

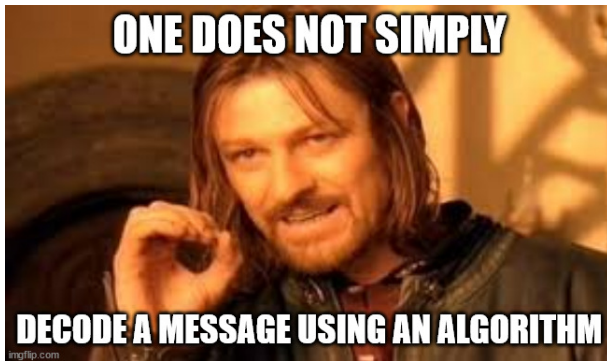
On the Correspondence Between Classic Coding Theory and Machine Learning

Julia Rechberger
Bachelor Thesis Presentation

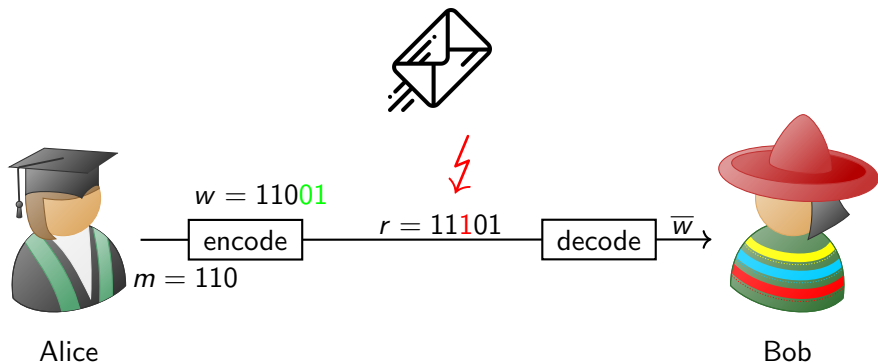
August 18, 2021

Key Question

How can a neural network advance decoding?



Coding Theory in a Nutshell



Message m is encoded to codeword w , w is corrupted by the channel and arrives as r , decoding produces approximation \bar{w} .

Channels: Error Pattern Modelling

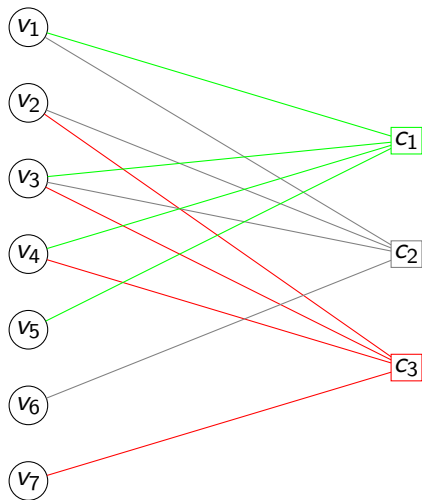
Additive White Gaussian Noise (AWGN)

- ① apply binary phase-shift keying (BPSK): 0-bit \rightarrow 1, 1-bit \rightarrow -1
- ② add Gaussian distributed independent noise
- ③ recover bit from a non-binary value with a hard decision HD

$$HD(\text{value}) = \begin{cases} 0, & \text{if value} > 0 \\ 1, & \text{if value} < 0. \end{cases}$$

Message Passing Decoding on a Tanner Graph

$$H = \begin{pmatrix} \textcolor{green}{1} & 0 & \textcolor{green}{1} & \textcolor{green}{1} & \textcolor{green}{1} & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & \textcolor{red}{1} & \textcolor{red}{1} & \textcolor{red}{1} & 0 & 0 & \textcolor{red}{1} \end{pmatrix} \Rightarrow$$



Probabilistic Decoding - Sum-Product Algorithm (SPA)

Idea: use probability for a certain bit value instead the bit value itself



Why use a neural network?

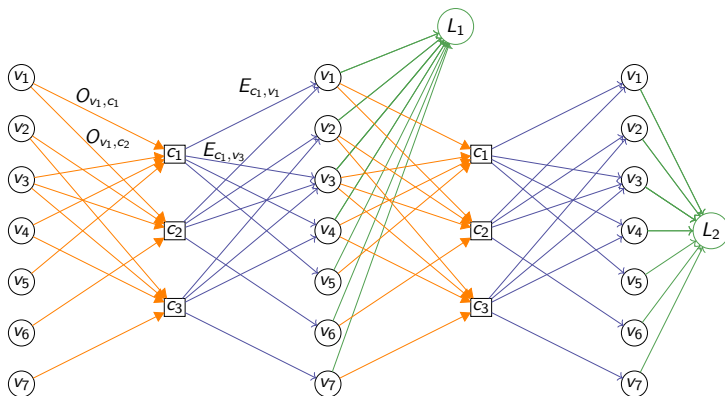
SPA...

- is correct if graph is a tree, otherwise approximative
- can't adapt to specific channel properties

Begin the transformation!

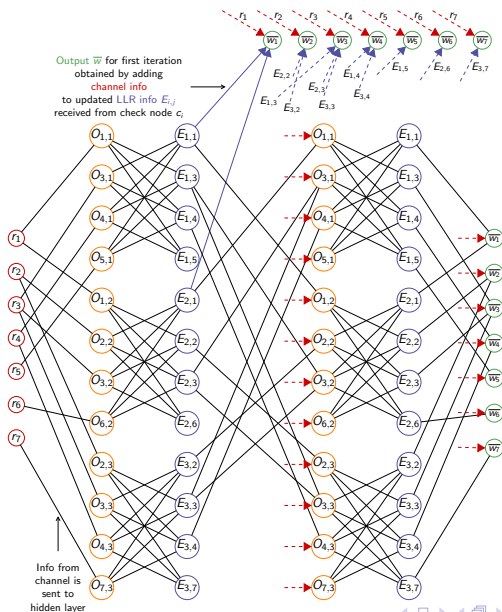
- 1 Unfold Tanner graph
- 2 Derive weight matrices representing edges from unfolded graph
- 3 Transform message calculation into activation functions

Unfolded Tanner graph for 2 SPA iterations



- O_{v_j, c_i} message sent from variable node to check node
- E_{c_i, v_j} message sent from check node to variable node

Neural network for 2 SPA iterations



Activation Functions - Odd Layer

$$O_{v_j, c_i} = \tanh \left(\frac{1}{2} \cdot \left(r_j + \sum_{c_t \in \text{adj}(v_j) \setminus c_i} E_{c_t, v_j} \right) \right)$$

Calculation for odd layer output vector \vec{X}_i of size $|E|$ for iteration T_i :

① $\vec{X}_0 = \vec{0}$ no initial information from check nodes

②

$$\vec{X}_i = \tanh \left(\frac{1}{2} \cdot \left(W_{in2odd}^T \cdot \vec{R} + W_{even2odd}^T \cdot \vec{X}_{i-1} \right) \right)$$

\vec{X}_{i-1} : output vector of previous hidden even layer and

\vec{R} : vector of size n with input LLR channel information for each bit.

Cross Entropy Loss Function (*CEL*)

Definition

$$CEL(\mathbb{P}(x = 0), y) = -\frac{1}{N} \sum_{n=1}^N y_n \cdot \ln(1 - \mathbb{P}(x_n = 0)) + (1 - y_n) \cdot \ln(\mathbb{P}(x_n = 0))$$

N : length of both codewords

y : binary vector of target codeword

y_n : n th bit of the target codeword

$\mathbb{P}(x_n = 0)$: prediction for n -th bit of input vector being a zero

Metrics

- Signal-to-Noise ratio in dB (SNR_{dB})

$$\text{SNR}_{\text{dB}} = 10 \cdot \log_{10} \left(\frac{P_{\text{signal}}}{P_{\text{noise}}} \right) \text{ dB}$$

- Bit-Error Rate (BER) divide incorrect bits b_e by total bits B :

$$\text{BER} \approx \frac{b_e}{B}$$

- Normalized-Validation Score (NVS) to compare SPA and Neural Network Decoder (NND)

$$\text{NVS} = \frac{1}{|S|} \cdot \sum_{s \in S} \frac{\text{BER}_{\text{NND}}(D(s))}{\text{BER}_{\text{SPA}}(D(s))}$$

S = a set of different SNR_{dB}

$D(s)$ = a dataset created by using $s \in S$.

Training Strategy and Hyperparameters

Fixed

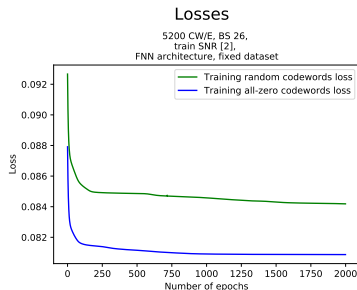
- 5 Iterations
- Supervised Learning
- Activation and Loss Functions
- Optimizer: RMSProp (Learning Rate 0.001, α 0.99)

Part of Experiments

- Training Dataset: Size, Codewords (CW), SNR_{dB}
- Batchsize
- Architecture: Feed Forward Neural Network, Recurrent Neural Network

All-Zero Codeword

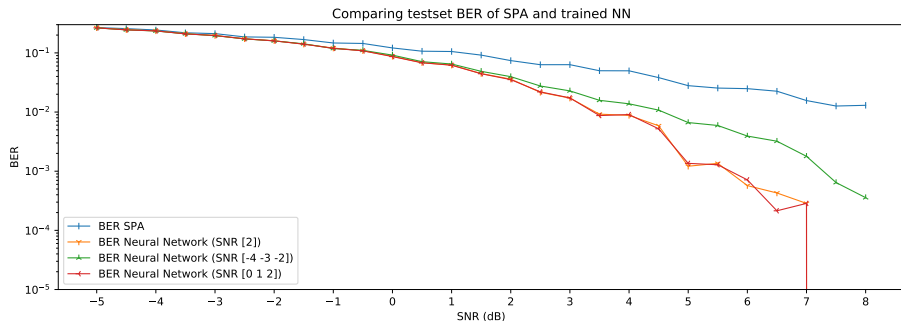
It does not matter if the all-zero codeword only or randomly picked codewords are used in the training data set.



- SPA independent of transmitted CW, NND keeps information flow
- slight difference: random NVS_{test} 0.487, zero NVS_{test} 0.486
- Culprit: Loss Function

SNR_{dB} Combinations

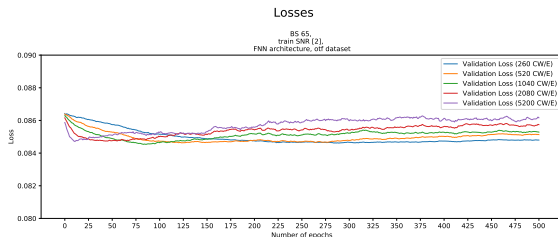
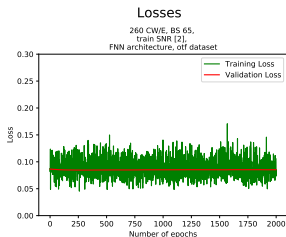
Training on only one SNR_{dB} leads to results equivalent to using a range of SNR_{dB} in the training data set.



- single SNR_{dB} (NVS_{test} 0.49) close to SNR_{dB} combo [0, 1, 2] with NVS_{test} 0.487
- optimal single SNR_{dB} is 2 or 3

On the fly (OTF) generated Training Dataset

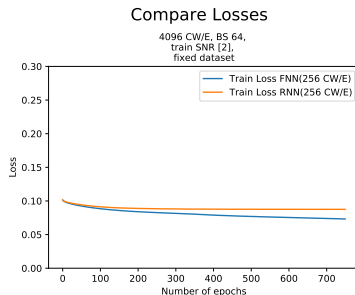
An on the fly created training dataset using one SNR_{dB} trains the network better than a fixed dataset.



- use validation set to cope with fluctuations
- 260 CW biggest improvement: fixed NVS_{test} 0.566, OTF NVS_{test} 0.49
- Neural network converges after few epochs

Feed-Forward (FNN) vs. Recurrent Neural Network (RNN)

FNN architecture leads to a better result than RNN architecture.



- RNN converges faster
- FNN more precise (FNN NVS_{test} 0.492, RNN NVS_{test} 0.498)

Why use a neural network for decoding?

- BER performance increase for short codes
- same computational complexity as Sum-Product algorithm (SPA)
- easy training data creation using the all-zero codeword
- learn channel and linear code at the same time

Thank you!



Main Sources

- AGRAWAL, Navneet: Machine Intelligence in Decoding of Forward Error Correction Codes. Stockholm, Sweden, KTH Royal Institute of Technology School of Electrical Engineering, Diss., 2017.
- JOHNSON, Sarah J.: Introducing Low-Density Parity-CheckCodes. In: School of Electrical Engineering and Computer Science
- MACKAY, David J.: Information Theory, Inference and Learning Algorithms. Cambridge University Press, 2003
- NACHMANI, Eliya ; MARCIANO, Elad ; LUGOSCH, Loren ; GROSS, Warren J. ; BURSHTAIN, David ; BE'ERY, Yair: Deep Learning Methods for Improved Decoding of Linear Codes. In: CoRR abs/1706.07043 (2017). <http://arxiv.org/abs/1706.07043>

Construct codes that correct maximal errors while using minimal redundancy with efficient encoding and decoding procedures.

- Encode message m of length k with canonical generator matrix $G_C = [I_k, A]$ of code C . Producing unique codewords w of length n .

$$m \cdot G_C = w \in C$$

- Check for errors in the received word r with canonical parity check matrix $H_C = [A, I_{n-k}]$.

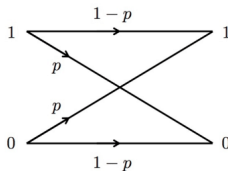
$$r \in C \Leftrightarrow H_C \cdot r = \vec{0}.$$

- Transform matrices by using the relationship

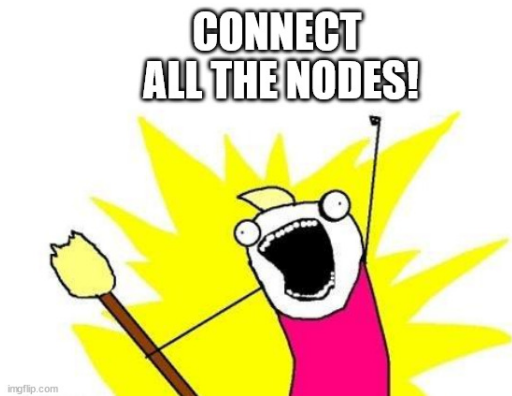
$$H_C = [A, I_{n-k}] \Leftrightarrow G_C = [I_k, A^T]$$

Binary Symmetric Channel (BSC)

Bit flip with probability p



Naive Approach: A Fully Connected Neural Network



Bit-Flipping-Algorithm

Idea: Variable node v_j bit is determined by a **Majority Vote** on values received from adjacent check nodes c_i .



Bit-Flipping-Algorithm

- ① Iteration counter $c_T = 0$, initialize graph with received word r
- ② Evaluate parity check equations, if ...
 - ① $\forall c_i = 0$, **return:** \bar{w}
 - ② $c_T = T_{\max}$, **return:** error.
 - ③ $\exists c_i \neq 0$: continue.
- ③ Calculate messages E_{c_i, v_j} , send from c_i to v_j

$$E_{c_i, v_j} = \sum_{n \in \text{adj}(c_i) \setminus v_j} v_n \mod 2$$

- ④ Majority vote for variable nodes v_j :

$$v_j = \begin{cases} 1, & \text{if } (\# E_{\text{adj}(v_j)} = 1) > (\# E_{\text{adj}(v_j)} = 0) \\ 0, & \text{if } (\# E_{\text{adj}(v_j)} = 1) < (\# E_{\text{adj}(v_j)} = 0) \\ v_j, & \text{if } (\# E_{\text{adj}(v_j)} = 1) = (\# E_{\text{adj}(v_j)} = 0). \end{cases}$$

- ⑤ $c_T = c_T + 1$, jump to 2.

SPA I - Initialization

- 1 Iteration counter $c_T = 0$
- 2 Determine a priori probability for every bit $r' = (r'_1, \dots, r'_n)$ to initialize Tanner graph by calculating LLR for received bits $r = (r_1, \dots, r_n)$.

$$LLR(bit(v_j)|r_j) = \ln \left(\frac{\mathbb{P}(bit(v_j) = 0|r_j)}{\mathbb{P}(bit(v_j) = 1|r_j)} \right)$$

$$priori(v_j) = LLR(bit(v_j)|r_j) = r'_j$$

Retrieve bit value from LLR by making a hard decision HD

$$HD(r'_j) = \begin{cases} 0, & \text{if } r'_j > 0 \\ 1, & \text{if } r'_j < 0. \end{cases}$$

$$bit(v_j) = HD(r'_j)$$

Initialize $value(v_j) = 0$ which holds LLR of current iteration.

SPA II - Iteration

- ③ Execute parity check, if ...
 - ① $c_T = T_{\max}$, **return:** error and approximation $\bar{w} = (\text{bit}(v_1), \dots, \text{bit}(v_n))$
 - ② $H \cdot (\text{bit}(v_1), \dots, \text{bit}(v_n)) = \vec{0}$, **return:** $\bar{w} = (\text{bit}(v_1), \dots, \text{bit}(v_n))$
 - ③ $H \cdot (\text{bit}(v_1), \dots, \text{bit}(v_n)) \neq \vec{0}$: continue.
- ④ Calculate E_{c_i, v_j} for all nodes and send them

$$E_{c_i, v_j} = \ln \left(\frac{\frac{1}{2} + \frac{1}{2} \cdot \prod_{v_t \in \text{adj}(c_i) \setminus v_j} \left(1 - 2 \cdot \frac{e^{-\text{value}(v_t)}}{1 + e^{-\text{value}(v_t)}} \right)}{\frac{1}{2} - \frac{1}{2} \cdot \prod_{v_t \in \text{adj}(c_i) \setminus v_j} \left(1 - 2 \cdot \frac{e^{-\text{value}(v_t)}}{1 + e^{-\text{value}(v_t)}} \right)} \right)$$

- ⑤ Update $\text{value}(v_j)$ with received messages E_{c_i, v_j}

$$\text{value}(v_j) = \text{priori}(v_j) + \sum_{c_i \in \text{Adj}(v_j)} E_{c_i, v_j}$$

- ⑥ Update $\text{bit}(v_j) = \text{HD}(\text{value}(v_j))$
- ⑦ Calculate O_{v_j, c_i} , send them, $c_T = c_T + 1$, and jump to ③

$$O_{v_j, c_i} = r'_j + \sum_{c_t \in \text{adj}(v_j) \setminus c_i} E_{c_t, v_j}$$

Weight matrices representing unfolded graphs edges

- W_{in2odd} for **input** layer to **odd** hidden layer
- $W_{odd2even}$ for **odd** hidden layer to **even** hidden layer
- $W_{even2odd}$ for **even** hidden layer to **odd** hidden layer
- $W_{even2out}$ for **even** hidden layer to **output** layer

e.g. $W_{odd2even}$: edges along which messages from variable nodes to check node are sent

Definition

$$W_{odd2even}(O_{(v_j, c_i)}, E_{(c'_i, v'_j)}) = \begin{cases} 1, & c_i = c'_i \text{ and } v_j \neq v'_j \\ 0, & \text{otherwise} \end{cases}$$

with rows $O_{(v_j, c_i)}$ where $(v_j, c_i) \in E$, E being the edges of the Tanner graph and columns $E_{(c'_i, v'_j)}$ where $(c'_i, v'_j) \in E$.

	E_{c_1, v_1}	E_{c_1, v_3}	E_{c_1, v_4}	E_{c_1, v_5}	E_{c_2, v_1}	E_{c_2, v_2}	E_{c_2, v_3}	E_{c_2, v_6}	E_{c_3, v_2}	E_{c_3, v_3}	E_{c_3, v_4}	E_{c_3, v_7}
O_{v_1, c_1}	0	1	1	1	0	0	0	0	0	0	0	0
O_{v_3, c_1}	1	0	1	1	0	0	0	0	0	0	0	0
O_{v_4, c_1}	1	1	0	1	0	0	0	0	0	0	0	0
O_{v_5, c_1}	1	1	1	0	0	0	0	0	0	0	0	0
O_{v_1, c_2}	0	0	0	0	0	1	1	1	0	0	0	0
O_{v_2, c_2}	0	0	0	0	1	0	1	1	0	0	0	0
O_{v_3, c_2}	0	0	0	0	1	1	0	1	0	0	0	0
O_{v_6, c_2}	0	0	0	0	1	1	1	0	0	0	0	0
O_{v_2, c_3}	0	0	0	0	0	0	0	0	0	1	1	1
O_{v_3, c_3}	0	0	0	0	0	0	0	0	1	0	1	1
O_{v_4, c_3}	0	0	0	0	0	0	0	0	1	1	0	1
O_{v_7, c_3}	0	0	0	0	0	0	0	0	1	1	1	0

Figure: If c_i is equal to c'_i it will be highlighted blue but only if v_j is **not** equal to v'_j a 1 is entered.

Activation Functions - Even Layer

Use relationship $2 \cdot \tanh^{-1}(p) = \ln\left(\frac{1+p}{1-p}\right)$ to transform E_{c_i, v_j}

$$E_{c_i, v_j} = 2 \cdot \tanh^{-1} \left(\prod_{v_t \in \text{adj}(c_i) \setminus v_j} (1 - 2 \cdot \mathbb{P}^{(\text{in})}(v_t = 1)) \right)$$

Calculation for even layer output vector \vec{X}_i of size $|E|$ for iteration T_i :

- 1 Repeat T_{i-1} odd layer output as a column $|E|$ times to get M_{i-1}
- 2 $M_{i-1}^* = W_{\text{odd2even}} \odot M_{i-1}$, where \odot is element-wise multiplication
- 3 Replace zeros in M_{i-1}^* with ones
- 4 Calculate \vec{X}_{i-1}^* by multiplying along the column elements of M_{i-1}^* .

5

$$\vec{X}_i = 2 \cdot \tanh^{-1} \left(\vec{X}_{i-1}^* \right)$$

Activation Functions - Output Layer

$$\overline{w}_j = r_j + \sum_{c_t \in \text{adj}(v_j)} E_{c_t, v_j}$$

Output layer $\overrightarrow{\overline{W}}_i$ of size n :

$$\overrightarrow{\overline{W}}_i = \sigma(\overrightarrow{R} + W_{\text{even2out}}^T \cdot \overrightarrow{X}_{i-1})$$

$\sigma(\cdot)$: sigmoid function

\overrightarrow{R} : vector of size n with input LLR channel information for each bit

\overrightarrow{X}_{i-1} : output of previous even hidden layer of size $|E|$.

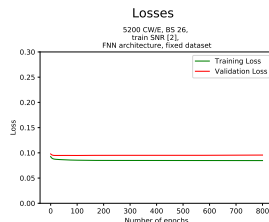
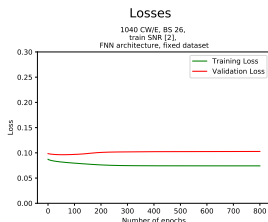
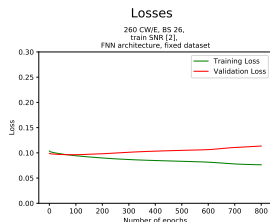
Datasets

- 1 Generate random message m
- 2 Encode m : $w = G \cdot m$
- 3 Map w 's bits: $0 \rightarrow 1$ and $1 \rightarrow -1$
- 4 Add Gaussian noise w to generate received word r , variance depending on SNR_{dB}

Dataset	Codewords	SNR_{dB}	CW Seed	Noise Seed
Training	Experiment dependent	Experiment dependent	0	11
Validation	500	same as experiment	4	5
Test S5	54000	$[-5.0, -4.5, \dots, 8]$	5	5
Test S4	54000	$[-5.0, -4.5, \dots, 8]$	4	4

Training Dataset Size

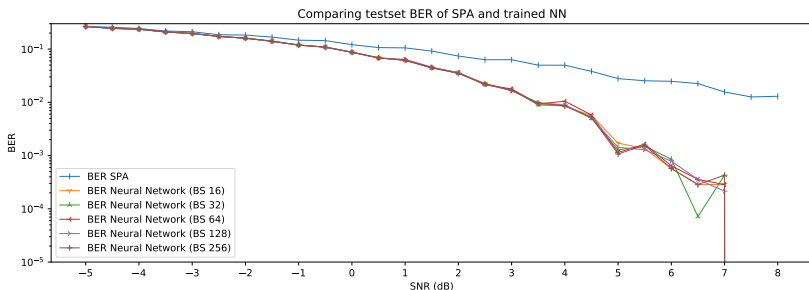
The bigger the fixed training dataset size the better the performance of the trained neural network.



- 260 CW: overfitting (NVS_{test} 0.566)
- 5200 CW for 48 learnable parameters too big (NVS_{test} 0.487)

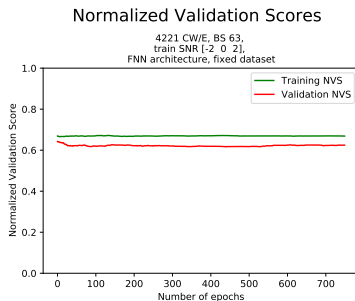
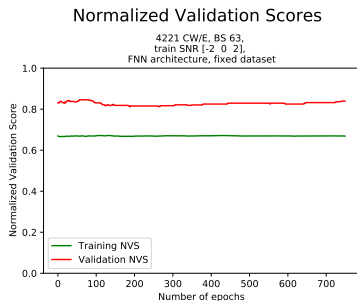
Batch Size

Code size and batch size share a linear dependency.



Batchsize	NVS testset S5	NVS testset S4
16	0.489	0.492
32	0.488	0.49
64	0.492	0.489
128	0.49	0.495
256	0.489	0.487

SNR_{dB} Combinations



- Left: Validation dataset with 140 CW per SNR_{dB} not representative
- Right: Validation dataset with 500 CW per SNR_{dB}

Implementation Lessons Learned

- don't use numpy in custom activation functions
- debug backwards pass using `torch.autograd.detect_anomaly`
- test edge cases
- use already existing modules and libraries

