

Deep Learning

Chapter 11 (+ 12...)
“Tricks of the trade”

Helge Langseth
helgel@ntnu.no

1

RefGrp

- We need at least 2 people...

3

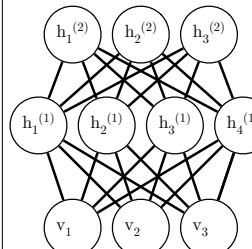
Disclaimer

- The content in this lecture is **mostly** based on the book, but also contains some “gut-feeling”, “rules of thumb”, and “urban legends”.
- Some of what I say may therefore be “controversial”:
 - You are allowed to disagree!
 - If you disagree, you may well be right!
- But then again, some slides are from the book, too, and those are obviously correct and just perfect!

2

What drives success in ML?

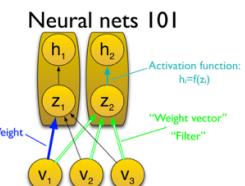
Arcane knowledge
of dozens of
obscure algorithms?



Mountains
of data?



Knowing how
to apply 3-4
standard techniques?



4

Three Step Process

- Use needs to define metric-based goals
- Build an end-to-end system
- Data-driven refinement

5

Choose Metrics... ... or combinations thereof

- Accuracy? (% of examples correct)
- Coverage? (% of examples processed)
- Precision? (% of detections that are right)
- Recall? (% of objects detected)
- Amount of error? (For regression problems)

7

Step 1: Identify needs

6

Metrics

- The system will try to “cheat” during learning:
 - If there is a short-cut to optimize the loss, the optimizer may well find it
 - Make sure a model with good results on your metric is a good model in production.
 - Ex.: Detection of a rare disease

8

What is required for this to be a success???

- High accuracy or low accuracy?
- Surgery robot: high accuracy
- Celebrity look-a-like app: low accuracy

9

Step 2: Build end-to-end system

10

End-to-end System

- Get up and running ASAP
- Build the simplest viable system first
- What baseline to start with though?
 - Copy state-of-the-art from related publication

11

Deep or Not?

- Lots of noise, little structure -> not deep
- Little noise, complex structure -> deep
- Good shallow baseline:
 - *Use what you know*
- Logistic regression, SVM, boosted tree are all good

12

Entangled models

- If your model contains several “sub-models” (which is not always recommendable...):
 - Evaluate them separately if possible
 - Where does the system error come from?
 - Prioritize improvements based on cost/benefit
- If your model does **not** consist of sub-models, it can **still** be beneficial to pipe out intermediate results and do “QA” if at all reasonable.

13

Step 3: Data driven refinement

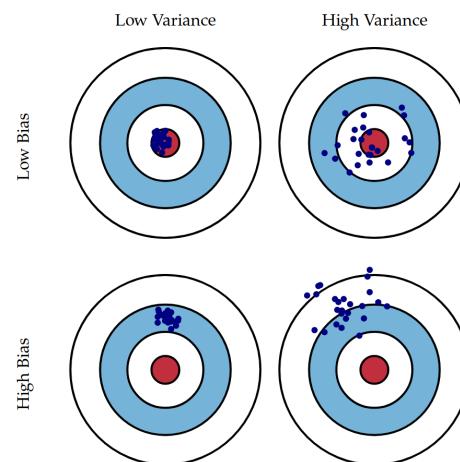
14

Data Driven Refinement

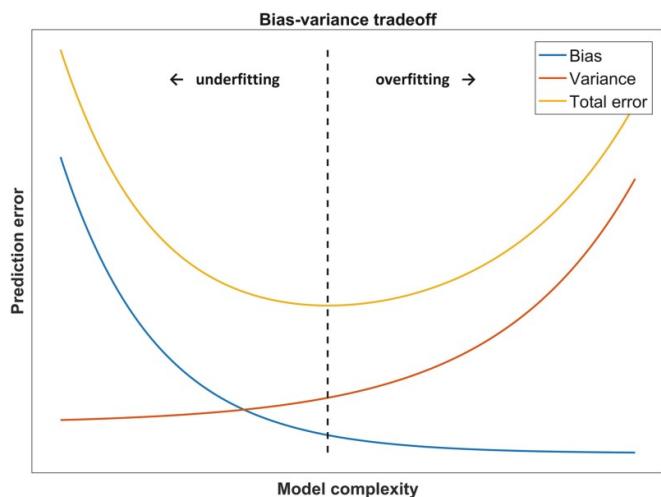
- Understand your model and its strengths/weaknesses
 - Measure error on train and test sets
 - Overfitting vs underfitting
- Understand the learning process
 - Use Tensorboard and/or other visualizations
 - Know what to look for, and how to fix problems
- Choose a models with relevant learning bias
 - Don't believe the hype
 - Feed Forward, CNN, RNN, more advanced tricks

15

Variance - Bias tradeoff



16



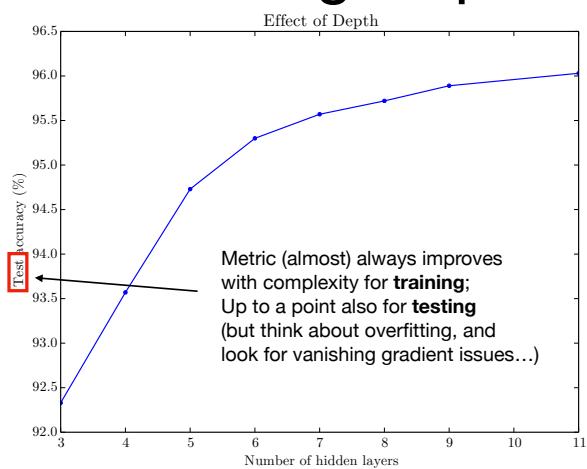
17



- This can be the high-bias case —> Under-fitting
 - Increase representational power (model depth, layer size, remove regularizers, ...)
- ... or a fundamental error somewhere:
 - Code issues (stick to what you know works — like TF...)
 - Data issues
 - Tuning of hyper-parameters

18

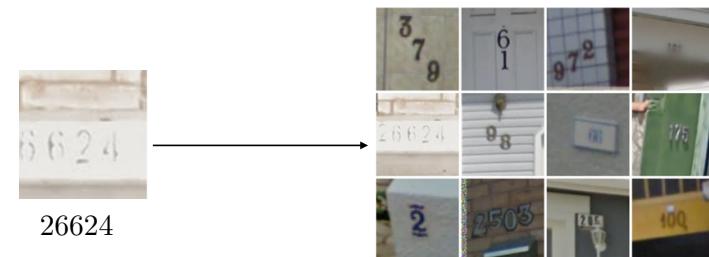
Increasing Depth



19

Checking Data for Defects

Can a human process it?



Just as important, but more tedious to check:
Are there systematic errors? Are there trends in mistakes (e.g., always on brick wall or similar).

20

Hyper-parameters

Alternative 1: Think long and hard about the HPs, and their effect...

Hyperparameter	Increases capacity when...	Reason	Caveats
Number of hidden units	increased	Increasing the number of hidden units increases the representational capacity of the model.	Increasing the number of hidden units increases both the time and memory cost of essentially every operation on the model.

Table 11.1

21

Hyper-parameters

Alternative 2: Automated search

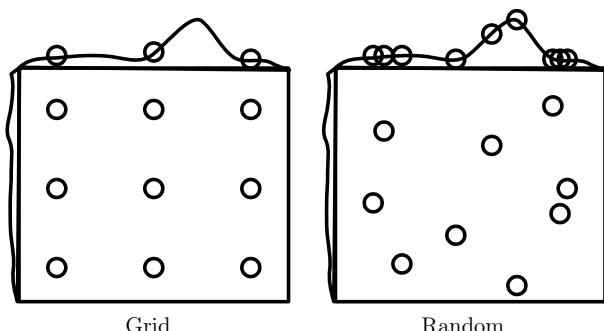


Figure 11.2

23

Tuning the Learning Rate

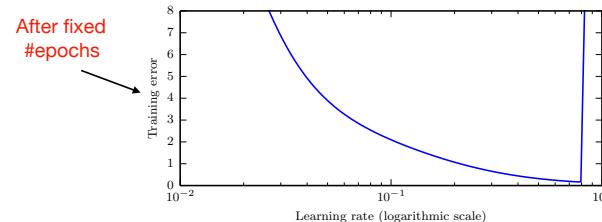


Figure 11.1

- **Too high lr:** Can go astray
- **Too low lr:**
 - Computationally costly (more epochs to get a decent move)
 - Potentially stuck in local minima of the loss landscape
- **Just right lr:** Works just right!!
- Either work very hard do tune the lr or use an lr scheduler (extremely simple with Keras: `tf.keras.callbacks.LearningRateScheduler`) ... or both!!

22

Practical Bayesian Optimization of Machine Learning Algorithms

Jasper Snoek
Department of Computer Science
University of Toronto
jasper@cs.toronto.edu

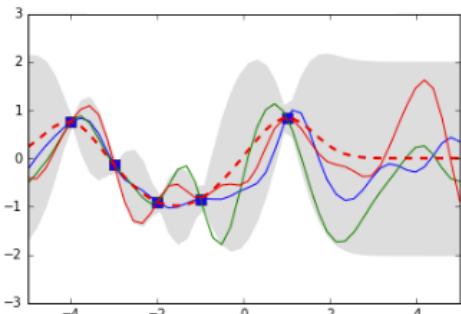
Hugo Larochelle
Department of Computer Science
University of Sherbrooke
hugo.larochelle@usherbrooke.edu

Ryan P. Adams
School of Engineering and Applied Sciences
Harvard University
rpa@seas.harvard.edu

24

Gaussian Processes

- A fancy way to make statistical distributions over **curves**.
- Assumes all possible points on the curve are jointly Gaussian, with a covariance between pairs depending on their distance in "input-space".
- We can sample by (infinitely many times) generating "y-value" for some "x-value", conditioned on previous samples



25

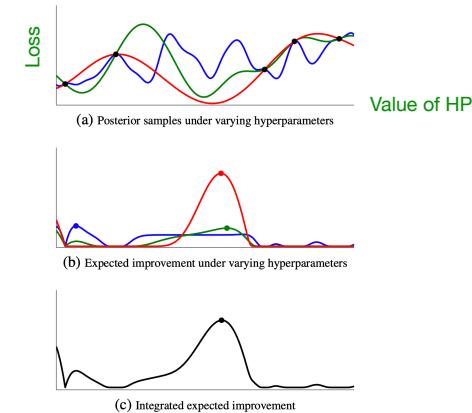


Figure 1: Illustration of integrated expected improvement. (a) Three posterior samples are shown, each with different length scales, after the same five observations. (b) Three expected improvement acquisition functions, with the same data and hyperparameters. The maximum of each is shown. (c) The integrated expected improvement, with its maximum shown.

26-1

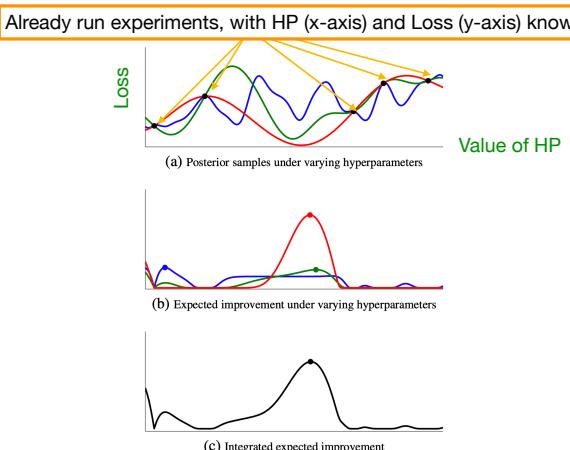


Figure 1: Illustration of integrated expected improvement. (a) Three posterior samples are shown, each with different length scales, after the same five observations. (b) Three expected improvement acquisition functions, with the same data and hyperparameters. The maximum of each is shown. (c) The integrated expected improvement, with its maximum shown.

26-2

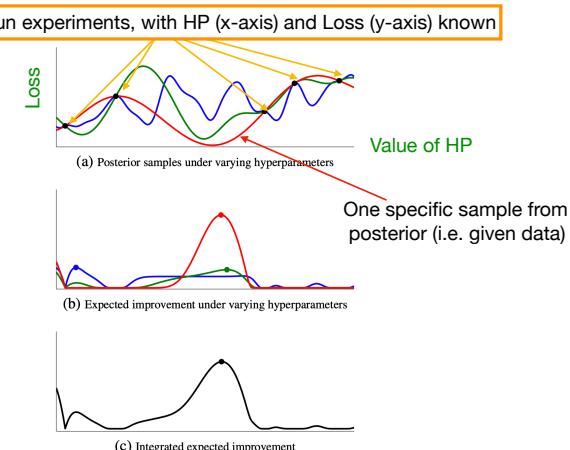


Figure 1: Illustration of integrated expected improvement. (a) Three posterior samples are shown, each with different length scales, after the same five observations. (b) Three expected improvement acquisition functions, with the same data and hyperparameters. The maximum of each is shown. (c) The integrated expected improvement, with its maximum shown.

26-3

Already run experiments, with HP (x-axis) and Loss (y-axis) known

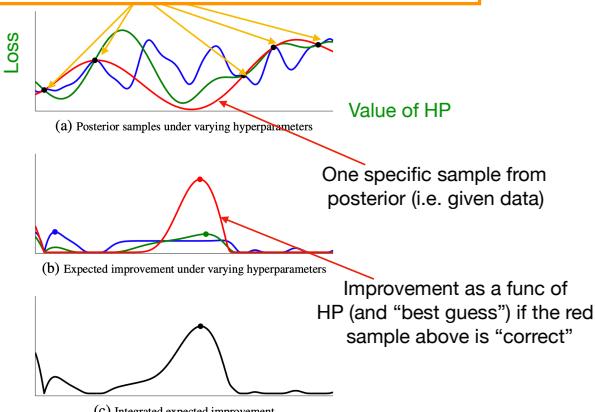


Figure 1: Illustration of integrated expected improvement. (a) Three posterior samples are shown, each with different length scales, after the same five observations. (b) Three expected improvement acquisition functions, with the same data and hyperparameters. The maximum of each is shown. (c) The integrated expected improvement, with its maximum shown.

26-4

Implementation in Keras tuner

```
model.add(
    Dropout(
        rate=hp.Float('dropout_1', min_value=0.0, max_value=0.5, default=0.25, step=0.05)
    )
)
model.add(
    Conv2D(
        filters=hp.Choice('num_filters', values=[32, 64], default=64),
        activation='relu',
        kernel_size=3
    )
)
bayesian_tuner = kerastuner.tuners.BayesianOptimization(
    hypermodel,
    objective='val_accuracy',
    seed=SEED,
    max_trials=MAX_TRIALS,
    executions_per_trial=EXECUTION_PER_TRIAL,
    directory=f'{directory}/bayesian',
    project_name=project_name
)
tuner.search(x_train, y_train,
             epochs=N_EPOCH_SEARCH,
             validation_split=0.1,
             verbose=0)

# Show a summary of the search
tuner.results_summary()

# Retrieve the best model.
best_model = tuner.get_best_models(num_models=1)[0]

# Evaluate the best model.
loss, accuracy = best_model.evaluate(x_test, y_test)
```

Show code?

Already run experiments, with HP (x-axis) and Loss (y-axis) known

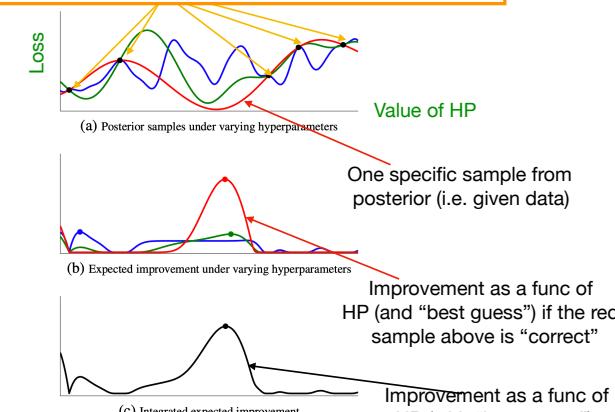
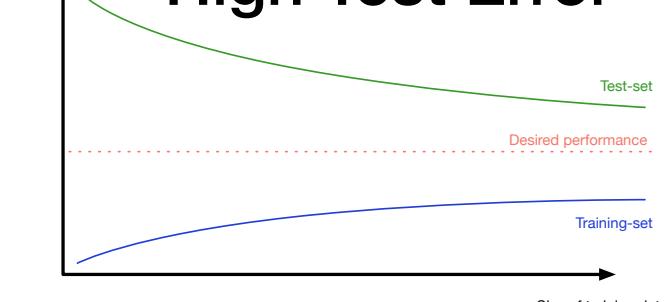


Figure 1: Illustration of integrated expected improvement. (a) Three posterior samples are shown, each with different length scales, after the same five observations. (b) Three expected improvement acquisition functions, with the same data and hyperparameters. The maximum of each is shown. (c) The integrated expected improvement, with its maximum shown.

26-5

High Test Error



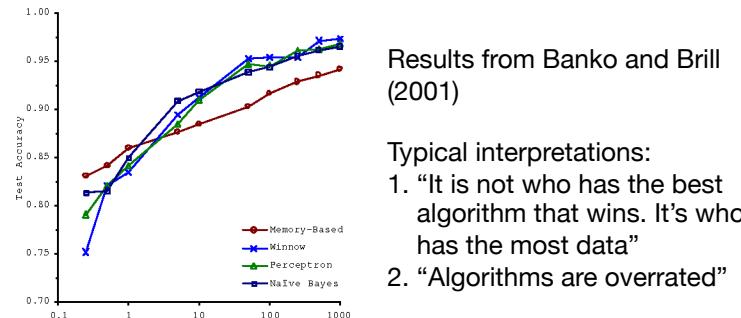
- This is the high-variance case —> Overfitting
- Regularization is good here (e.g., Dropout, L1, L2, ...)
- Standard trick: **MORE DATA** (or data augmentations)

27

28

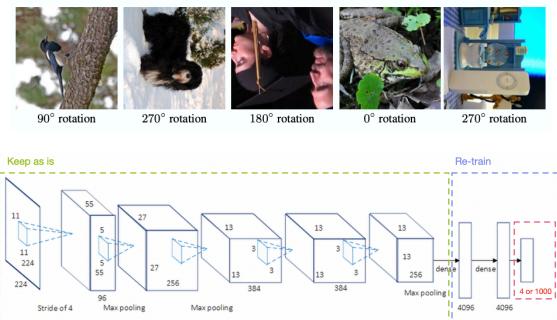
Data, data, data ...

- Big dataset cannot be bad (if you have it, and time to use it):
 - Will be able to learn whatever is to be learned
 - Also: Less danger of overfitting



29

Learning rotations



31

Self-supervision

- **Self-supervision:** Supervised learning where the training data is automatically labelled.
- Typically done with the expectation that the internal representations learned will be beneficial for a different (sparsely labelled) problem. —> Transfer-learning into the original problem domain.
- **Standard example:** Filling in missing image patches.

UNSUPERVISED REPRESENTATION LEARNING BY PREDICTING IMAGE ROTATIONS

Spyros Gidaris, Praveer Singh, Nikos Komodakis

30

Does it work?

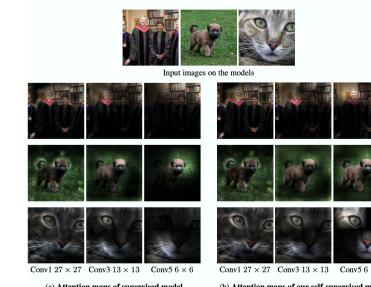


Figure 3: Attention maps generated by an AlexNet model trained (a) to recognize objects (supervised), and (b) to recognize image rotations (self-supervised). In order to generate the attention map of a conv. layer we first compute the feature maps of this layer, then we raise each feature activation on the layer to a power p , and then we average the results. For the conv. layers 1, 2, and 3 we used the powers $p = 1$, $p = 2$, and $p = 4$ respectively. For visualization of our self-supervised model’s attention maps for all the rotated versions of the images see Figure 6 in appendix A.

Table 3: Evaluation of unsupervised feature learning methods on CIFAR-10. The *Supervised NIN* and the *(Ours) RotNet + conv* entries have exactly the same architecture but the first was trained fully supervised while on the second the first 2 conv. blocks were trained unsupervised with our rotation prediction task and the 3rd block only was trained in a supervised manner. In the *Random Init. + conv* entry a conv. classifier (similar to that of *(Ours) RotNet + conv*) is trained on top of two NIN conv. blocks that are randomly initialized and stay frozen. Note that each of the prior approaches has a different ConvNet architecture and thus the comparison with them is just indicative.

Method	Accuracy
Supervised NIN	92.80
Random Init. + conv	72.50
<i>(Ours) RotNet + conv-linear</i>	89.96
<i>(Ours) RotNet + conv</i>	91.16
<i>(Ours) RotNet + non-linear (fine-tuned)</i>	91.73
<i>(Ours) RotNet + conv (fine-tuned)</i>	92.17
Roto-Sar + SVM Oyallon & Mallat (2015)	82.3
ExemplarCNN Dosovitskiy et al. (2014)	84.3
DCGAN Radford et al. (2015)	82.8
Scattering Oyallon et al. (2017)	84.7

32

Data Driven Refinement

- Understand your model
 - Measure error on train and test sets
 - Overfitting vs underfitting
- Understand the learning process
 - Use Tensorboard and/or other visualizations
 - Know what to look for, and how to fix problems
- Choose a models with relevant learning bias
 - Don't believe the hype
 - Feed Forward, CNN, RNN, more advanced tricks

33

Example: Suboptimal learning in (too) deep FFNNs

- I found that layer activations from ReLUs were had:
 - Exhibiting problems with vanishing gradient
 - High fraction dead nodes
 - Some nodes potentially “overwhelming” the others in next layer
 - Behaviour changing considerably between epochs
- BatchNorm only partly solved it — still dead nodes
- LeakyReLU can also partly solve it, but vanishing gradients remains

35

Visualization

- Plot all the usual stuff, like loss, accuracy, other metrics (train and validation sets); both actual value and deltas between epochs can be useful. Log-scale sometimes beneficial when plotting
- Per layer:
 - Scalars (Mean, SD) + histograms of: Weights, biases, activations, gradients
 - Track sparsity of: Outputs, Gradients (`tf.nn.zero_fraction`)
- Spit out values of all “moving parts” (e.g., learning-rate), somewhere so it is easy to align with other per-epoch result

34

Self-Normalizing Neural Networks

Günter Klambauer

Thomas Unterthiner

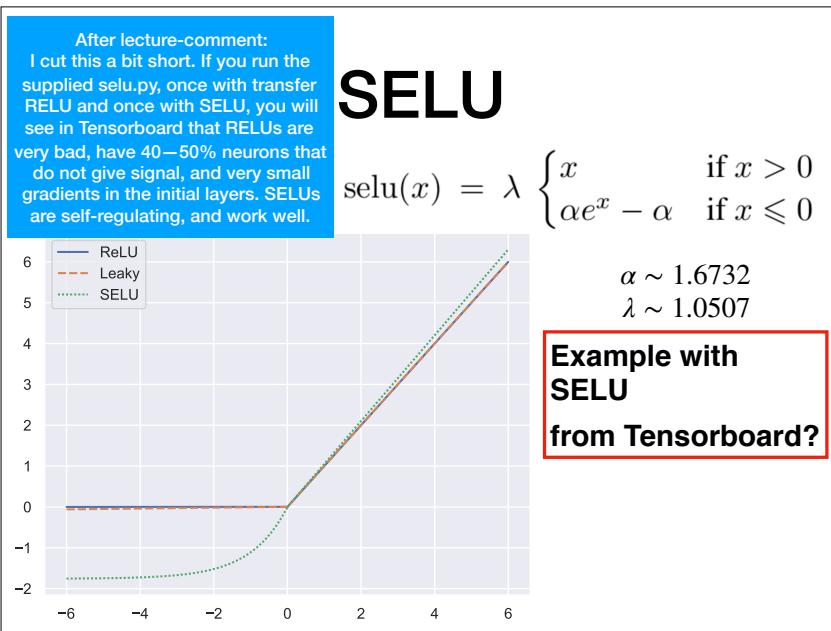
Andreas Mayr

Sepp Hochreiter

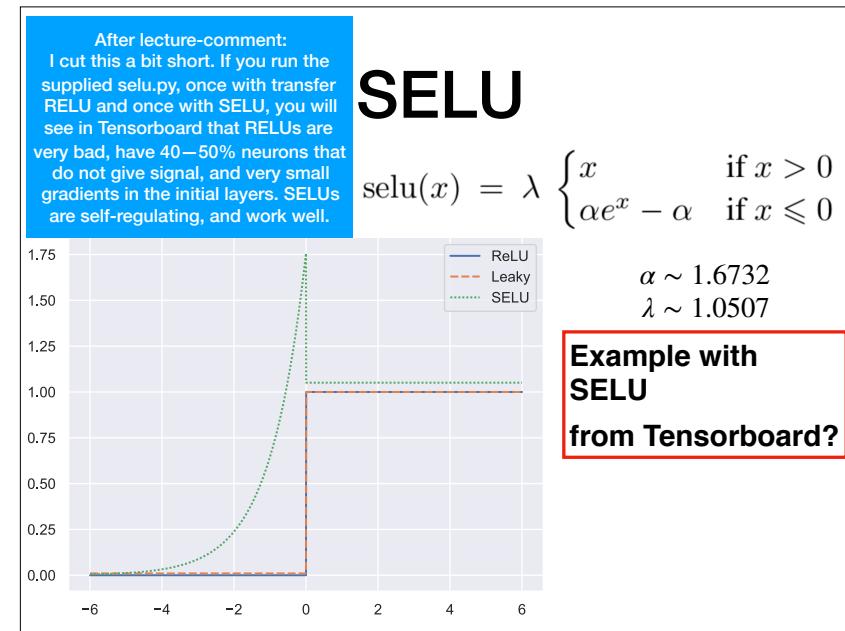
LIT AI Lab & Institute of Bioinformatics,
Johannes Kepler University Linz
A-4040 Linz, Austria

{klambauer, unterthiner, mayr, hochreit}@bioinf.jku.at

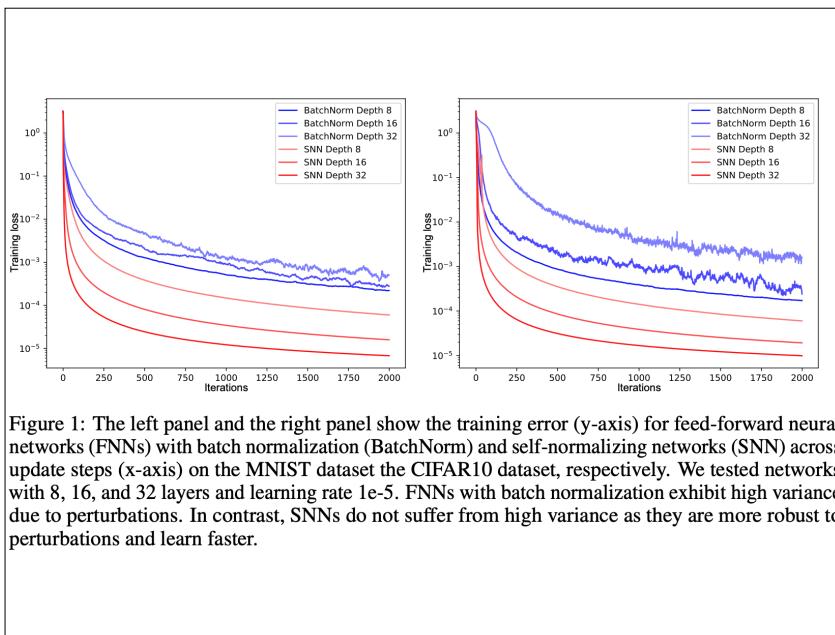
36



37-1



37-2



38

Data Driven Refinement

- Understand your model
 - Measure error on train and test sets
 - Overfitting vs underfitting
- Understand the learning process
 - Use Tensorboard and/or other visualizations
 - Know what to look for, and how to fix problems
- Choose a models with relevant learning bias
 - Don't believe the hype
 - Feed Forward, CNN, RNN, more advanced tricks

39

CNN, RNN, FFNN, ???

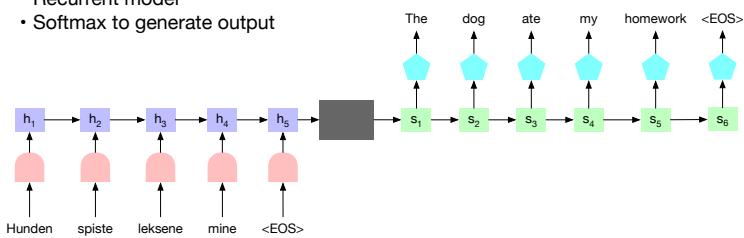
- CNNs are very good at data with localized correlation structure (e.g., neighbouring pixels in picture, or larger scale in hierarchy, i.e., [neighbours of ...] * N pixels)
 - RNN useful for sequences in particular when we can think of the domain as “Markov-model-like”. CNNs also much used in sequence data (... or combos of both)
 - FFNNs can be useful if there is no known structure to relate to; rarely state-of-the-art (e.g., Kaggle) unless a bit deep (>5 layers).
- USE THE DATA TO SELECT YOUR MODEL, AND CHECK IF LEARNING BIAS AND INFORMATION FLOW FITS!!!**

40

Encoder - Decoder

Building blocks:

- Word embeddings
- Recurrent model
- Softmax to generate output



42-1

Sequence-to-Sequence

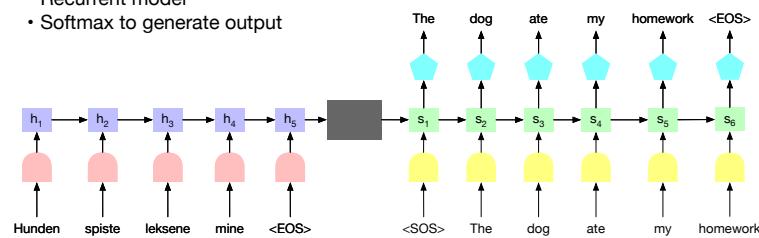
- Standard setup for, e.g., language translation
- For quite some time the “encode-decoder” setup with RNNs was the state-of-the-art.
- Now, moving towards attention-models...

41

Encoder - Decoder

Building blocks:

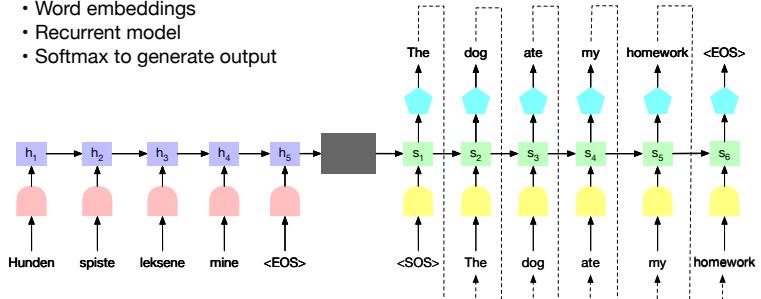
- Word embeddings
- Recurrent model
- Softmax to generate output



42-2

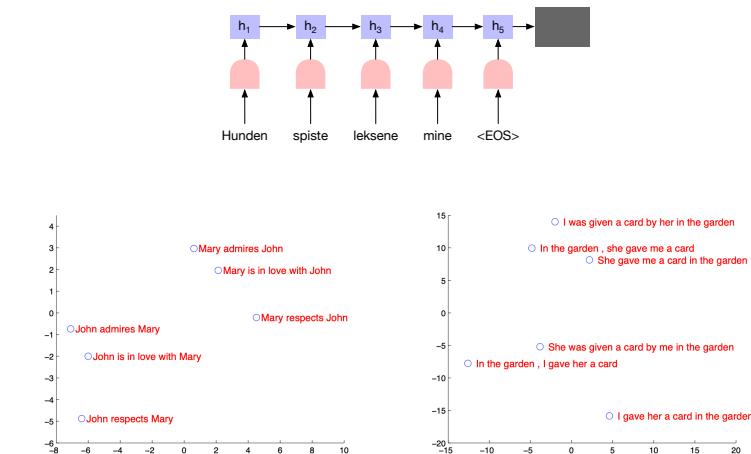
Encoder - Decoder

- Building blocks:
- Word embeddings
 - Recurrent model
 - Softmax to generate output



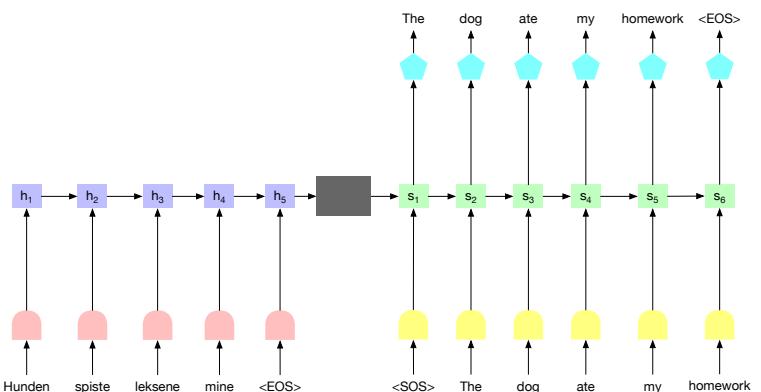
42-3

Does it work?



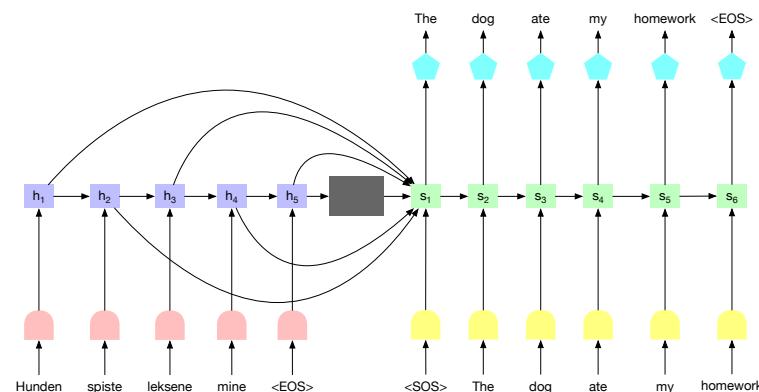
43

What's the problem



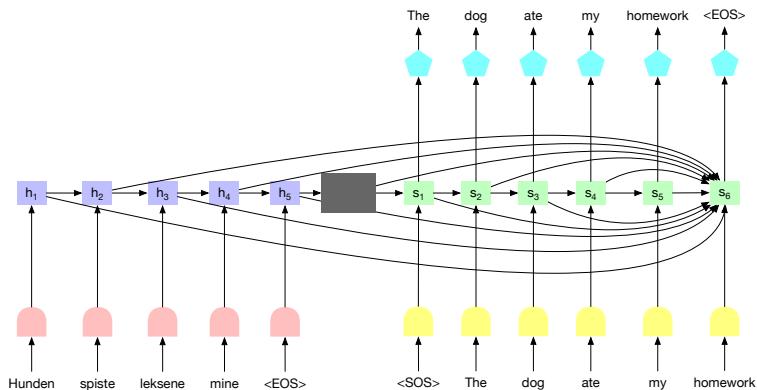
44-1

What's the problem

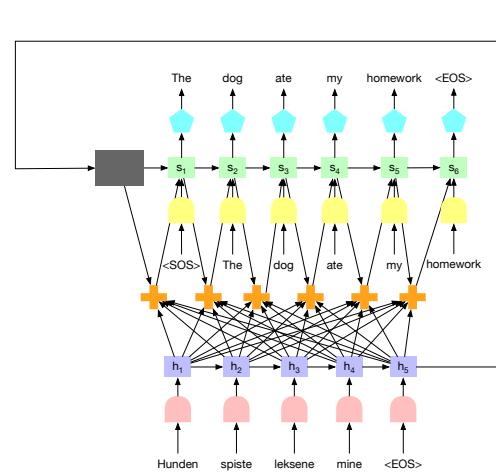


44-2

What's the problem



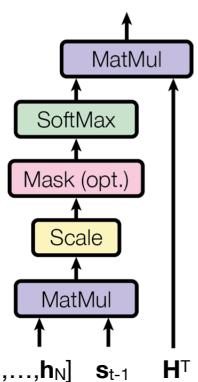
44-3



45

What goes on?

- Attention to vector \mathbf{h}_i depends on vector \mathbf{s}_{t-1} “pointing in similar direction” as \mathbf{h}_i .
 - Random vectors are approximately orthogonal in high-dim.
 - Typically only a few $e_i(t)$ -values dominate
 - This is further strengthened by softmax
- Attention works like a differentiable lookup-table that averages over a (sparse) set. $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]$ \mathbf{s}_{t-1} \mathbf{H}^T



46

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

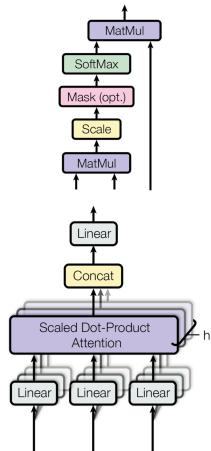
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

47

Multi-head attention

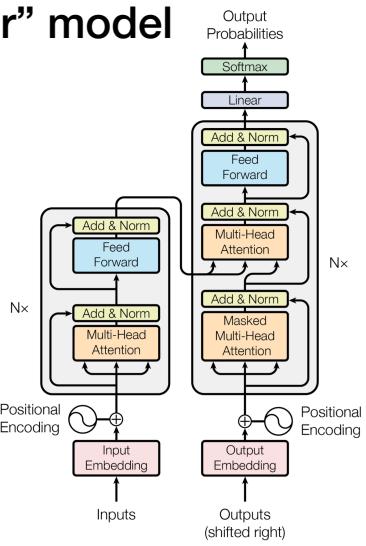


- The (single head) attention module typically picks out one or a few words as important
 - If more than one: Scaled average
- Sometimes several (sets of) words may be useful in “unrelated ways”:
 - One (set) gives the word to translate
 - Another the associated grammar
 - ... etc
- Solution: Several sets of attention, concatenated

48

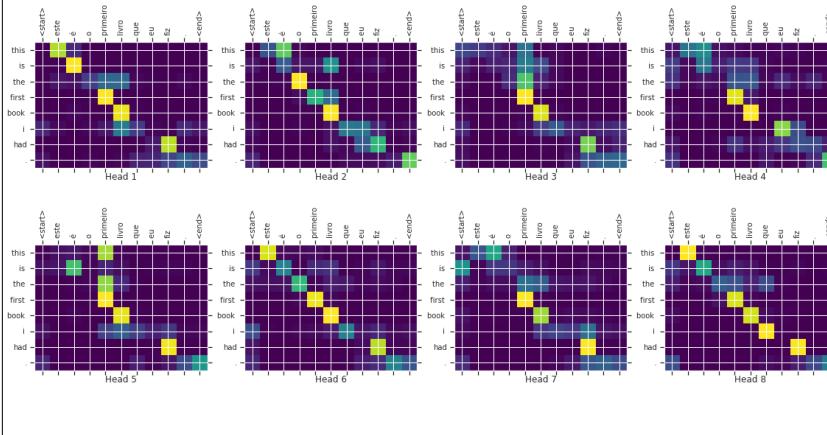
The full “Transformer” model

- No recurrence or convolutions
- Positional encoding + stacked substructures of attention and normalization
- Self-attention at encoder
- Masked self-attention and external attention at decoder
- Skip-connections to maintain gradients
- State of the art results in many (most) seq2seq problems



50

Attention maps Going from Portuguese to English



49

Summary

- Use needs to define metric-based goals
 - Be sure that you get what you need
- Build an end-to-end system
 - Start simple, check what to improve — cost/benefit
- Data-driven refinement
 - Understand your model and its strengths/weaknesses.
 - Optimize hyper-parameters
 - Understand the learning process
 - Choose a models with relevant learning bias

51