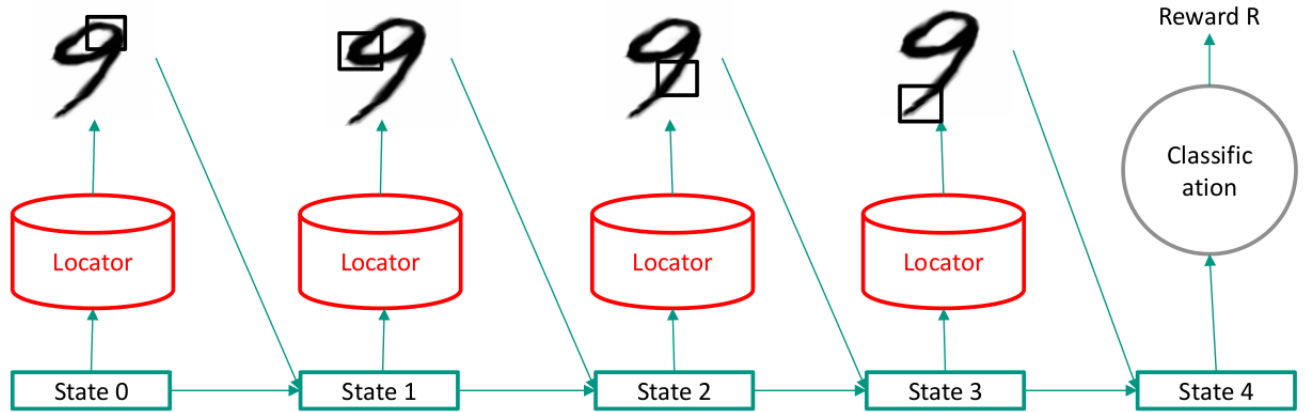
- Policy Gradients algorithm
- Online: update after episode
- Offline: update while in episodes



- Where do you look to see if this is a 9 or an 8?
- We look at the image with some “attention” mechanism: focusing on a certain position at one time
- Changing our focus on multiple time steps until we have information to classify the image
- Reinforcement Learning can help learn where to look to classify handwritten numbers

For the image shown we can use a RNN to determine what our current state is given the last state. RL for where should we look next, given our current states.

RL gives use an reward on which we can calculate the backward pass (reward represents the gradient).



$$\blacksquare \frac{\Delta R}{\Delta p} = a * (R - b) * \frac{\Delta \ln f(x,p)}{\Delta p}$$

- a : Scaling factor
- b : Baseline reward used to reduce the variance of the gradient
- x : The sampled value (position in the image)
- $f(x,p)$: the Probability density function, learnt by the neural network

8.1 Bellman Equation

$$Q^\pi(s, a) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots | s, a]$$

Recursively:

$$Q^\pi(s, a) = \mathbb{E}_{s'}[r_{t+1} + \gamma Q^\pi(s', a') | s, a]$$

Optimal value function:

$$Q^*(s, a) = \mathbb{E}_{s'}[r_{t+1} + \gamma Q^*(s', a') | s, a]$$

Value iteration solve the Bellman Equation:

$$Q_{i+1}(s, a) = \mathbb{E}_{s'}[r_{t+1} + \gamma Q_i(s', a') | s, a]$$

8.2 Q-Learning

Objective function by *mean-squared error* in Q-values:

$$L(w) = \mathbb{E}[(r + \gamma \max_{a'} Q(s', a', w') - Q(s, a, w))^2]$$

Q-learning gradient:

$$\frac{\delta L(w)}{\delta w} = \mathbb{E}[(r + \gamma \max_{a'} Q(s', a', w') - Q(s, a, w)) \frac{\delta Q(s, a, w)}{\delta w}]$$

8.3 Temporal Difference Learning

General update rule:

$$Q(s_t, a_t) += learning_rate \cdot (td_target - Q(s_t, a_t))$$

8.3.1 SARSA

$$Q(s_t, a_t) += \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Q-value for a state-action is updated by an error, adjusted by the learning rate alpha.

TD Target for SARSA:

$$R_{t+1} + discount_factor \cdot Q(s_{t+1}, a_{t+1})$$

TD Target for Q-Learning:

$$R_{t+1} + discount_factor \cdot \max Q(s_{t+1}, a_{t+1})$$

8.4 Challenges

- Data is sequential
- policy changes rapidly with small changes to Q-values
- scale of gradients and reward us unknown at the onset

9 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

10 Neural Network Applications in Machine Translation

NLP vs MT: NLP covers all part up until we have a sentence. MT then takes this sentences and tries to translate it into another language.

Challenges of NLP:

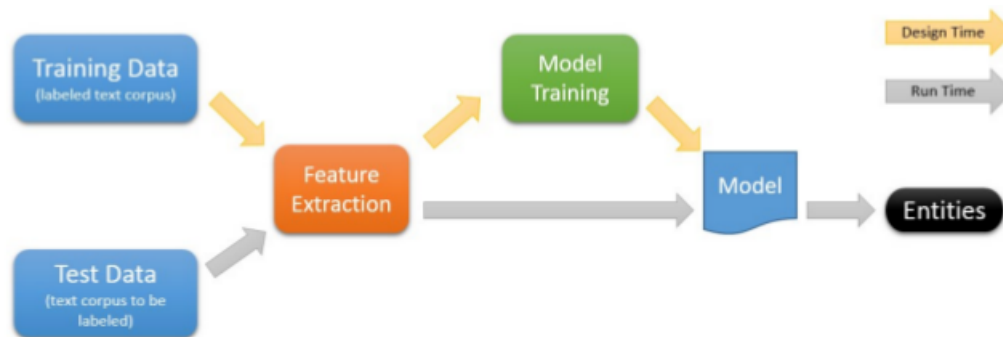
- Disambiguation: word sense, structural
- co-references: referring to objects within and across sentence boundaries
 - anaphora: The man goes to work. He takes the bus
 - deictic references: here, now, I, you
 - references: same object using a synonym, hypernym

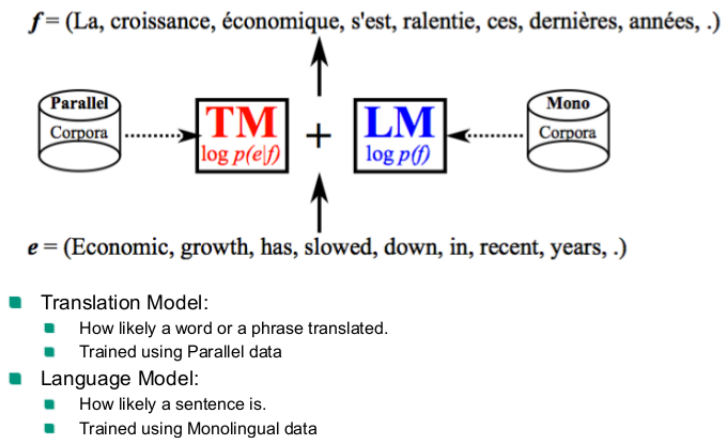
Challenges of MT:

- translation divergences: different meaning of components in different languages
 - structural divergences: word order
 - thematic divergences: changes of grammatical role
 - head switching: ??
 - lexicalization: swim across -> to cross swimming
 - categorial: a little bread -> un poco de pan
 - collocational: make a decision -> eine Entscheidung treffen
- translation mismatches: En: fish -> Spa: pez, pezcado

10.1 Conventional Statistical MT

The conventional statistical MT learns automatically from the data by breaking the big problem into smaller problems. MT: Tranlation models, language models, reordering. Problems:



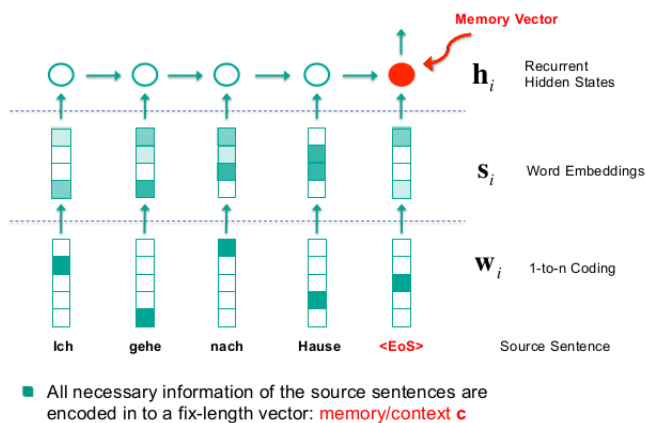


- requires in-depth expertise
- limited ability to generalize
- simple learning methods: cannot model well

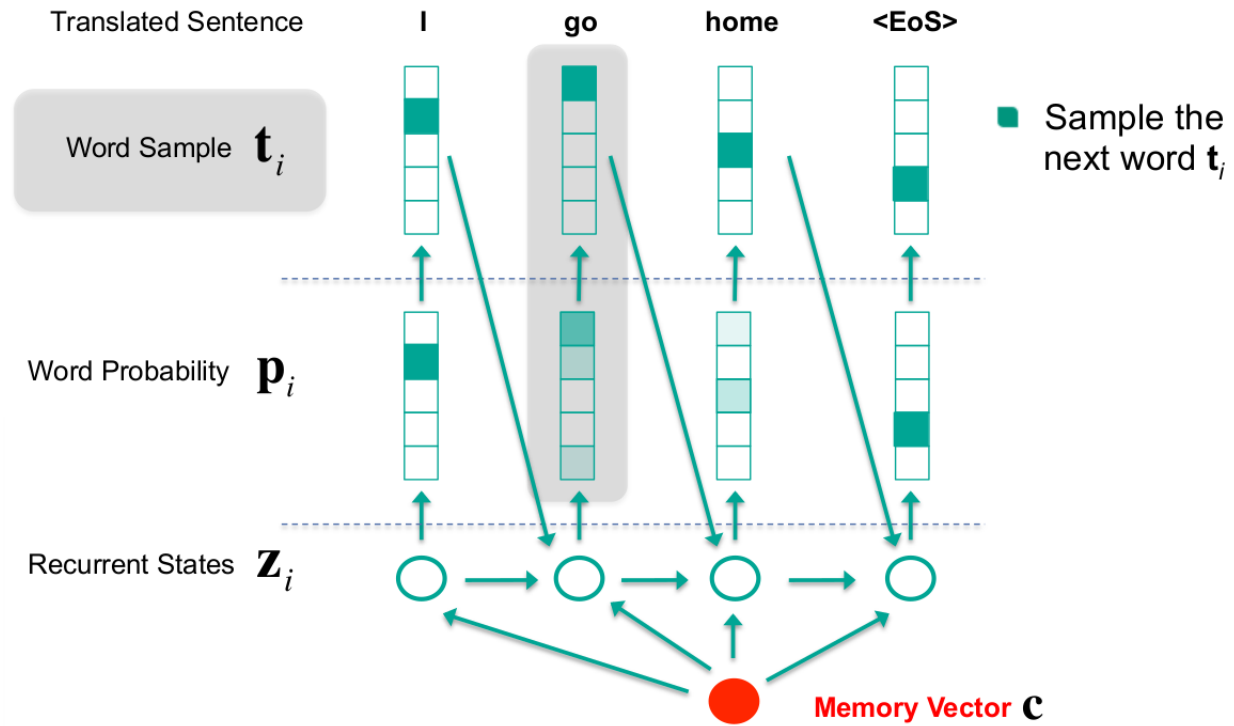
10.2 Neural MT

Idea: train model to represent source sentence and generate target sentence from source representation.

Encoder: read input, represent content in hidden fix dimension vector (LSTM-based model): **De-**



coder: generate output, use fix dimension vector as input: Attention mechanism to set focus on important parts:



Context vector: Represent source sentence with focus /attention on some part of the sentence

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

Attention: Probability that this part of the source sentence is important (based on energy function)

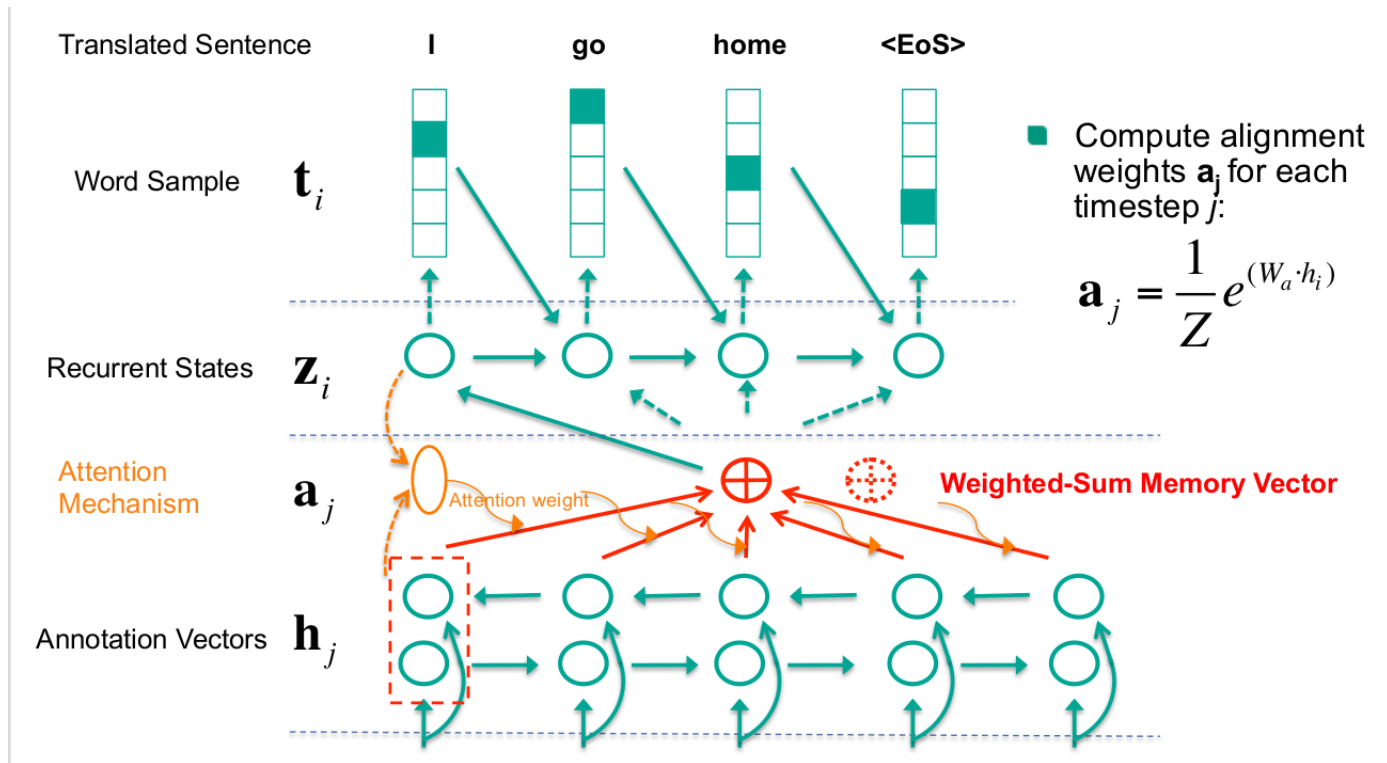
$$\alpha_{ij} = \frac{e^{e_{ij}}}{\sum_{k=1}^{T_x} e^{e_{ik}}}$$

Energy function:

$$e_{ij} = a(s_{i-1}, h_j)$$

Bi-directional Recurrent Encoder (BiRNN): look similar to the normal RNN but the hidden state goes in both directions. Two hidden states are concatenated by the *annotation vector*.

Encoder-Decoder-Training: the source sentence is an English sentence, the target sentence is a German sentence. We use mini-batch training (like 100 sentences) and average the update! We can



also calculate all update from all sentences in parallel. However we might have different sentence length.

Sentence probability:

$$p(e|f) = \prod_{j=1}^n p(e_j|f, e_1^{j-1})$$

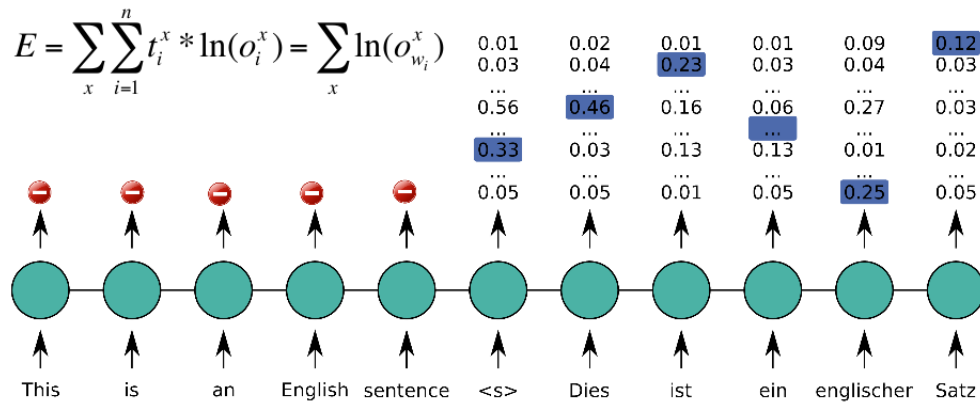
Translation:

- Calculate output probabilities
- Select best n translations
- Extend all hypothesis in beam
- Prune hypothesis not in beam

■ Calculate Forward path

■ Measure error

■ Cross entropy error:



11 Speaker Independence

Normal TDNNs are time invariance. However men, women and children differ in *frequency*. Men normally have a darker voice than women and children. Therefore we have to compensate this invariance.

Observations on TDNNs

- TDNNs develop linguistically plausible features in the hidden units
- TDNNs developed alternate internal representations that can link quite different acoustic realizations to the same higher level concept (because of multilayer arrangement)
- hidden units fire synchrony because they operate independent of precise time alignment or segmentation (time invariant)
- small network output may not be useful in complex task (but the internal abstractions may be valuable)
- complex concepts -> use stages with different knowledge
- new learning strategies should be build in existing knowledge

Model invariance

- frequency shift, tilt, compression

Variability

- adaption
 - slow adaption -> modify weights
 - fast adaption -> pretrained specific submodels
- normalization
 - environment: to the room
 - speaker: mapping new speaker to standard speaker (with standard sentences)

Combine **two standard TDNNs** to one (overall better classification): one with MSE and the other with CFM. Those combination yields the correct classification.
Even better than two TDNNs: **three TDNNs**! The third TDNN uses CE.

11.1 Frequency Invariance

Convolutional Acoustic Models:

- parameter sharing across spectrum and time -> exploit 2D structure of features
- upper layer fully connected
- pooling gives more robustness and less overfitting

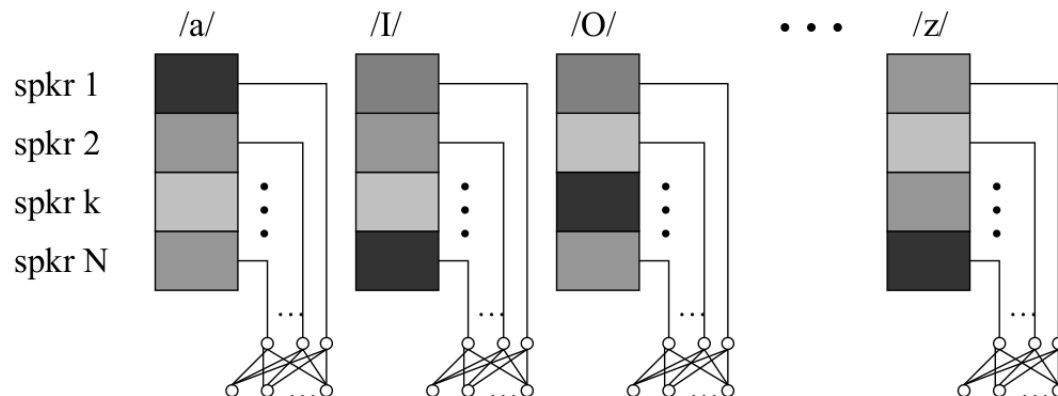
11.2 Multi-Speaker Reference Model

Idea: A speaker-specific reference model is composed from several well trained reference models:
Meta-Pi-Net:

A separate neural net *Meta-Pi-Net* controls the activation / deactivation of the single nets. The whole net is a combination of one net per speaker. In general we look which speaker net is closest to the input (which nets classifies the best). The Meta-Pi-Net produces weights for each net which control how much each net is taking into account for the classification.

It is not important that the correct speaker is choosing. Sometimes a combination of two or more speaker might fit better to the input. The actual correct speaker net might not be even taking into account (results in a mixture of the other nets).

Speaker Normalization: we have speaker dependent nets. Now when we use those nets to classifier what another speaker says (no dependent net for this speaker!) the error rate is 41.9%. With 40 text-dependent training sentences the error rate is reduced to 6.8%.



phoneme specific
reference model selection networks

i-vectors: identity vectors (i-vectors) describe the speaker-characteristic offset to an universal background model. Those i-vectors are used to train a recognition system.

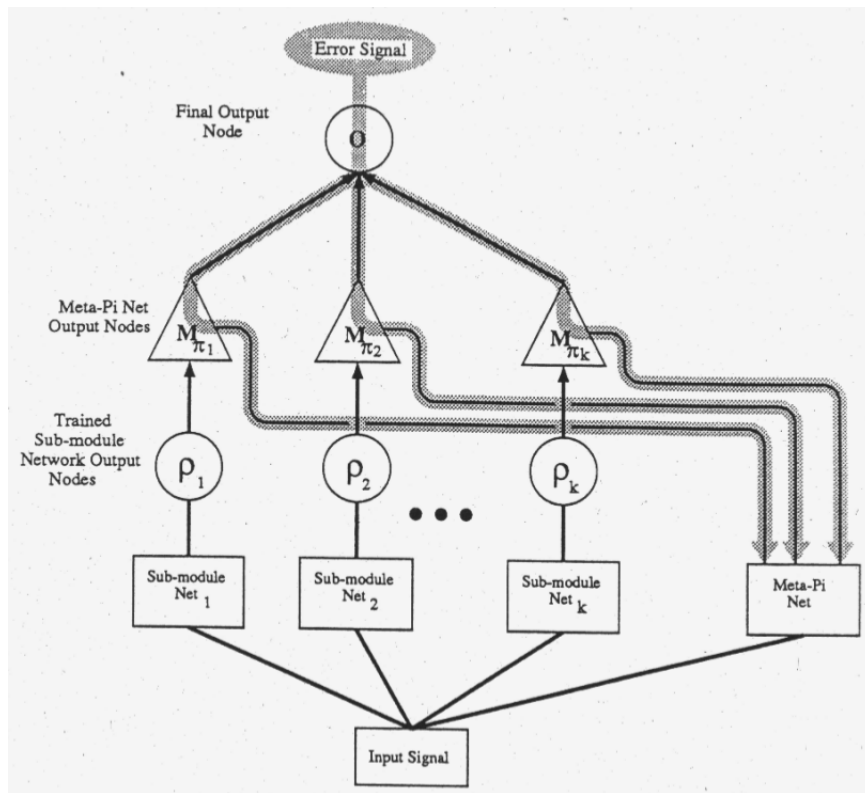
Speaker Adaptive Training of DNN:

1. train initial DNN with (and keep it fixed)
 - SI features: fbank
 - SA features: Feature space Maximum Likelihood Linear Regression (fMLLR)
2. Train an i-vector NN
 - inputs: i-vectors
 - outputs: linear shift to the original feature vectors
 - added features become speaker-normalized
3. update the DNN in the new feature space, i-vector NN fixed -> yields SAT-DNN

11.3 Cross-Language DNNs

Multilingual Bottleneck Features:

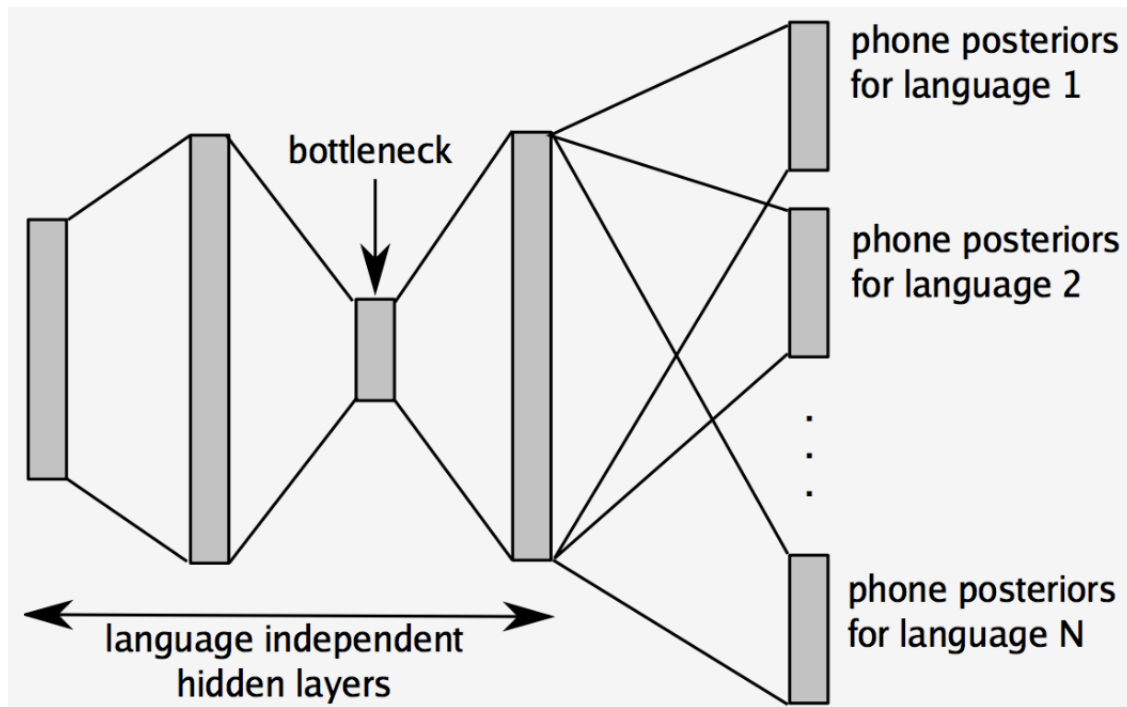
- Humans can only produce a finite amount of different sounds



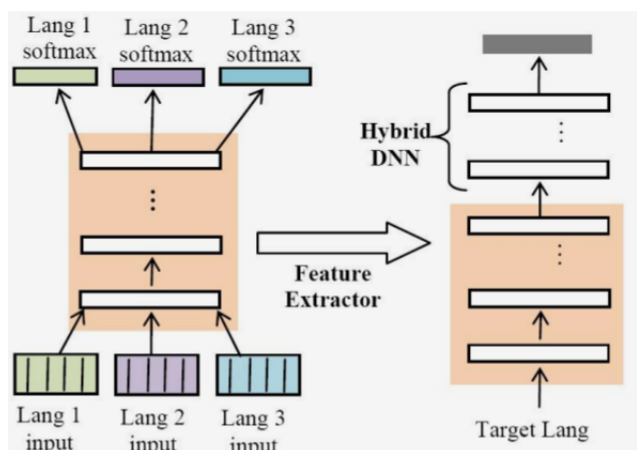
- Subset of sounds is used in individual languages
- Some sounds are used in different languages
- Share data of the same sounds from different languages
- Extract more robust features using data with more variability

Cross-Language DNNs with Language-Universal Feature Extractors:

DNNs can be trained on multiple languages. The hidden layers (*language-universal feature extractor*) are used in other DNNs.



12 Hand Writing



13 Natural Language Processing

NN applications in NLP

- Classification/Sequence labeling problems
- Language Modeling and Word Embeddings
- Sentence Modeling
- Machine Translation
- Dialog Systems

Tasks Of NLP:

- Part-Of-Speech: label each word
- Chunking: label segments of a sentence (noun or verb phrase)
- Named Entity Recognition: label atomic elements (person, locatin)
- Semantic Role Labeling: label elements of a sentence with a semantic role

Window of n-gram: k words before, k words after currentl considering word: $n = 2 * k + 1$ **1-of-N coding(one-hot):** a vocabulary-size-dimensional vector, 1 in the index of the word, 0 otherwise.

13.1 Language Model

Model how fluent a sentence is. Given sentence S , calculate $P(S)$. Given incomplete sentence H , calculate $P(\text{wordX} | H)$. Probability of the next word given a history (MLE):

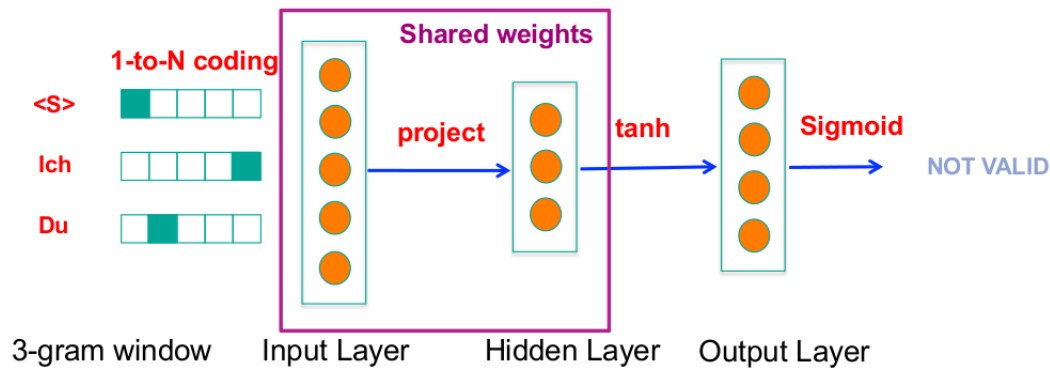
$$P(w_i | H) = P(w_i | w_1 \dots w_{i-1}) = \frac{\text{count}(w_1 \dots w_i)}{\text{count}(w_1 \dots w_{i-1})}$$

Conventional n-gram approach (only n previous words):

$$P(w_i|H) \approx P(w_i|w_{i-n+1}...w_{i-1}) = \frac{\text{count}(w_1...w_i)}{\text{count}(w_{i-n+1}...w_{i-1})}$$

Probability of the whole sentence S:

$$P(S) = P(w_1...w_N) \approx \prod_{i=1}^N P(w_i|w_{i-n+1}...w_{i-1})$$



- From one-hot (orthogonal, binary, sparse $|V|$ -dim space)
 - Car=(0, 1, 0, 0, 0, 0) ### Automobile=(0, 0, 0, 0, 0, 1)
- To a (dim-reduced) shared word embedding space
 - Shared projection (Lookup Table)

13.1.1 Word Embeddings

Spatial distance corresponds to word similarity: $\text{we}(\text{car})$ close to $\text{we}(\text{automobile})$.

Vector operations can be used to combine word meanings: $\text{we}(\text{king}) - \text{we}(\text{queen}) \approx \text{we}(\text{man}) - \text{we}(\text{woman})$

13.1.2 Word2Vec

needed??

13.2 Sentence Modeling

Representation of word sequences: Need fixed-length representation vectors for variable-length sequences.

Approaches:

- Simple aggregation: Sum/Average/Max?
- Recurrent architectures: One word vector at a time step
- Recursive NN: Compositions using syntax parsed tree
- Time-Delay NN or Convolutional NN

Convolutional:

input sequence: $\mathbf{s} \in \mathbf{R}^s$

Filter vector: $\mathbf{m} \in \mathbf{R}^m$ (?????)

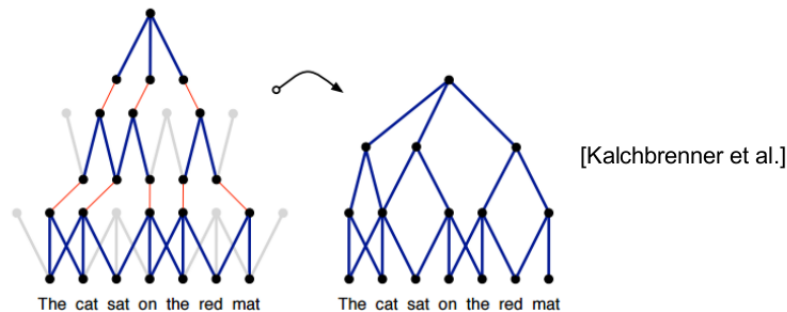
Take the dot product from \mathbf{m} with each m-gram in \mathbf{s} to obtain a sequence \mathbf{c}_j :

$$\mathbf{c}_j = \mathbf{m}^T \cdot \mathbf{s}_{j-m+1:j}$$

After convolution step we have a matrix \mathbf{c}_j (sequence of vectors) but still be a variable-length sequence (on s): $\mathbf{c}_j \in \mathbf{R}^{m \times s}$

Max Pooling: Take the max value of each row in \mathbf{c}_j . Now we have a fixed-length vector (on m). Exactly the same as Time-Delay Neural Networks (TDNN)

13.2.1 Dynamic k-max Convolutional NN



- Series of convolutional and pooling layers
- Dynamic k-max pooling:
 - k-max: returns k max values instead of one.
 - Dynamic: k is dynamically chosen based on some aspects of the network

Advantages

- A filter is a linguistic feature detector of a class of n-grams
- Features are extracted independently from their positions.
- Higher filters can capture syntactic and semantic info.

Dynamic k-max pooling:

$$k_I = \max(k_{top}, \lceil \frac{L-I}{L} S \rceil)$$

k_{top} : fixed k_{top} pooling layer on top

I, L : convolutional layer I over total L convolutional layers where the k^{th} applied on.

Properties:

- Sensitivity to word and n-gram order
- Have the ability to induce feature graph

14 Gradient Optimizations and 2nd order Methods

Idea: Start from a point close to the wanted solution and iteratively move to the point that makes the gradient of the function approach zero (hence *decent*). All weights are initialised with small random numbers.

Important: update the parameters in the **opposite** direction of the gradient:

$$\mathbf{w} = \mathbf{w} - \eta \nabla_{\mathbf{w}} L(\mathbf{x}, \mathbf{w})$$

14.1 Logistic Regression

$$L(\mathbf{x}, \mathbf{w}) = - \sum_k [t_k^x \log(o_k^x) + (1 - t_k^x) \log(1 - o_k^x)]$$

$$o_k^x = \sigma(\mathbf{w}\mathbf{x}_k) = \frac{1}{1 + e^{-\mathbf{w}\mathbf{x}_k}}$$

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \sigma} \frac{\partial \sigma}{\partial \mathbf{w}} \\ &= (\mathbf{o} - \mathbf{t})\mathbf{x} \end{aligned}$$

Update rule:

$$\mathbf{w} = \mathbf{w} - \eta(\mathbf{o} - \mathbf{t})\mathbf{x}$$

14.1.1 Gradient Descent - General Approach

- init \mathbf{w} , chose learning rate η
- iterate:
 - $\text{grad} = \mathbf{0}$
 - iterate for every instance of the training data $x \in X$:
 - * compute output $o_x = \sigma(\mathbf{w}\mathbf{x})$
 - * compute error $E = t_x - o_x$
 - * $\text{grad} = \text{grad} + E\mathbf{x}$
 - update $\mathbf{w} = \mathbf{w} + \eta \text{grad}$
- until a convergence condition is reached (small error or all epochs done)

14.1.2 Batch Gradient Descent

Works like the normal gradient descent but the weights are updated every **epoch**!

- - iterate whole dataset to perform one update
- - slow and memory-intensive
- - can not be used for *online* training

14.1.3 Stochastic Gradient Descent

Update the weights right after we have seen *one training instance*. Before each epoch *shuffle* the training set!

- + converges much faster
- + no huge requirement of memory
- + can be used for online training
- - approximation of the gradient
- - high variance in updating

14.1.4 Mini-Batch Gradient Descent

Compromise between Batch and Stochastic Gradient descent. Perform one weight update after k training instances (called 1 **iteration**).

- + faster than batch gradient descent
- + overcome memory intensity (depends on k)
- + can be used for online training
- + Faster and more stable than Stochastic Gradient Descent
- + fewer smaller weight updates
- very easy to parallelize

14.2 Learning Rate Scheduling

Decay learning rates, making it smaller as the training is going. From time step t_d (the t_d^{th} iteration), calculate: $\eta_{t+1} = \beta \eta_t$. Choosing learning rates and learning rate scheduling can be tricky.

14.2.1 Adagrad - ADaptive GRADient method

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \nabla_{\mathbf{w}} L(\mathbf{x}, \mathbf{w}_{t-1})$$

$G_t \in \mathbb{R}^{d \times d}$: diagonal matrix, each diagonal entry is the sum of the squares of the gradients with respect to w_i up to time step t

ϵ : small smoothing term that avoids division by zero

- + good for sparse data
- + no manual tuning of the learning rate (normally $\eta = 0.01$, $\epsilon = 1e^{-8}$)
- - have to store G_t
- - G_t is accumulated \Rightarrow adaptive learning rate is smaller over time

14.2.2 Adadelta

Same as Adagrad but only stores a limited history of the gradients (normally 2).

$$\begin{aligned}\mathbf{w}_t &= \mathbf{w}_{t-1} - \frac{RMS[\Delta w]_{t-1}}{RMS[g]_t} g_t \\ g_t &= \nabla_{\mathbf{w}} L(\mathbf{x}, \mathbf{w}_{t-1}) \\ RMS[\Delta w]_t &= \sqrt{E[\Delta w^2]_t + \epsilon} \\ E[\Delta w^2]_t &= \gamma E[\Delta w^2]_{t-1} + (1 - \gamma) \Delta w_t^2\end{aligned}$$

- das $E[\Delta w^2]_t$???

- + not memory intensiv
- + no picking of the learning rate needed

14.2.3 RMSprop

$$\begin{aligned}\mathbf{w}_t &= \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \\ E[g^2]_t &= \gamma E[g^2]_{t-1} + g_t^2\end{aligned}$$

14.2.4 Adam - ADaptive Moment estimation

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\widehat{\mathbf{v}}_t + \epsilon}} \widehat{\mathbf{m}}_t$$

$$\widehat{\mathbf{m}}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\widehat{\mathbf{v}}_t = \frac{v_t}{1 - \beta_2^t}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

-beta1, beta2??

- + faster convergence
- - practically more overfitting
- - needs resetting
- - needs to store to matrices of past gradients

15 Error Functions

15.1 Binary Cross Entropy or Negative Log Likelihood

$$L(\mathbf{x}, \mathbf{w}) = - \sum_k [t_k^x \log(o_k^x) + (1 - t_k^x) \log(1 - o_k^x)]$$

t_k^x : target label of instance \mathbf{x}

o_k^x : real output

$$o_k^x = P(t_k^x = 1 | x_k; w) = g(\mathbf{w} \cdot \mathbf{x}_k) \quad = P(t_k^x = 0 | x_k; w) = 1 - g(\mathbf{w} \cdot \mathbf{x}_k)$$

$g(z)$: sigmoid functions (or Logistic functions)

Mean Square Error Cross Entropy Classification Figure of Merit

16 Activation Functions

The activation functions should have the following properties:

- continuous
- bounded
- monotonically increasing
- differentiable

16.1 Step Function

$$\varphi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

Derivative is always 0

16.2 Linear Function

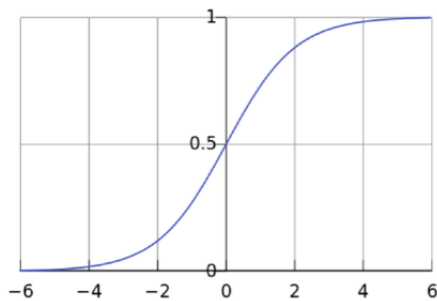
$$\varphi(x) = x$$

Linear functions alone can only solve linear separable problems but can be used in a combination node for function approximation problems.

16.3 Logistic / Sigmoid Function

$$\varphi(x) = \text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$
$$\frac{\delta\varphi(x)}{\varphi(x)} = \varphi(x)(1 - \varphi(x))$$

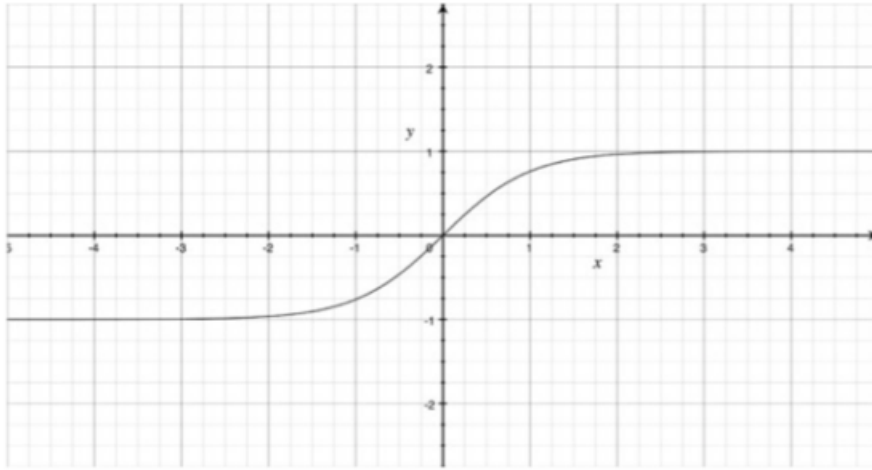
Good for internal nodes, bad for output nodes.



16.4 Hyperbolic Tangent function

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

If the input has a mean of 0 then so will the output



16.5 Softmax Function

$$\varphi(x_j) = \frac{e^{x_j}}{\sum_k e^{x_k}}$$
$$\frac{\delta\varphi(x_j)}{\varphi(x_j)} = \varphi(x_j) - \varphi(x_j)^2 = \varphi(x_j)(1 - \varphi(x_j))$$

Outputs an a posteriori probability $p(c|x)$ and is good for classification tasks.

16.6 Rectified Linear Unit

s