

INTRODUCTION TO DEEP LEARNING

Seminar @ UPC TelecomBCN Barcelona (3rd edition). 22-28 January 2020.



Instructors



Xavier
Giró-i-Nieto



Marta R.
Costa-Jussà



Noé
Casas



Verónica
Vilaplana



Ramon
Morros



Javier
Ruiz

Organizers



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Supporters



Google Cloud

GitHub Education

+ info: <http://bit.ly/idl2020>

[Course Site](#)



#DLUPC

Day 2 Lecture 4

Methodology



Javier Ruiz Hidalgo

javier.ruiz@upc.edu

Associate Professor

Universitat Politècnica de Catalunya
Technical University of Catalonia



Outline

- **Data**
 - training, validation, test partitions
 - Augmentation
- **Capacity of the network**
 - Underfitting
 - Overfitting
- **Prevent overfitting**
 - Dropout, regularization
- **Strategy**

Outline



It's all about the data...

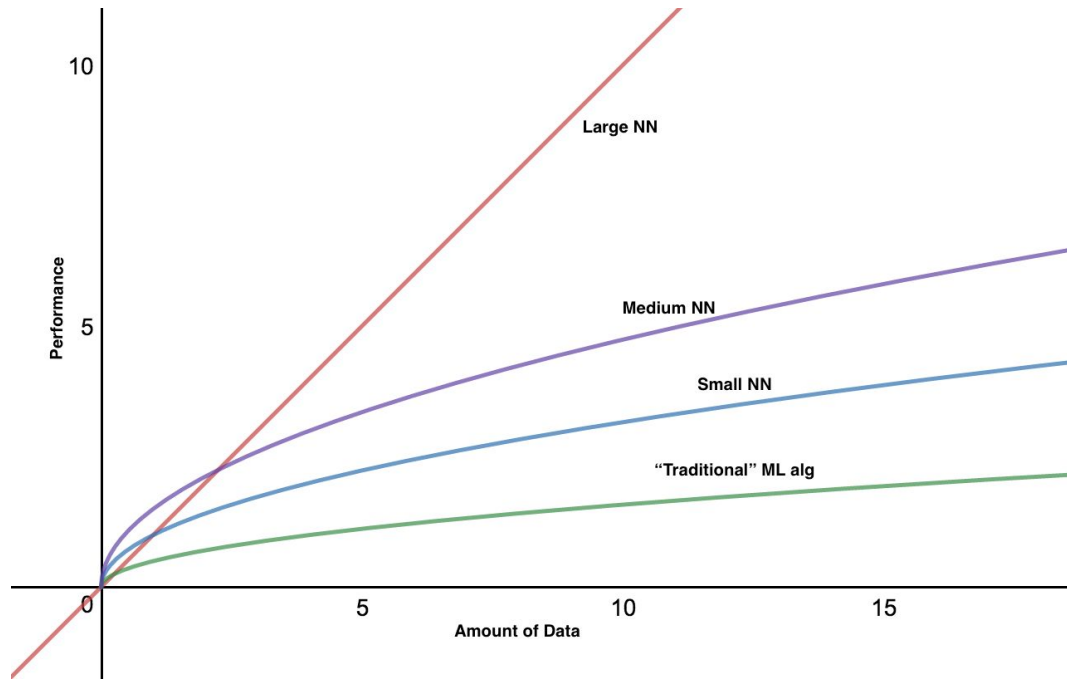
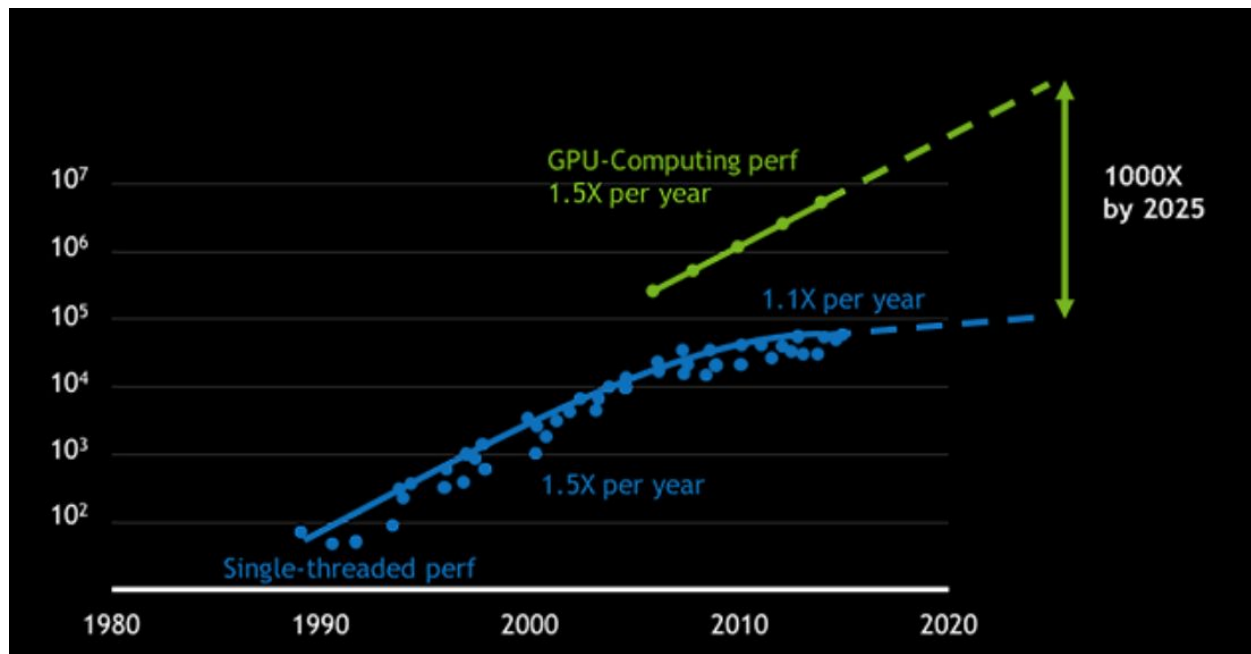


Figure extracted from Kevin Zakka's Blog, [Nuts and Bolts of Applying Deep Learning](#), 2016.

well, not only data...

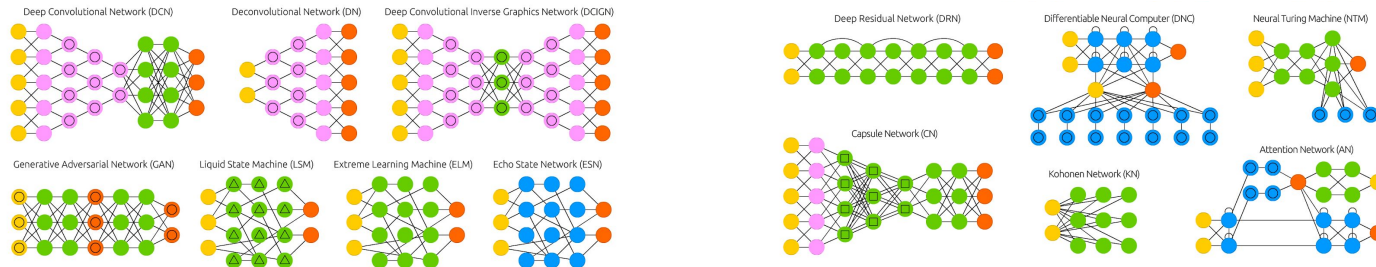
- Computing power: GPUs



Source: NVIDIA 2017

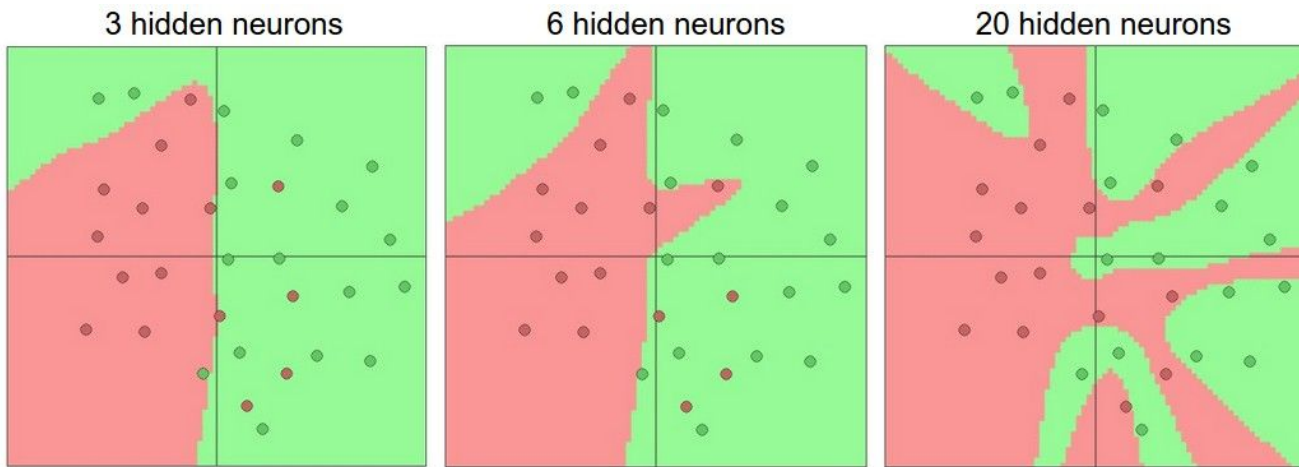
well, not only data...

- Computing power: GPUs
- New learning architectures
 - CNN, RNN, LSTM, DBN, GNN, GAN, Transformers, etc.



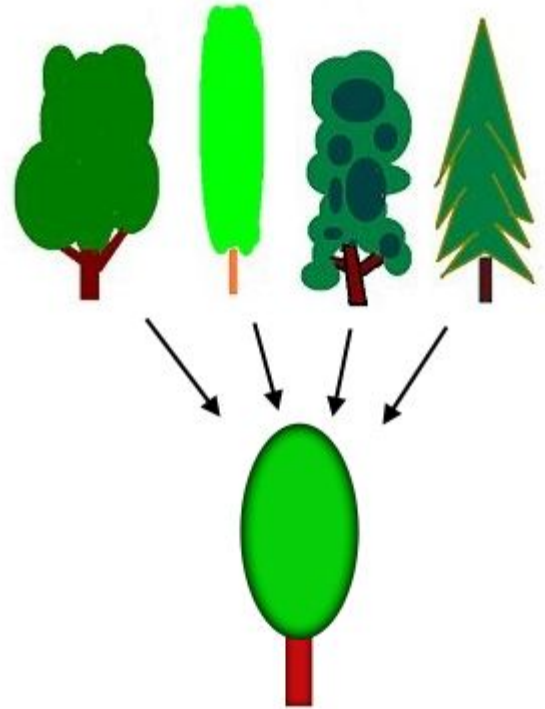
Network capacity

- Space of representable functions that a network can potentially learn:
 - Number of layers / parameters



Generalization

The network needs to **generalize** beyond the training data to work on new data that it has not seen yet

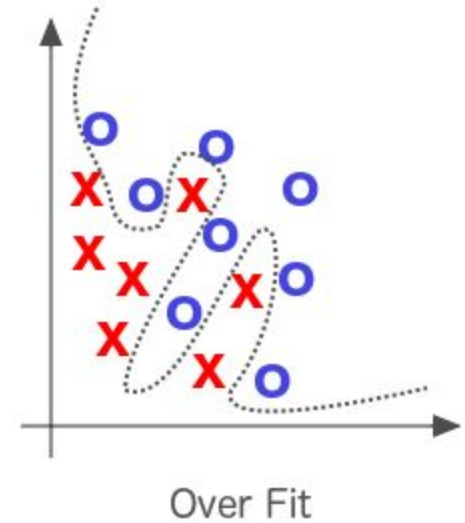
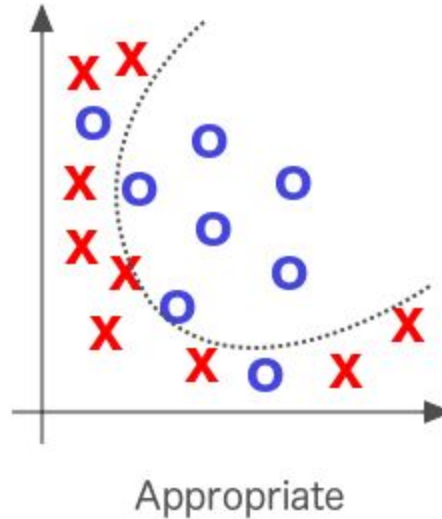
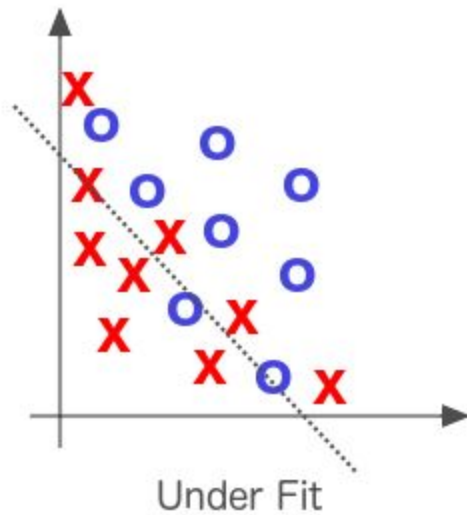


Underfitting vs Overfitting

- **Overfitting:** network fits training data too well
 - Excessively complicated model
- **Underfitting:** network does not fit training data well enough
 - Excessively simple model

Both underfitting and overfitting lead to poor predictions on new data and they do not generalize well

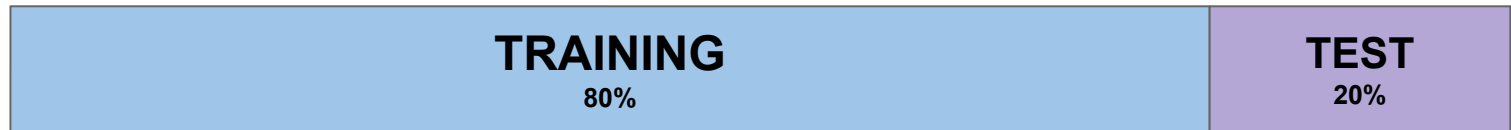
Underfitting vs Overfitting



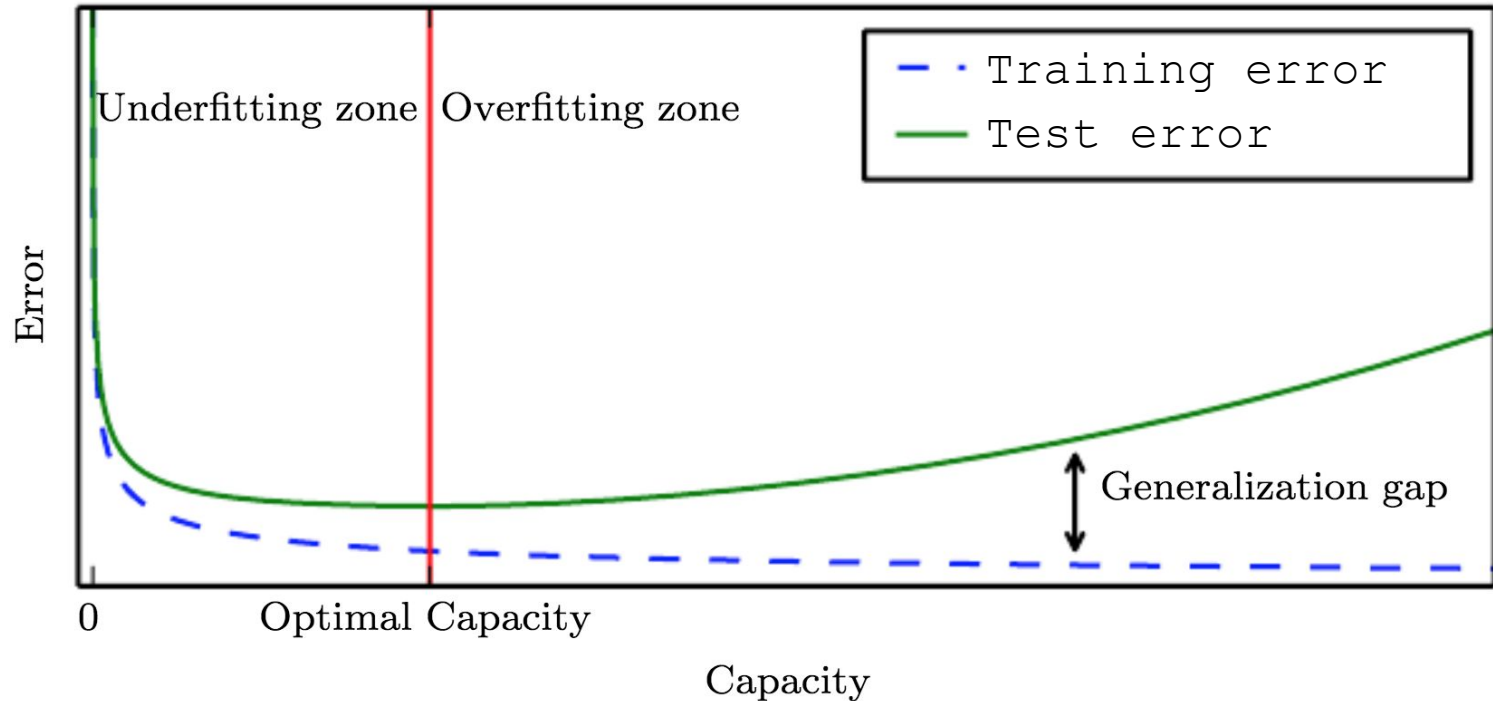
Data partition

How do we measure the generalization instead of how well the network does with the memorized data?

Split your data into two sets: training and test



Underfitting vs Overfitting



Data partition revisited

- Test set **should not** be used to tune your network
 - Network architecture
 - Number of layers
 - Hyper-parameters
- Failing to do so will overfit the network to your test set!
 - <https://www.kaggle.com/c/higgs-boson/leaderboard>

Data partition revisited (2)

- Add a **validation** set!



- Lock away your test set and use it only as a last validation step

The bigger the better?

- Larger networks
 - More capacity / More data
 - Prone to overfit
- Smaller networks
 - Lower capacity / Less data
 - Prone to underfit



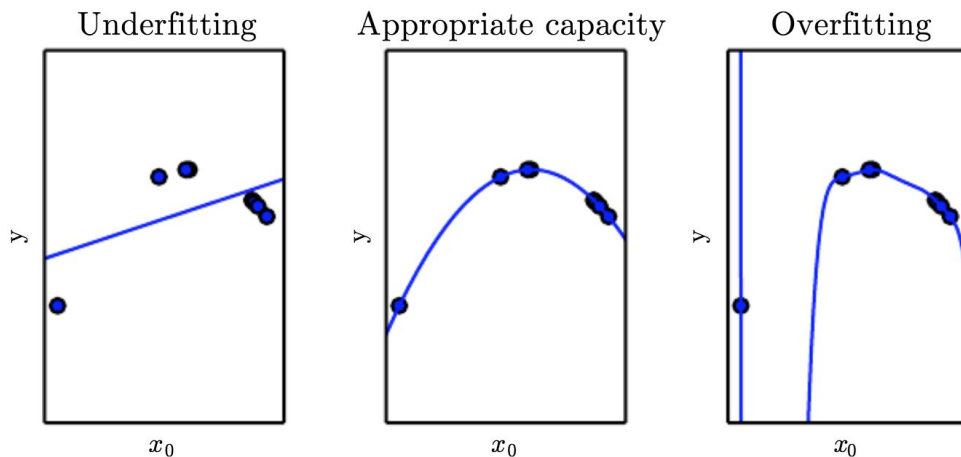
The bigger the better?

- In large networks, most local minima are equivalent and yield similar performance.
- The probability of finding a “bad” (high value) local minimum is non-zero for small networks and decreases quickly with network size.
- Struggling to find the global minimum on the training set (as opposed to one of the many good local ones) is not useful in practice and may lead to overfitting.

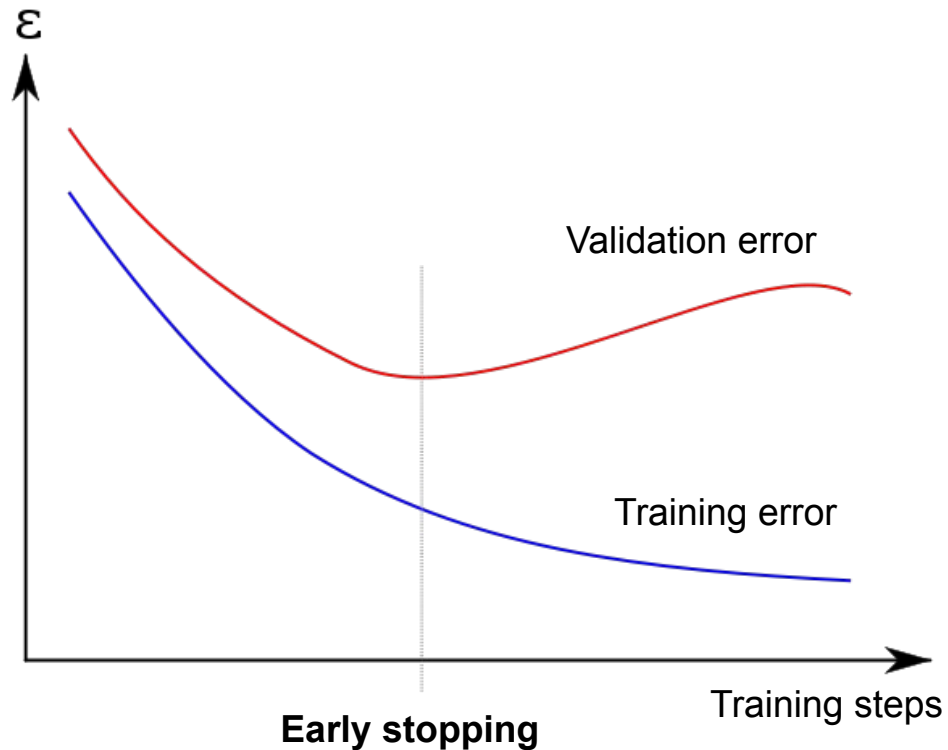
Better large capacity networks and prevent overfitting

Prevent overfitting

- Early stopping
- Loss regularization
- Data augmentation
- Dropout

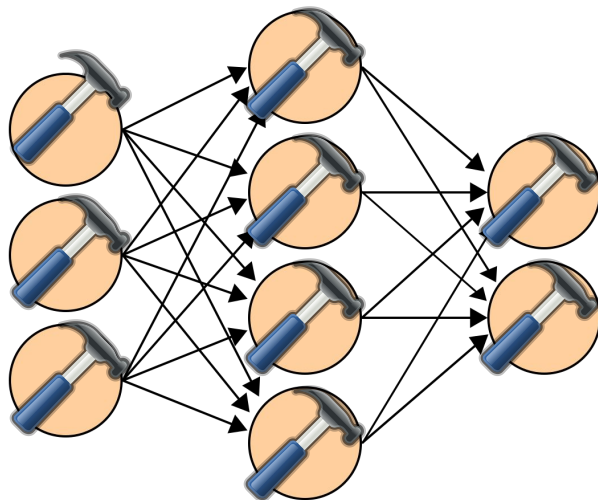


Early stopping



Loss regularization (1)

- Control the capacity of the network to **prevent overfitting**
- Large weights tend to cause sharp transitions in node functions → large changes in output for small changes in inputs
 - **Penalize the weights** of the nodes in the network
 - **Discourages** learning a more **complex** or flexible **model**



Loss regularization (2)

- L2-regularization (weight decay):

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} W^2$$

regularization
hyper-parameter
↙

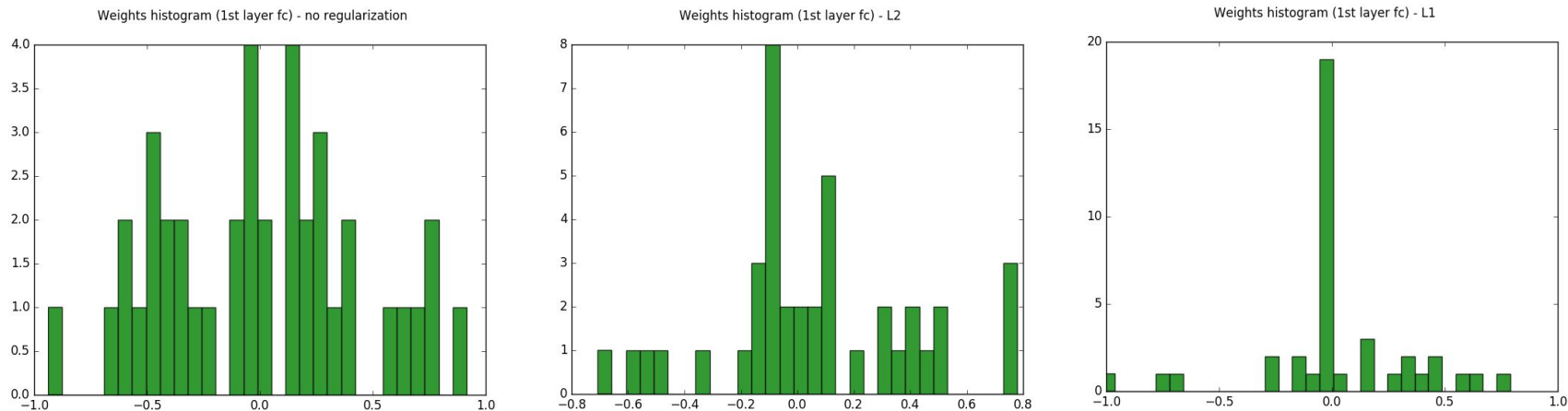
- L1-regularization:

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} |W|$$

Loss regularization (3)

- Limit the values of parameters in the network
 - L2 vs L1 regularization

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} W^2 \quad \mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} |W|$$



Loss regularization (4)

- L2 regularization heavily penalizes peaky weights and prefers diffuse / low value weights
- L1 regularization leads weights to become sparse (i.e. very close to exactly zero)

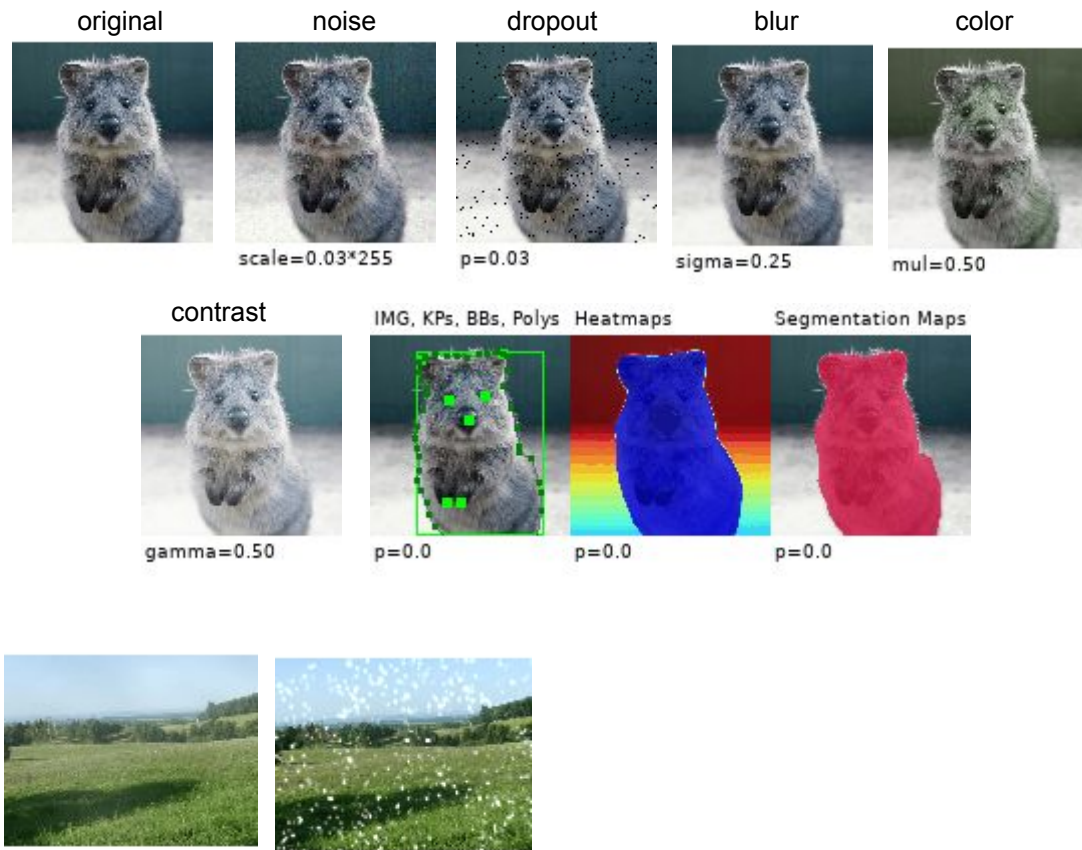
Data augmentation (1)

- Modify input samples artificially to increase the data size
- On-the-fly while training
 - Inject Noise
 - Transformations
- Not used in testing/validation



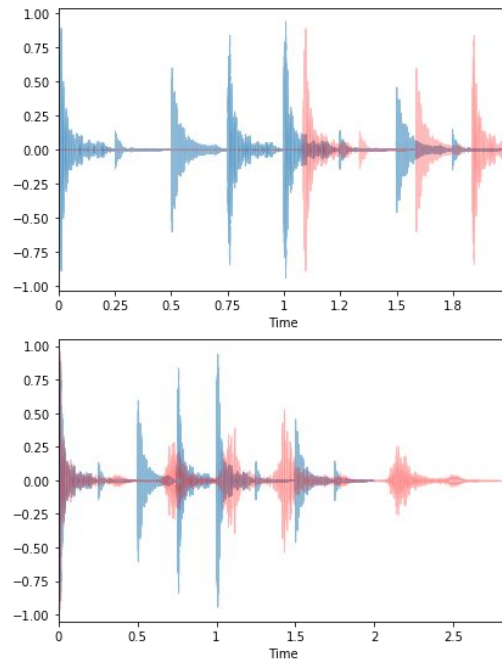
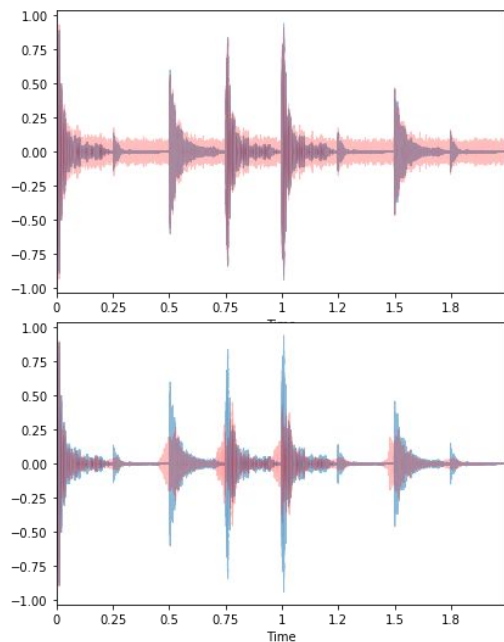
Data augmentation (2): Image

- Noise injection
- Dropout
- Blurs
- Color changes
- Contrast
- Transformations
 - GT transformed!
- Crops, shifts
- Application specific
 - Clouds, snow, etc.



Data augmentation (3): Audio

- Noise injection
- Shifting time
- Changing pitch
- Changing speed
- Crops
- Loudness
- Masks



Data augmentation (4)

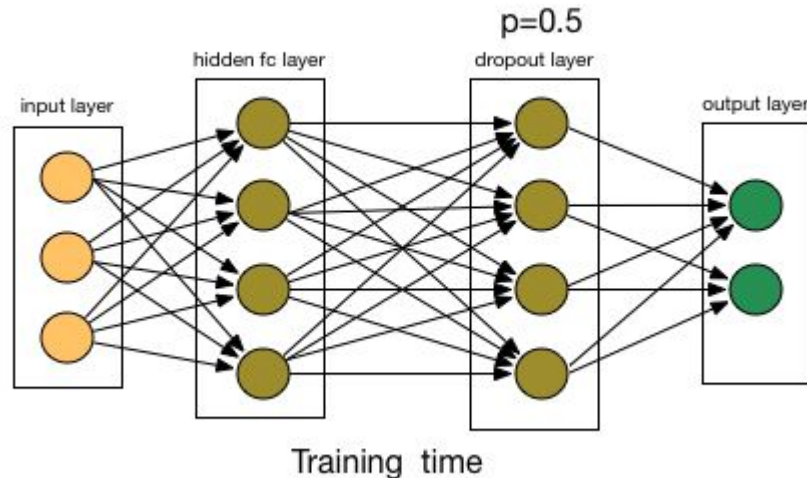
- Synthetic data: Generate new input samples



A. Palazzi, [Learning to Map Vehicles into Bird's Eye View](#), ICIAP 2017
DeepGTAV plugin: <https://github.com/ai-tor/DeepGTAV>
[CARLA Simulator](#).

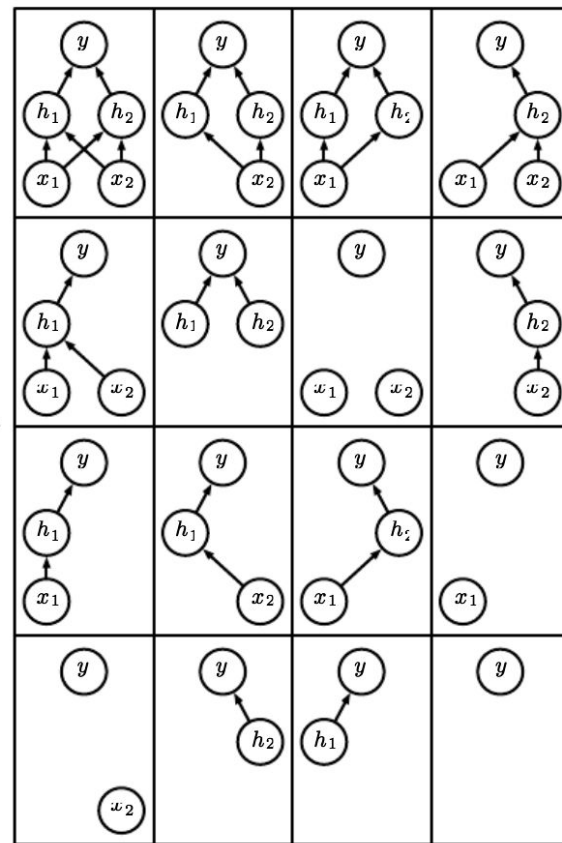
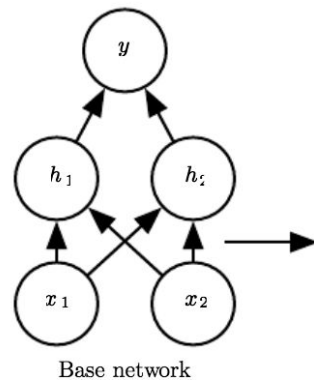
Dropout (1)

- At each training iteration, randomly remove some nodes in the network along with all of their incoming and outgoing connections ([N. Srivastava, 2014](#))



Dropout (2)

- Why dropout works?
 - Nodes become more insensitive to the weights of the other nodes → more robust.
 - Averaging multiple models → **ensemble**.
 - Training a collection of 2^n thinned networks with parameters sharing

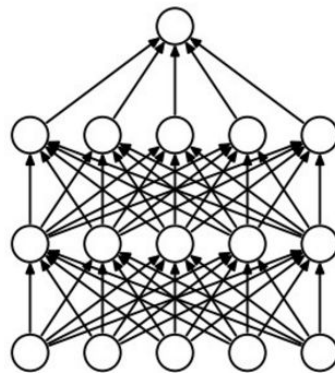


Ensemble of subnetworks

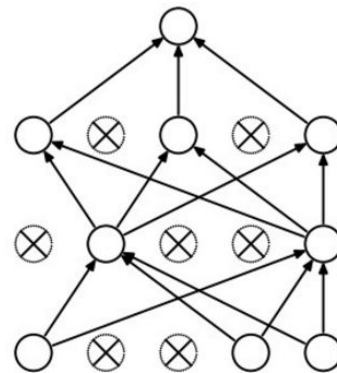
Dropout (3)

- Every forward pass, network slightly different.
- Reduce co-adaptation between neurons
- More robust features
- Dropout is **removed** in validation/testing

✗ **More iterations** for convergence



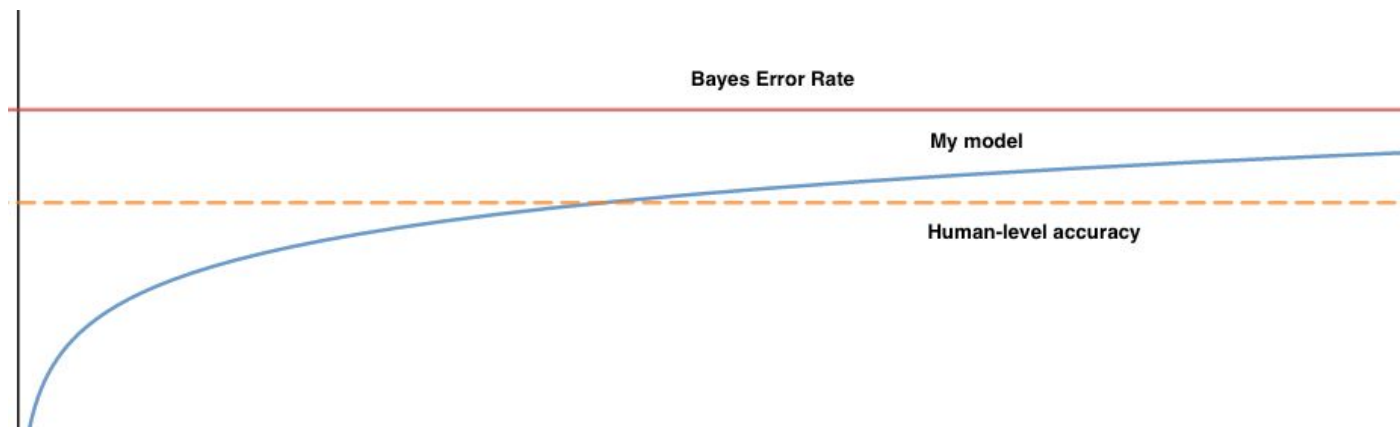
(a) Standard Neural Net



(b) After applying dropout.

Strategy for machine learning (1)

Human-level performance can serve as a very reliable proxy which can be leveraged to determine your next move when training your model.



Strategy for machine learning (2)

TRAINING 60%	VALIDATION 20%	TEST 20%
-----------------	-------------------	-------------

Human level error . . 1%

Training error . . . 19%

Validation error . . 20%

Test error 21%



Underfitting

Strategy for machine learning (3)

TRAINING 60%	VALIDATION 20%	TEST 20%
-----------------	-------------------	-------------

Human level error . . 1%

Training error . . . 1.1%

Validation error . . 20%

Test error 21%



Overfitting training

Strategy for machine learning (4)

TRAINING 60%	VALIDATION 20%	TEST 20%
-----------------	-------------------	-------------

Human level error . . 1%

Training error . . . 1.1%

Validation error . . 2%

Test error 21%



Overfitting validation

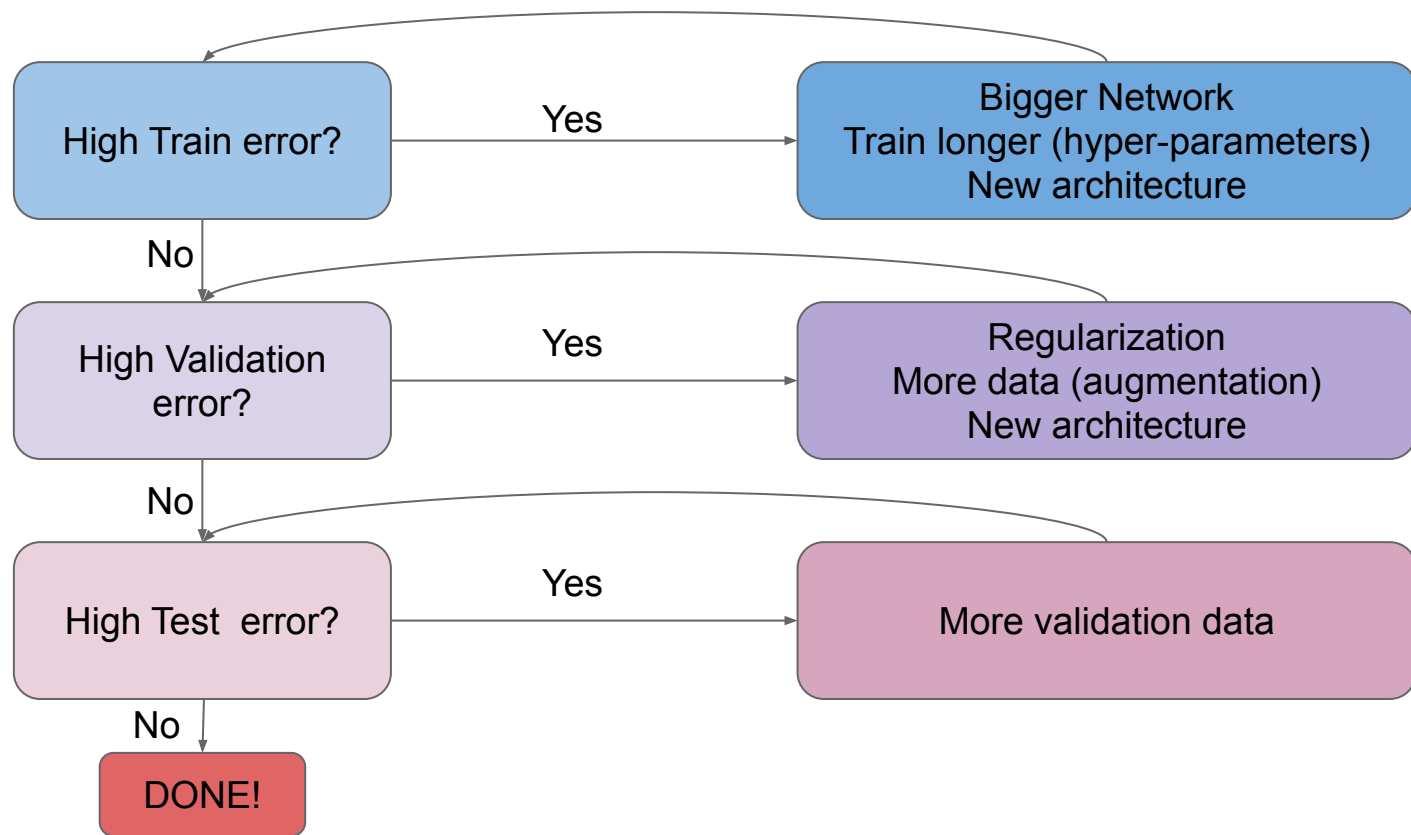
Strategy for machine learning (5)

TRAINING 60%	VALIDATION 20%	TEST 20%
-----------------	-------------------	-------------

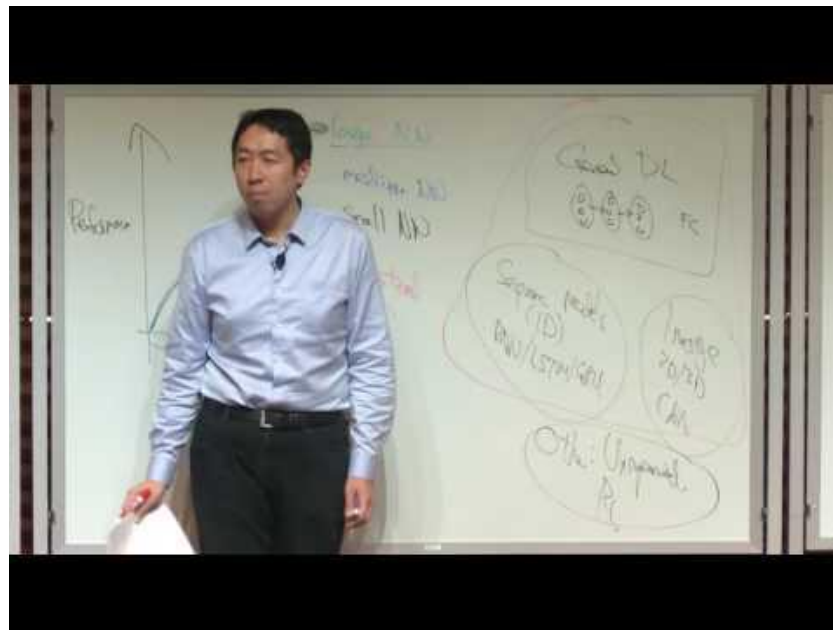
Human level error	.	.	1%
Training error	.	.	1.1%
Validation error	.	.	1.2%
Test error	.	.	1.2%



Strategy for machine learning (5)



References



Nuts and Bolts of Applying Deep Learning by Andrew Ng

<https://www.youtube.com/watch?v=F1ka6a13S9I>

Thanks! Questions?



UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH

**Department of Signal Theory
and Communications**

Image Processing Group

<https://imatge.upc.edu/web/people/javier-ruiz-hidalgo>