



Parallel and Scalable Machine Learning

Introduction to Machine Learning Models

Prof. Dr. – Ing. Morris Riedel

Associated Professor

School of Engineering and Natural Sciences, University of Iceland, Reykjavik, Iceland

Research Group Leader, Juelich Supercomputing Centre, Forschungszentrum Juelich, Germany

LECTURE 5

 @Morris Riedel

 @MorrisRiedel

 @MorrisRiedel

Introduction to Statistical Learning Theory

February 18, 2020

Juelich Supercomputing Centre, Germany



UNIVERSITY OF ICELAND
SCHOOL OF ENGINEERING AND NATURAL SCIENCES
FACULTY OF INDUSTRIAL ENGINEERING,
MECHANICAL ENGINEERING AND COMPUTER SCIENCE



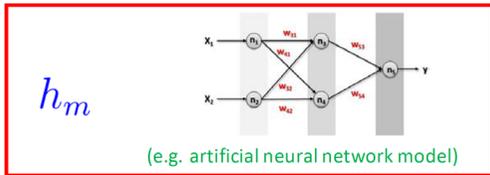
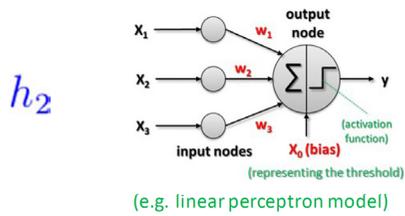
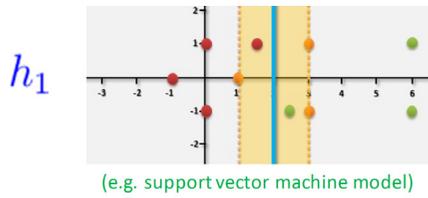
JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE



HELMHOLTZAI | ARTIFICIAL INTELLIGENCE
COOPERATION UNIT

Review of Lecture 4 – Artificial Neural Networks (ANNs)

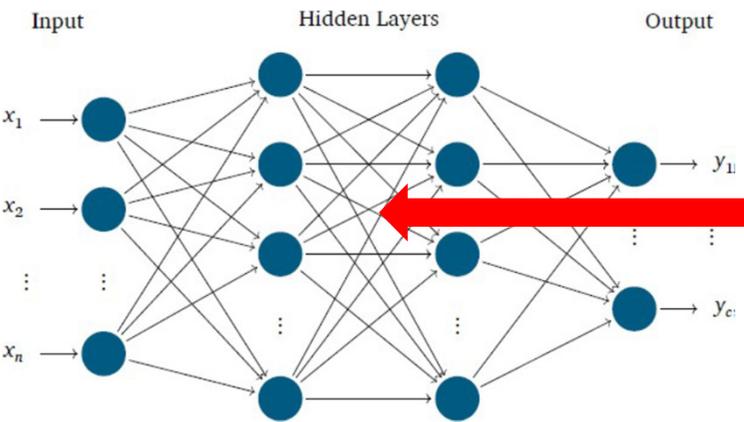
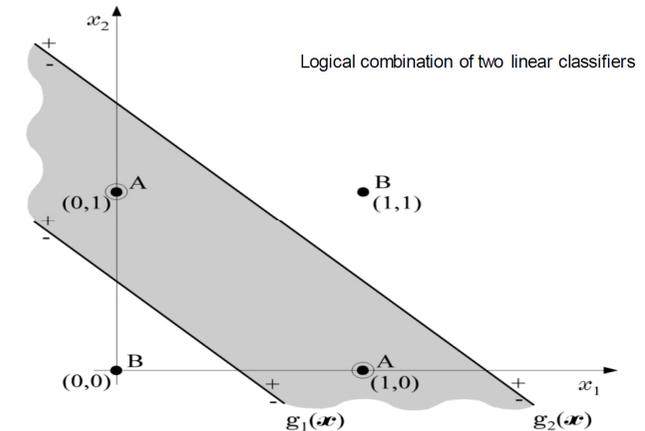
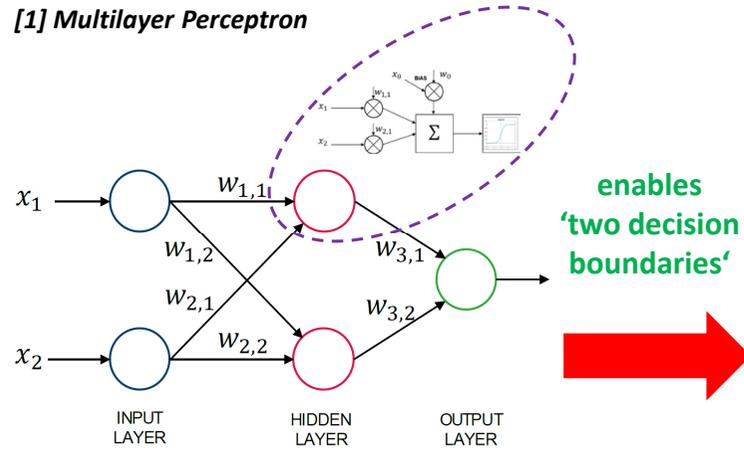


```
import numpy as np
from keras.datasets import mnist

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

introducing hidden layers concept

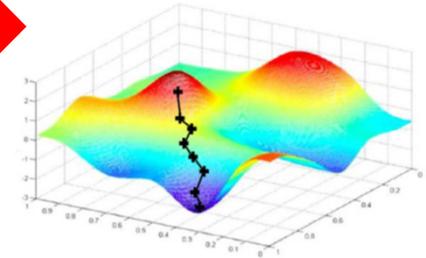
[1] Multilayer Perceptron



$$W := W - \eta \frac{\partial \mathcal{L}(W)}{\partial W}$$

[2] MIT Course

How to set the learning rate?



Outline of the Training Course

1. Parallel & Scalable Machine Learning driven by HPC
2. Introduction to Machine Learning Fundamentals
3. Supervised Learning with a Simple Learning Model
4. Artificial Neural Networks (ANNs)
5. Introduction to Statistical Learning Theory
6. Validation and Regularization
7. Pattern Recognition Systems
8. Parallel and Distributed Training of ANN
9. Supervised Learning with Deep Learning
10. Unsupervised Learning – Clustering
11. Clustering with HPC
12. Introduction to Deep Reinforcement Learning



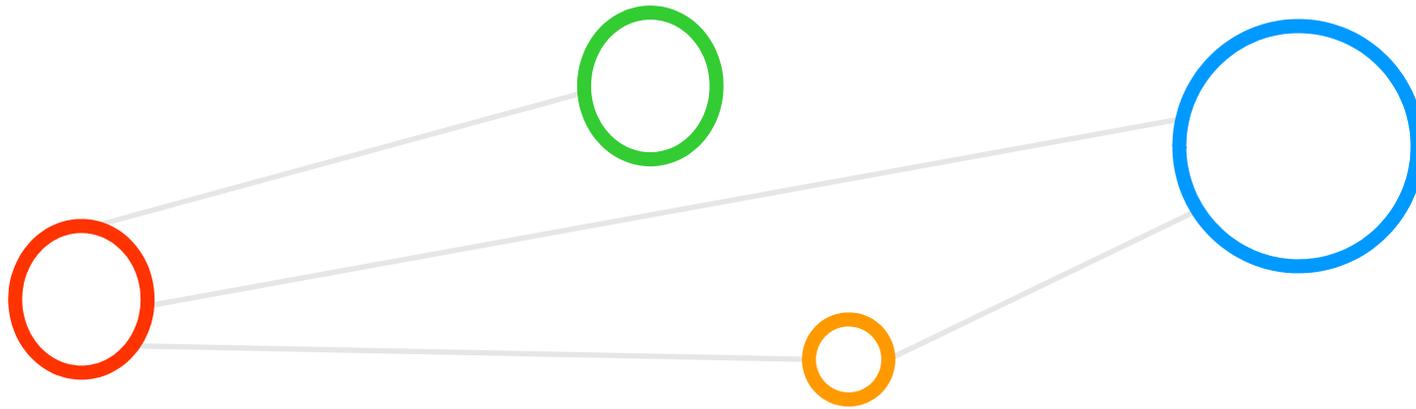
Outline

- Supervised Learning & Statistical Learning Theory
 - Statistical Learning Theory is Best Understood for Supervised Learning
 - Refining Mathematic Building Blocks for Supervised Learning
 - Feasibility of Learning from Data & Hoeffding Inequality
 - Probably Approximately Correct (PAC) Learning
 - Learning Model Perceptron Examples in Context

- Vapnik – Chervonenkis (VC) Inequality & Dimension
 - Theory of Generalization – Revisited & Reviewed
 - Infinite Learning Models Challenge & Union Bound Solution
 - Vapnik – Chervonenkis (VC) Inequality & Dimension
 - Relationship of Number of Samples & Model Complexity
 - Overfitting as Key Problem in Machine Learning



Supervised Learning & Statistical Learning Theory



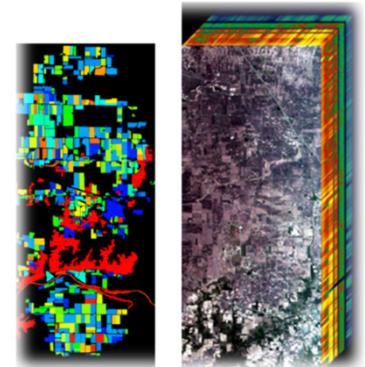
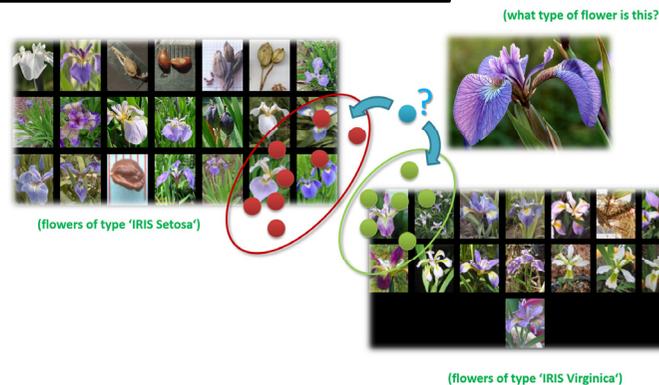
Learning Approaches – What means Learning from data – Revisited

- The basic meaning of learning is ‘to use a set of observations to uncover an underlying process’
- The three different learning approaches are supervised, unsupervised, and reinforcement learning

[3] Image sources: Species Iris Group of North America Database, www.signa.org

Supervised Learning

- Majority of methods follow this approach in this course
- Example: credit card approval based on previous customer applications

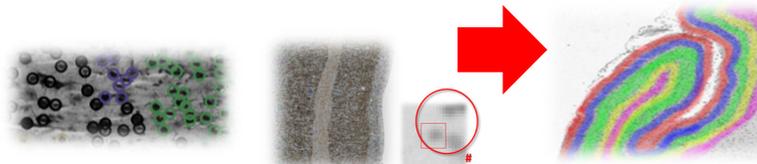


Unsupervised Learning

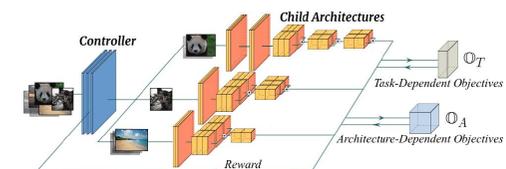
- Often applied before other learning → higher level data representation
- Example: Coin recognition in vending machine based on weight and size

Reinforcement Learning

- Typical ‘human way’ of learning
- Example: Toddler tries to touch a hot cup of tea (again and again)



[4] A.C. Cheng et al., ‘InstaNAS: Instance-aware Neural Architecture Search’, 2018



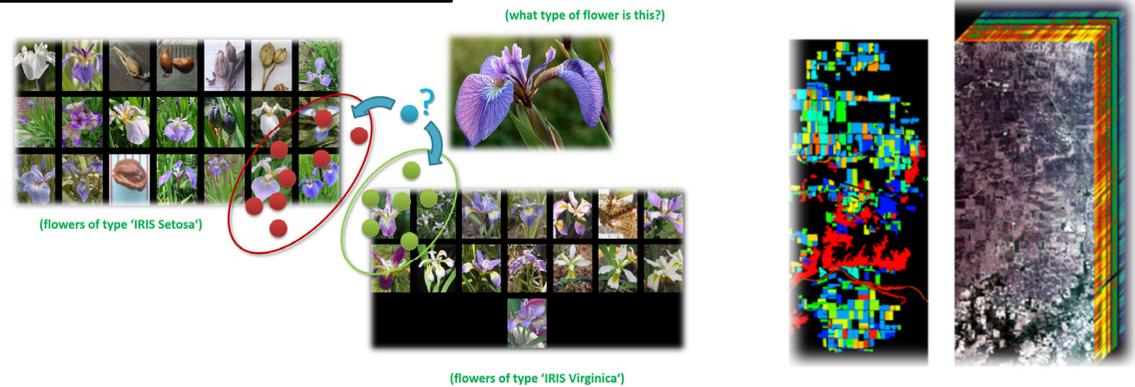
Supervised Learning – Statistical Learning Theory Perspective

- The basic meaning of learning is ‘to use a set of observations to uncover an underlying process’
- The three different learning approaches are supervised, unsupervised, and reinforcement learning

[3] Image sources: Species Iris Group of North America Database, www.signa.org

Supervised Learning

- Majority of methods follow this approach in this course
- Example: credit card approval based on previous customer applications



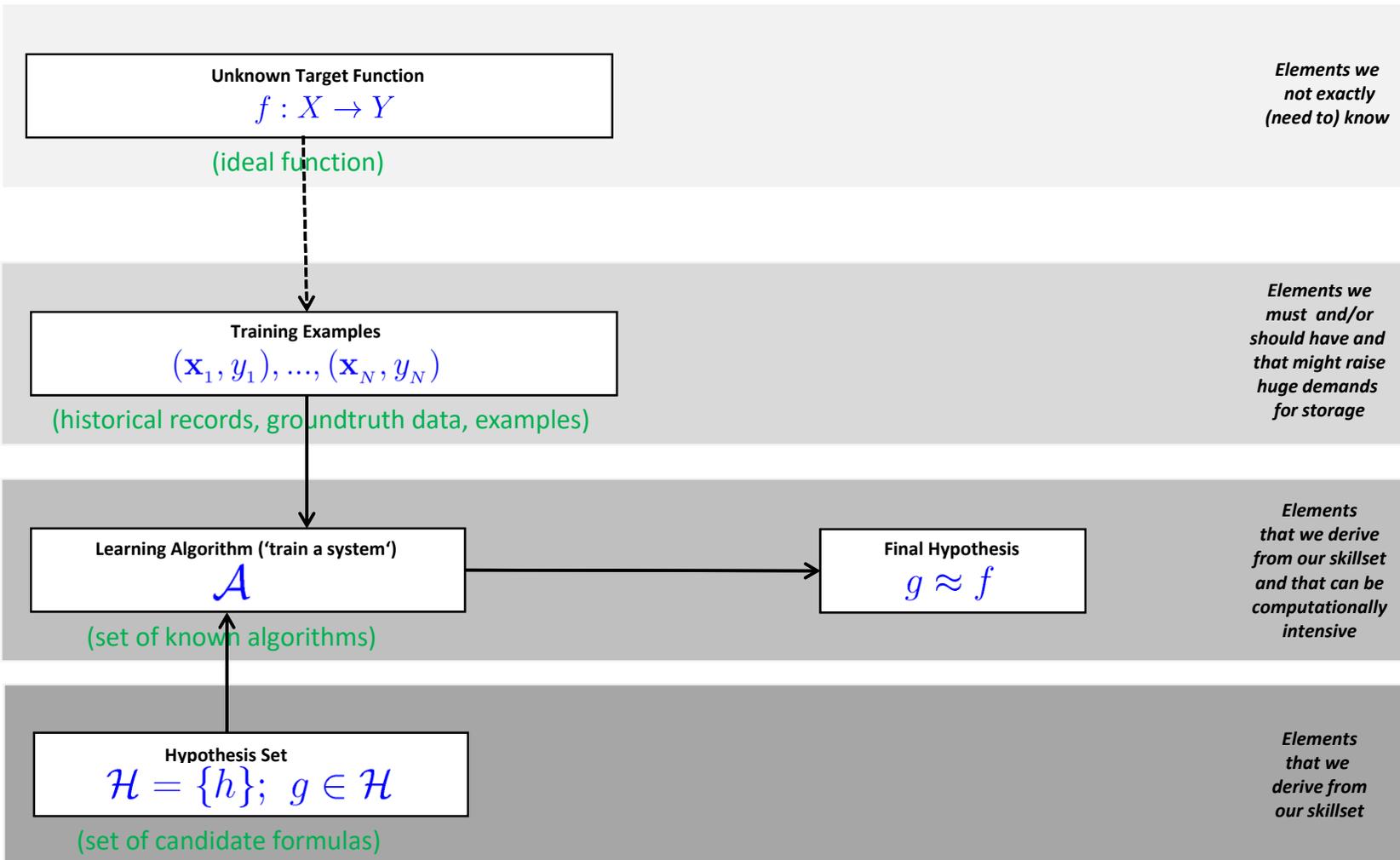
Statistical Learning Theory Perspective

- Supervised Learning is best understood
- Represents a ‘framework’ in which learning from **examples can be studied in a principled way**
- Gives indicators for selected questions:
How much data we need for the model we want to use?
- **This lecture just introduces the topic** (much more deeper & complex)

Statistical Learning Theory: A Primer

THEODOROS EVGENIOU, MASSIMILIANO PONTIL AND TOMASO POGGIO
*Center for Biological and Computational Learning, Artificial Intelligence Laboratory, MIT,
Cambridge, MA, USA*

[5] *Evgeniou T., Pontil M., Poggio T.,
'Statistical Learning Theory: A Primer'*



Different Models – Hypothesis Set & Unlimited ‘Degrees of Freedom’ – Revisited

Hypothesis Set

$$\mathcal{H} = \{h\}; g \in \mathcal{H}$$

$$\mathcal{H} = \{h_1, \dots, h_m\};$$

(all candidate functions derived from models and their parameters)

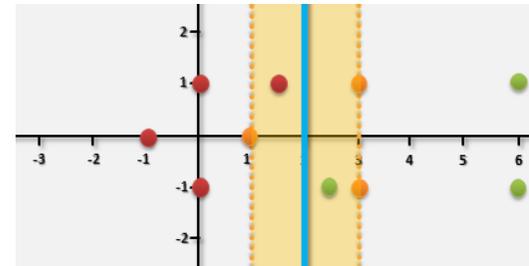
- Choosing from various model approaches h_1, \dots, h_m is a different hypothesis
- Additionally a change in model parameters of h_1, \dots, h_m means a different hypothesis too

‘select one function’ that best approximates

Final Hypothesis

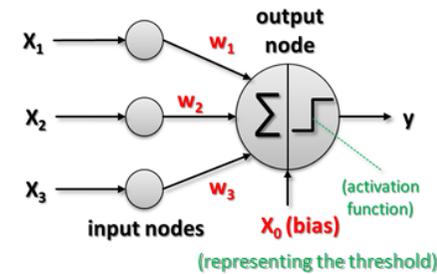
$$g \approx f$$

h_1



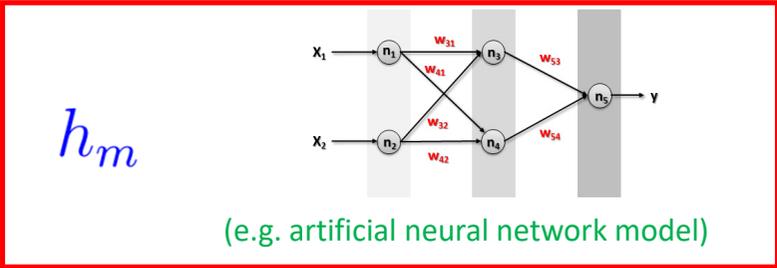
(e.g. support vector machine model)

h_2



(e.g. linear perceptron model)

h_m



(e.g. artificial neural network model)

Solutions – Train on Testing Dataset & Test on Training Dataset & Increase Epochs

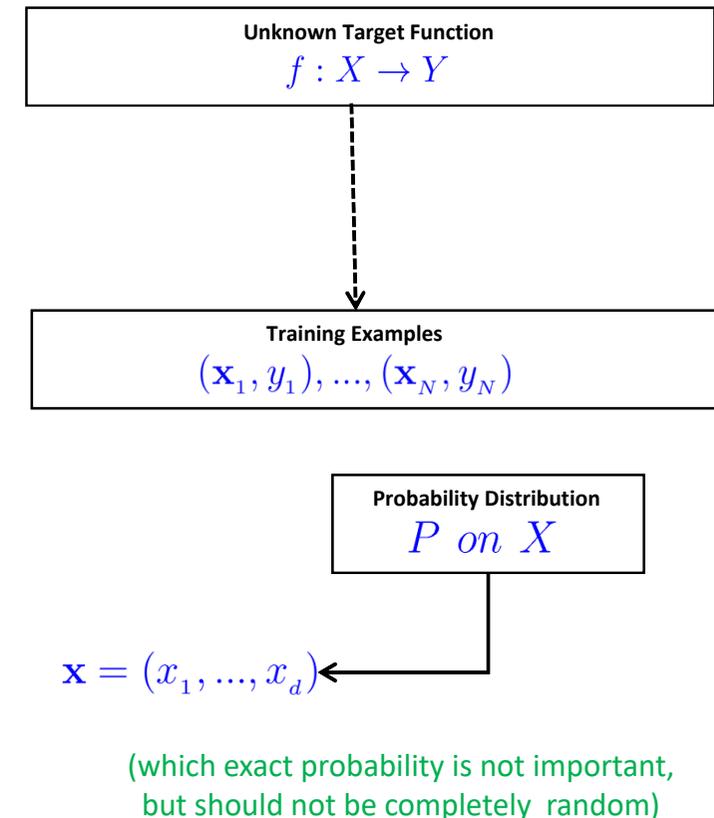
- It seems the number of N samples matter in learning – why?



Feasibility of Learning – Probability Distribution

- Predict output from future input (fitting existing data is not enough)
 - In-sample ‘1000 points’ fit well
 - Possible: Out-of-sample \geq ‘1001 point’ doesn’t fit very well
 - Learning ‘any target function’ is not feasible (can be anything)
- Assumptions about ‘future input’
 - Statement is possible to define about the data outside the in-sample data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
 - All samples (also future ones) are derived from same ‘unknown probability’ distribution P on X

■ Statistical Learning Theory assumes an unknown probability distribution over the input space X

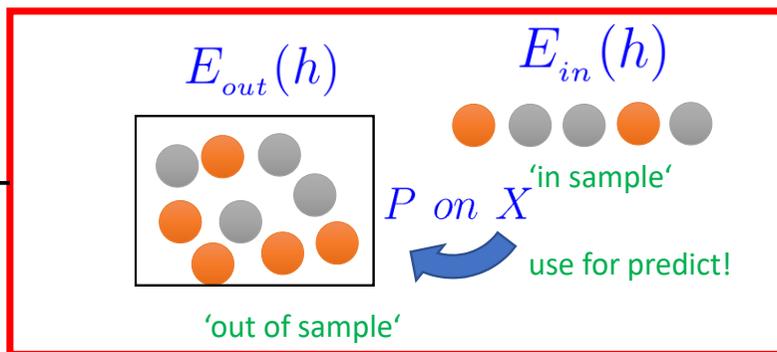


Feasibility of Learning – In Sample vs. Out of Sample

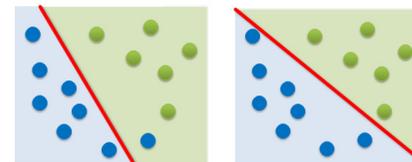
- Given ‘unknown’ probability P on X
 - Given large sample N for $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
 - There is a probability of ‘picking one point or another’
 - ‘Error on in sample’ is known quantity (using labelled data): $E_{in}(h)$
 - ‘Error on out of sample’ is unknown quantity: $E_{out}(h)$
 - In-sample frequency is likely close to out-of-sample frequency

Statistical Learning Theory part that enables that learning is feasible in a probabilistic sense (P on X)

depend on which hypothesis h out of M different ones



$$\mathcal{H} = \{h_1, \dots, h_m\};$$



$$E_{in}(h) \approx E_{out}(h)$$

use $E_{in}(h)$ as a proxy thus the other way around in learning

$$E_{out}(h) \approx E_{in}(h)$$

Feasibility of Learning – Union Bound & Factor **M**

- Assuming no overlaps in hypothesis set
 - Apply very ‘poor’ mathematical rule ‘union bound’
 - (Note the usage of g instead of h , we need to visit all)

Final Hypothesis
 $g \approx f$

Think if E_{in} deviates from E_{out} with more than tolerance ϵ it is a ‘bad event’ in order to apply union bound

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq \Pr [| E_{in}(h_1) - E_{out}(h_1) | > \epsilon$$

$$\text{or } | E_{in}(h_2) - E_{out}(h_2) | > \epsilon \dots$$

$$\text{or } | E_{in}(h_M) - E_{out}(h_M) | > \epsilon]$$

‘visiting **M**
different
hypothesis’

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq \sum_{m=1}^M \Pr [| E_{in}(h_m) - E_{out}(h_m) | > \epsilon]$$

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq \sum_{m=1}^M 2e^{-2\epsilon^2 N}$$

fixed quantity for each hypothesis
obtained from Hoeffdings Inequality

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

problematic: if **M** is too big we loose the link
between the in-sample and out-of-sample

▪ The union bound means that (for any countable set of m ‘events’) the probability that at least one of the events happens is not greater than the sum of the probabilities of the m individual ‘events’

Feasibility of Learning – Modified Hoeffding’s Inequality – PAC Learning

- Errors in-sample $E_{in}(g)$ track errors out-of-sample $E_{out}(g)$
 - Statement is made being ‘Probably Approximately Correct (PAC)’
 - Given M as number of hypothesis of hypothesis set \mathcal{H}
 - ‘Tolerance parameter’ in learning ϵ
 - Mathematically established via ‘modified Hoeffdings Inequality’: (original Hoeffdings Inequality doesn’t apply to multiple hypothesis)

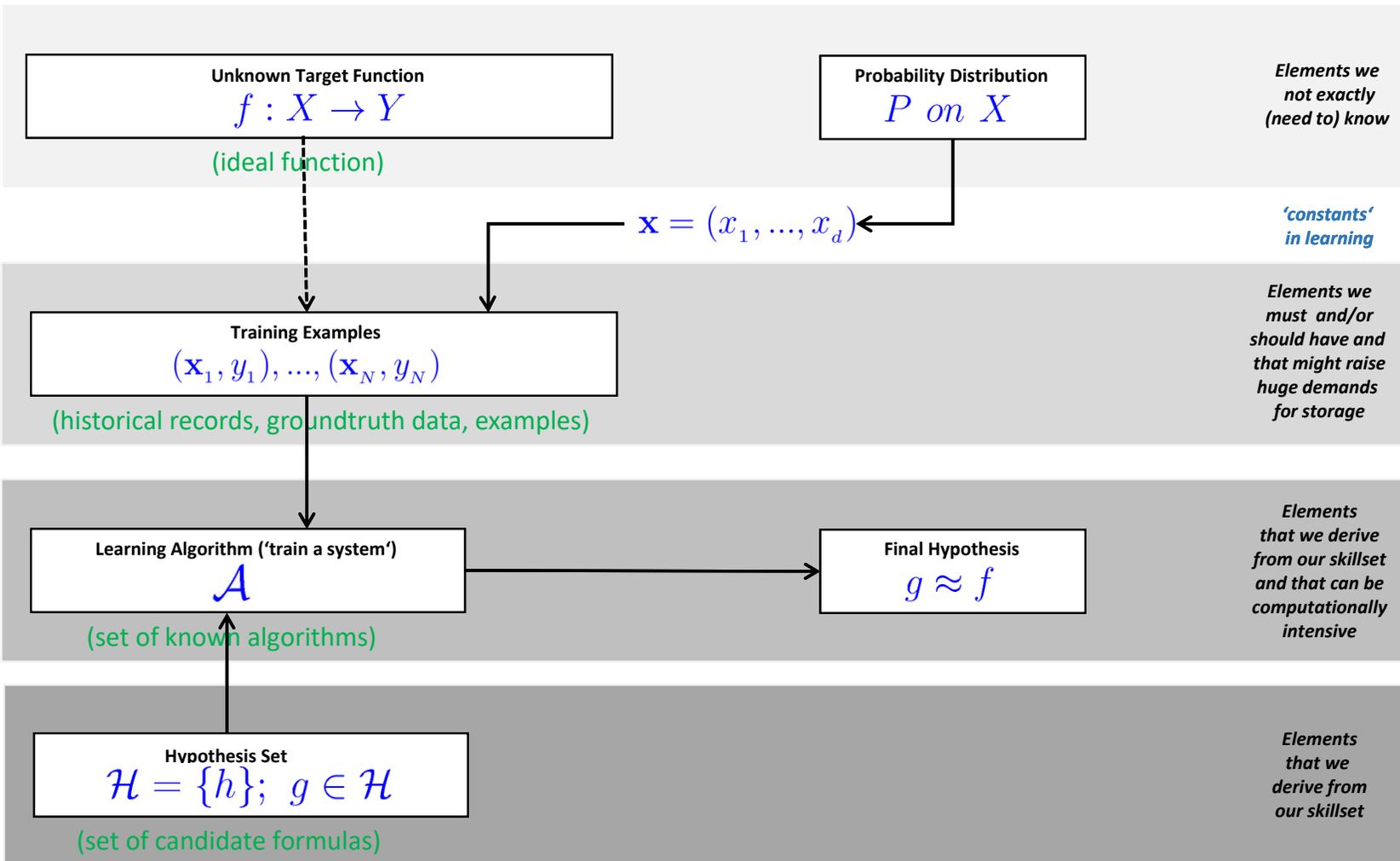
[1] Valiant, ‘A Theory of the Learnable’, 1984

$$\Pr \left[\underset{\text{‘Approximately’}}{\left| E_{in}(g) - E_{out}(g) \right|} > \epsilon \right] \leq \underset{\text{‘Probably’}}{2Me^{-2\epsilon^2 N}}$$

‘Probability that E_{in} deviates from E_{out} by more than the tolerance ϵ is a small quantity depending on M and N ’

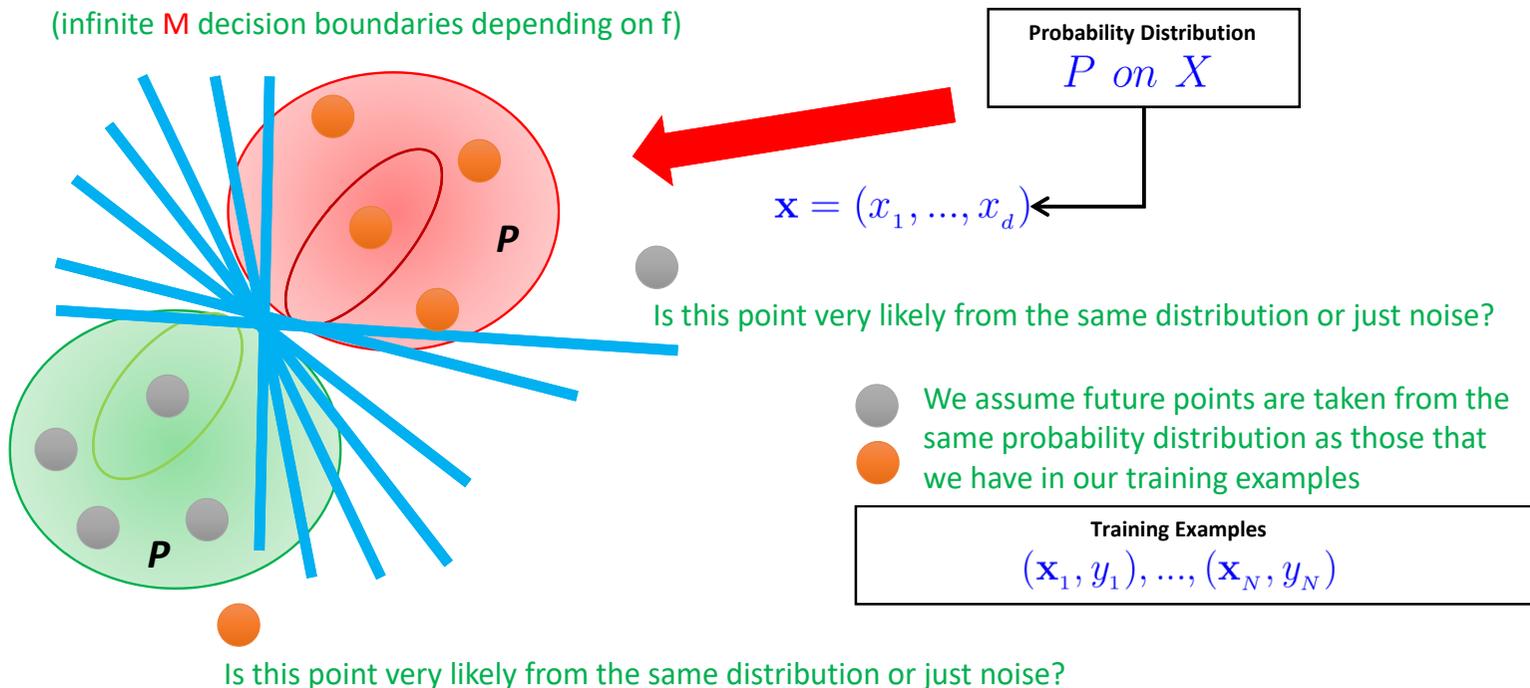
▪ Statistical Learning Theory part describing the Probably Approximately Correct (PAC) learning

- Theoretical ‘Big Data’ Impact \rightarrow more $N \rightarrow$ better learning
 - The more samples N the more reliable will track $E_{in}(g) E_{out}(g)$ well
 - (But: the ‘quality of samples’ also matter, not only the number of samples)
 - For supervised learning also the ‘label’ has a major impact in learning (later)



Mathematical Building Blocks – Our Linear Example

(infinite M decision boundaries depending on f)



(we help here with the assumption for the samples)

(we do not solve the M problem here)

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

(counter example would be for instance a random number generator, impossible to learn this!)

Statistical Learning Theory – Error Measure & Noisy Targets

- Question: How can we learn a function from (noisy) data?
- 'Error measures' to quantify our progress, the goal is: $h \approx f$

- Often user-defined, if not often 'squared error':

$$e(h(\mathbf{x}), f(\mathbf{x})) = (h(\mathbf{x}) - f(\mathbf{x}))^2$$

Error Measure
 α

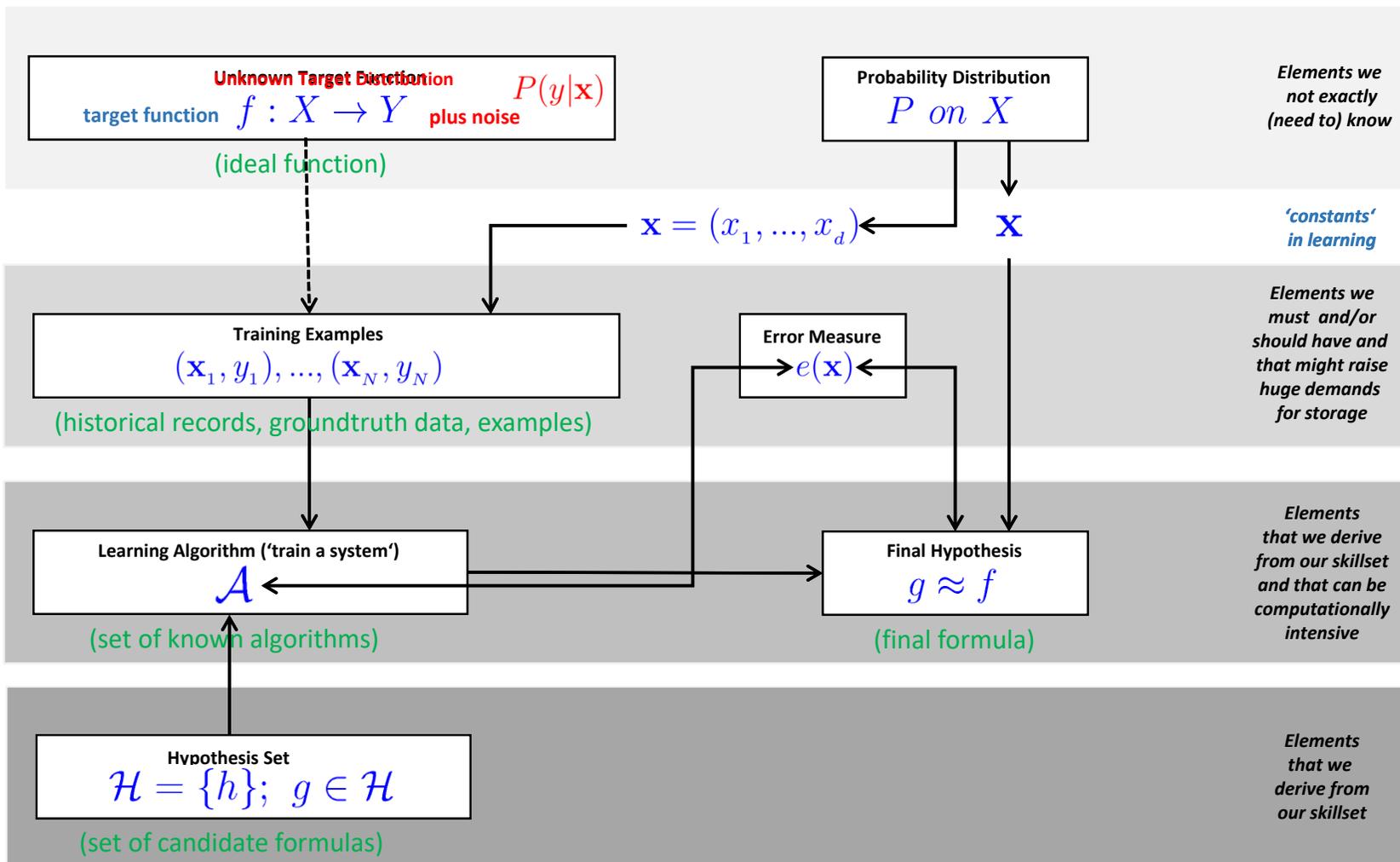
- E.g. 'point-wise error measure'

- '(Noisy) Target function' is not a (deterministic) function (e.g. think movie rated now and in 10 years from now)

- Getting with 'same x in' the 'same y out' is not always given in practice
- Problem: 'Noise' in the data that hinders us from learning
- Idea: Use a 'target distribution' instead of 'target function'
- E.g. credit approval (yes/no)

Unknown Target Distribution $P(y|\mathbf{x})$
target function $f : X \rightarrow Y$ plus noise
(ideal function)

Statistical Learning Theory refines the learning problem of learning an unknown target distribution



Mathematical Building Blocks – Our Linear Example – Error Measures

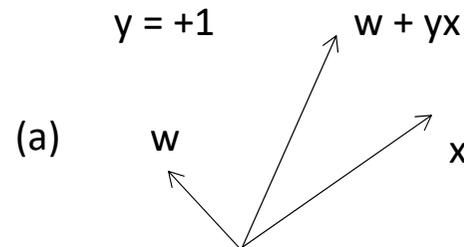
- Iterative Method using (labelled) training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

(one point at a time is picked)

- Pick one misclassified training point where:

$$\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$$

Error Measure
α

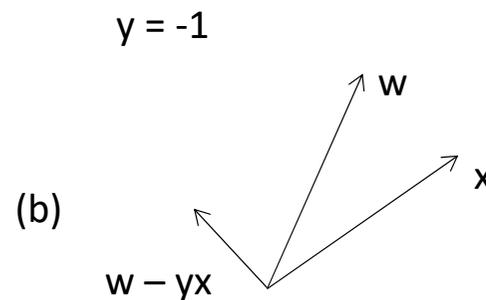


- Update the weight vector:
 - adding a vector or
 - subtracting a vector

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

(y_n is either +1 or -1)

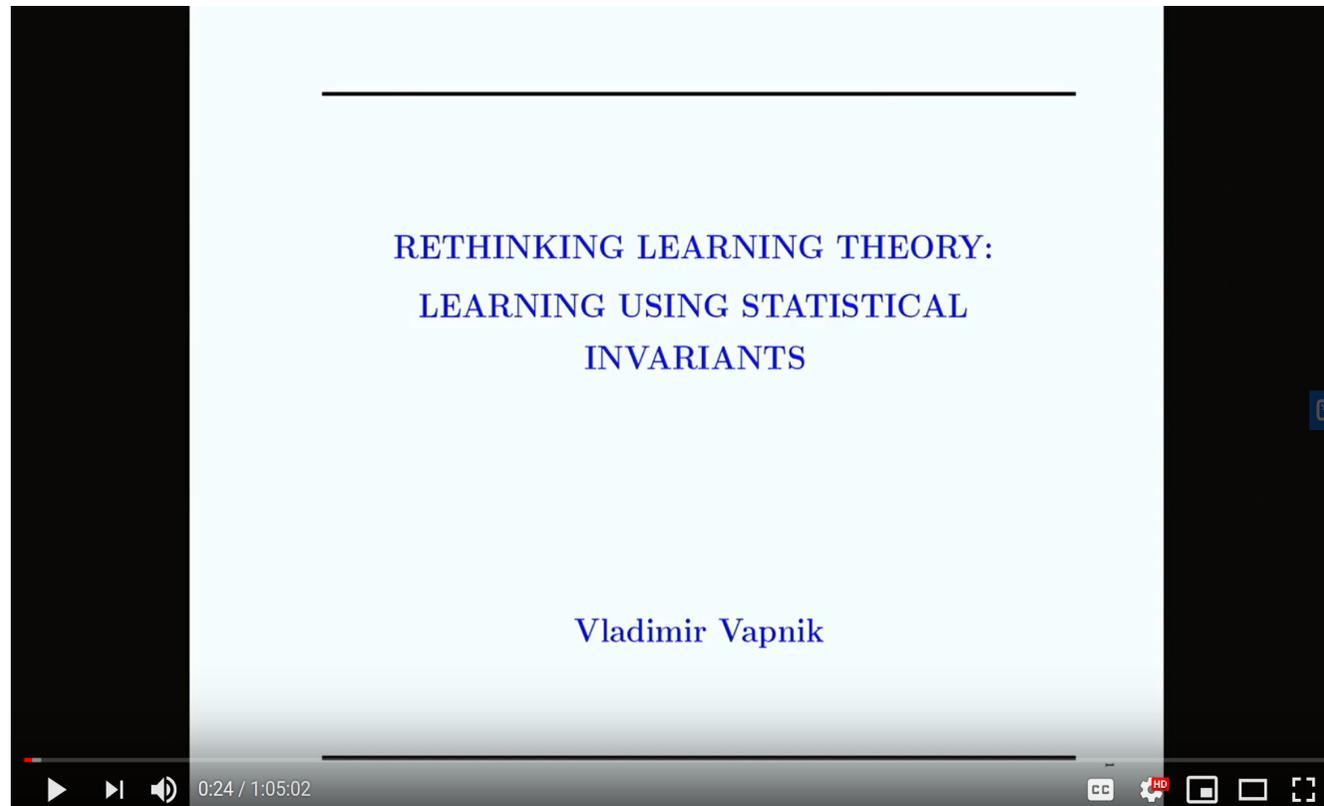
Error Measure
α



- Terminates when there are no misclassified points

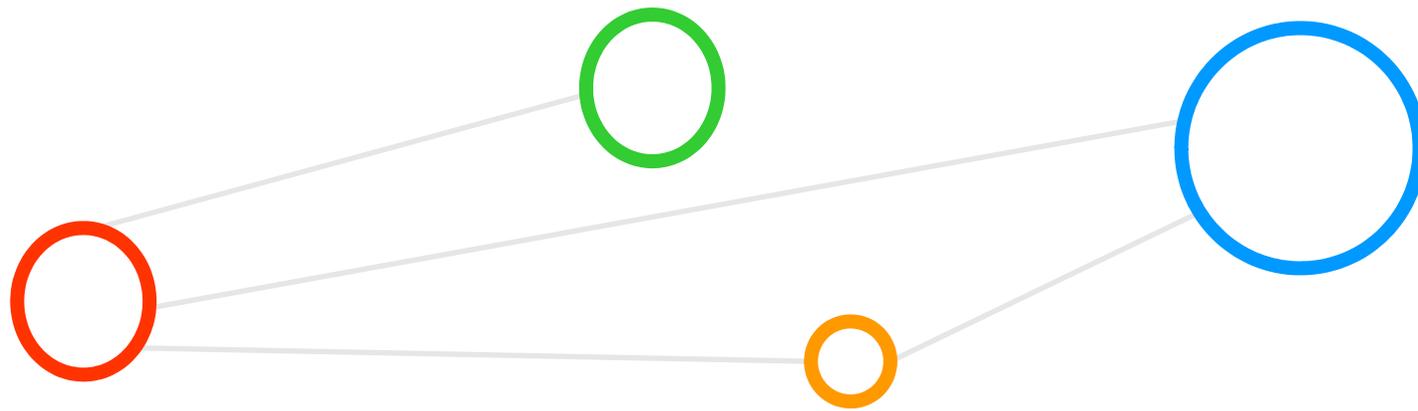
(converges only with linearly separable data)

[Video] Rethinking Statistical Learning by Prof. Vladimir Vapnik (2018)



[7] YouTube Video, Rethinking Learning Theory

Vapnik – Chervonenkis (VC) Inequality & Dimension



Training and Testing – Influence on Learning

- Mathematical notations

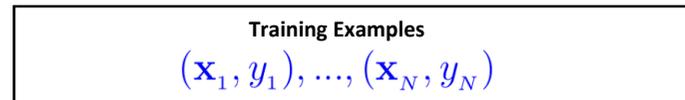
- **Testing** follows: (hypothesis clear) $\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2 e^{-2\epsilon^2 N}$

- **Training** follows: (hypothesis search) $\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$

- Practice on ‘**training examples**’

(e.g. student exam training on examples to get E_{in} ,down’, then test via exam)

- Create **two disjoint datasets**
- One used **for training only** (aka training set)
- Another **used for testing only** (aka test set)



(historical records, groundtruth data, examples)

- Training & Testing are **different phases in the learning process**

- **Concrete number of samples in each set often influences learning**

Theory of Generalization – Initial Generalization & Limits

- Learning is feasible in a probabilistic sense

- Reported final hypothesis – using a ‘generalization window’ on $E_{out}(g)$
- Expecting ‘out of sample performance’ tracks ‘in sample performance’
- Approach: $E_{in}(g)$ acts as a ‘proxy’ for $E_{out}(g)$

$$E_{out}(g) \approx E_{in}(g)$$

This is not full learning – rather ‘good generalization’ since the quantity $E_{out}(g)$ is an unknown quantity

- Reasoning

- Above condition is not the final hypothesis condition:
- More similar like $E_{out}(g)$ approximates 0 (out of sample error is close to 0 if approximating f)
- $E_{out}(g)$ measures how far away the value is from the ‘target function’
- Problematic because $E_{out}(g)$ is an unknown quantity (cannot be used...)
- The learning process thus requires ‘two general core building blocks’

Final Hypothesis

$$g \approx f$$

Theory of Generalization – Learning Process Reviewed

■ ‘Learning Well’

- Two core building blocks that achieve $E_{out}(g)$ approximates 0

■ First core building block

- **Theoretical result** using Hoeffdings Inequality $E_{out}(g) \approx E_{in}(g)$
- Using $E_{out}(g)$ directly is not possible – it is an unknown quantity

■ Second core building block

- **Practical result** using tools & techniques to get $E_{in}(g) \approx 0$
- e.g. **linear models with the Perceptron Learning Algorithm (PLA)**
- Using $E_{in}(g)$ is possible – it is a known quantity – ‘so lets get it small’
- Lessons learned from practice: **in many situations ‘close to 0’ impossible**

(try to get the ‘in-sample’ error lower)

- Full learning means that we can make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$ [from theory]
- Full learning means that we can make sure that $E_{in}(g)$ is small enough [from practical techniques]

Complexity of the Hypothesis Set – Infinite Spaces Problem

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

theory helps to find a way to deal with infinite M hypothesis spaces

Tradeoff & Review

- Tradeoff between ϵ , M , and the ‘complexity of the hypothesis space H ’
- Contribution of detailed learning theory is to ‘understand factor M ’
- M Elements of the hypothesis set \mathcal{H} M elements in H here
 - Ok if N gets big, but problematic if M gets big \rightarrow bound gets meaningless
 - E.g. classification models like perceptron, support vector machines, etc.
 - Challenge:** those classification models have continuous parameters
 - Consequence:** those classification models have infinite hypothesis spaces
 - Approach:** despite their size, the models still have limited expressive power

Many elements of the hypothesis set H have continuous parameter with infinite M hypothesis spaces

Factor **M** from the Union Bound & Hypothesis Overlaps

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq \Pr [| E_{in}(h_1) - E_{out}(h_1) | > \epsilon$$

assumes no overlaps, all probabilities happen disjointly

$$\text{or } | E_{in}(h_2) - E_{out}(h_2) | > \epsilon \dots$$

$$\text{or } | E_{in}(h_M) - E_{out}(h_M) | > \epsilon]$$

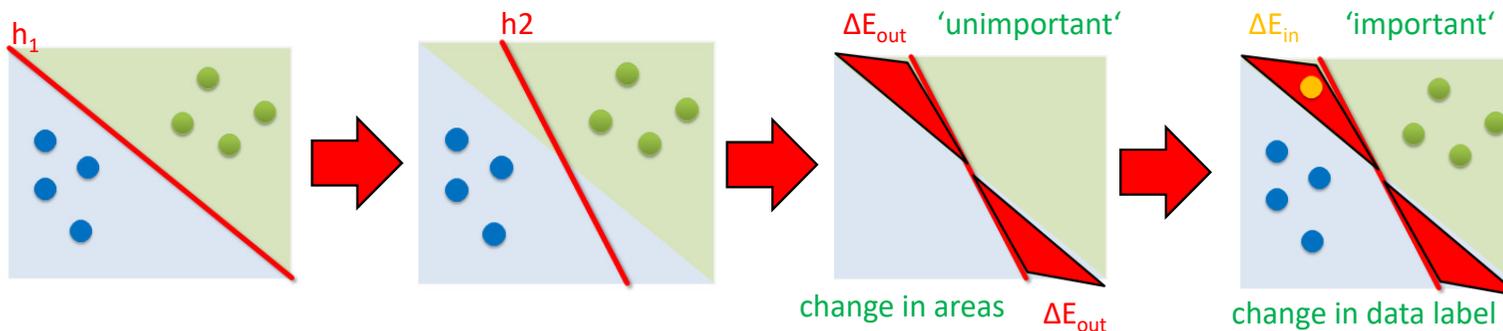
$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

takes no overlaps of **M** hypothesis into account

- Union bound is a ‘poor bound’, ignores correlation between **h**
 - Overlaps are common: the interest is shifted to data points changing label

$$| E_{in}(h_1) - E_{out}(h_1) | \approx | E_{in}(h_2) - E_{out}(h_2) |$$

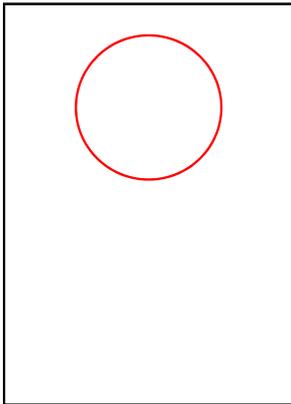
(at least very often, indicator to reduce **M**)



▪ Statistical Learning Theory provides a quantity able to characterize the overlaps for a better bound

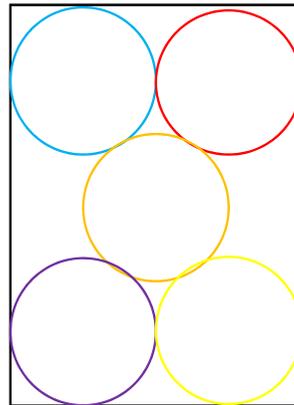
Replacing **M** & Large Overlaps

(Hoeffding Inequality)



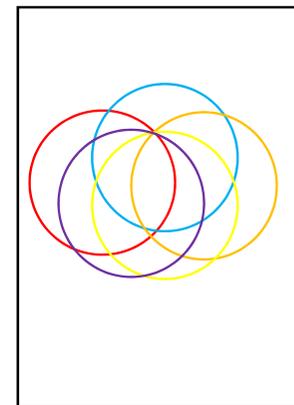
(valid for 1 hypothesis)

(Union Bound)



(valid for M hypothesis, worst case)

(towards Vapnik Chervonenkis Bound)



(valid for m(N) as growth function)

- **Characterizing the overlaps** is the idea of a ‘growth function’
 - **Number of dichotomies:**

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N} |\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$$
 Number of hypothesis but on finite number N of points
 - Much redundancy: Many hypothesis will **reports the same dichotomies**
- The mathematical proofs that $m_{\mathcal{H}}(N)$ can replace M is a key part of the theory of generalization

Complexity of the Hypothesis Set – VC Inequality

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N} |\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$$

■ Vapnik-Chervonenkis (VC) Inequality

- Result of mathematical proof when replacing M with growth function m
- $2N$ of growth function to have another sample ($2 \times E_{in}(h)$, no $E_{out}(h)$)

$$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 4m_{\mathcal{H}}(2N)e^{-1/8\epsilon^2 N}$$

(characterization of generalization)

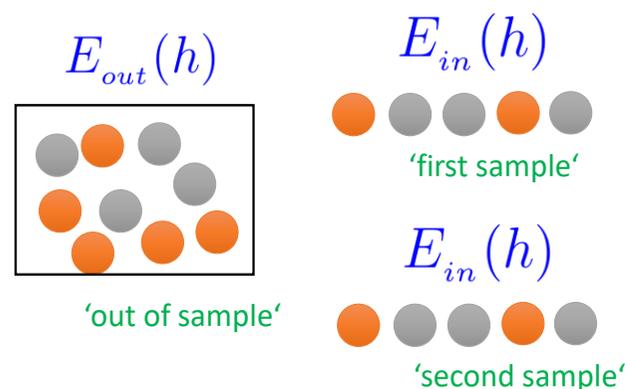
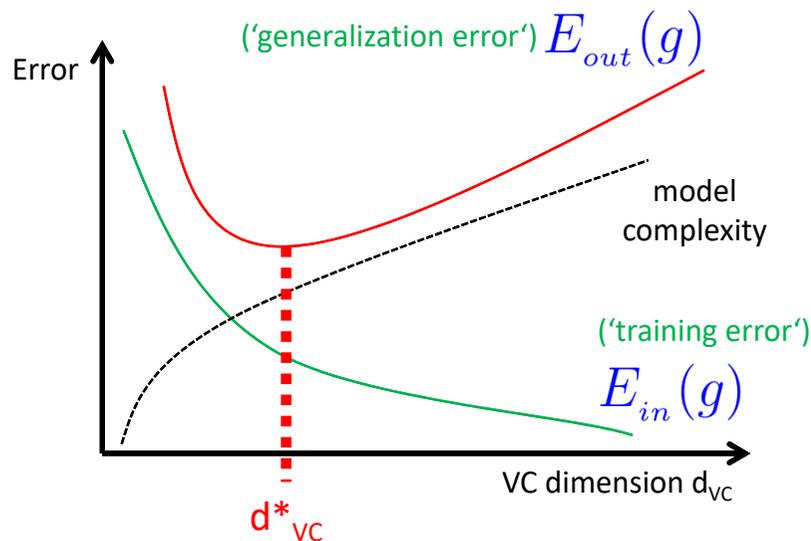
- In Short – finally : We are able to learn and can generalize ‘ouf-of-sample’

- The Vapnik-Chervonenkis Inequality is the most important result in machine learning theory
- The mathematical proof brings us that M can be replaced by growth function (no infinity anymore)
- The growth function is dependent on the amount of data N that we have in a learning problem

Complexity of the Hypothesis Set – VC Dimension & Model Complexity

- Vapnik-Chervonenkis (VC) Dimension over instance space X
 - VC dimension gets a 'generalization bound' on all possible target functions

Issue: unknown to 'compute' – VC solved this using the growth function on different samples



- Complexity of Hypothesis set H can be measured by the Vapnik-Chervonenkis (VC) Dimension d_{VC}
- Ignoring the model complexity d_{VC} leads to situations where $E_{in}(g)$ gets down and $E_{out}(g)$ gets up

Different Models – Hypothesis Set & Model Capacity

Hypothesis Set

$$\mathcal{H} = \{h\}; g \in \mathcal{H}$$

$$\mathcal{H} = \{h_1, \dots, h_m\};$$

(all candidate functions
derived from models
and their parameters)

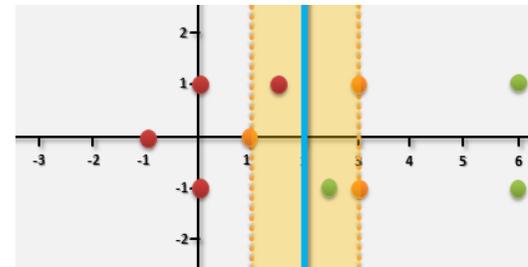
- Choosing from various model approaches h_1, \dots, h_m is a different hypothesis
- Additionally a change in model parameters of h_1, \dots, h_m means a different hypothesis too
- The model capacity characterized by the VC Dimension helps in choosing models
- Occam's Razor rule of thumb: 'simpler model better' in any learning problem, not too simple!

'select one function'
that best approximates

Final Hypothesis

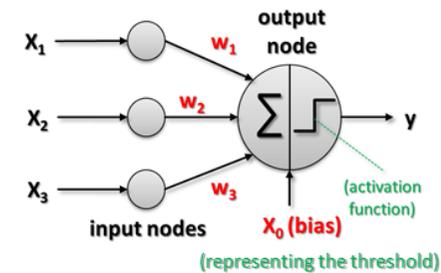
$$g \approx f$$

h_1



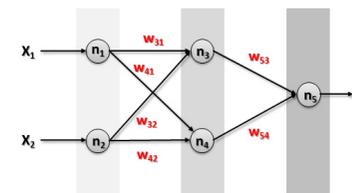
(e.g. support vector machine model)

h_2



(e.g. linear perceptron model)

h_m



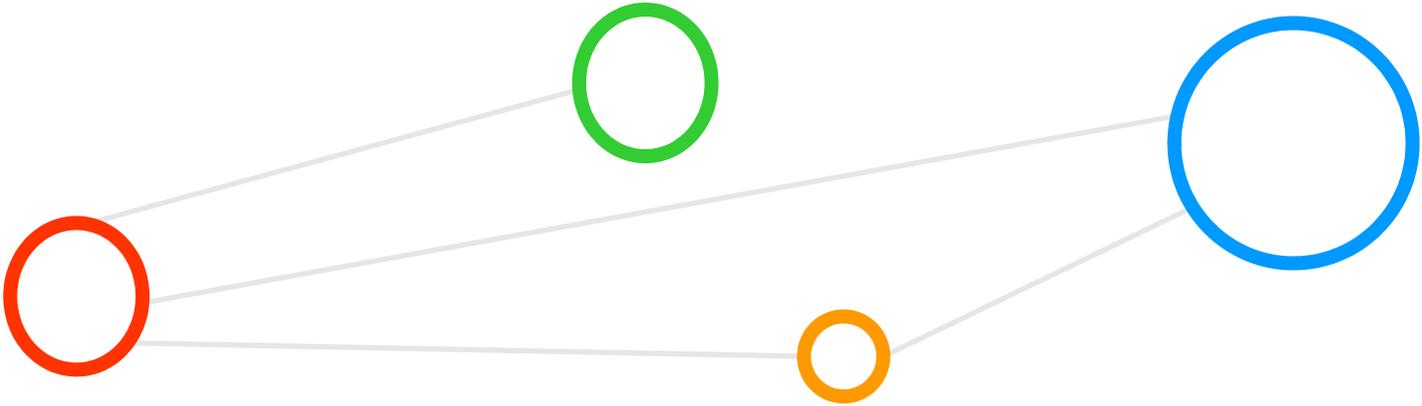
(e.g. artificial neural network model)

[Video] Prevent Overfitting for better Generalization



[6] YouTube Video, Stop Overfitting

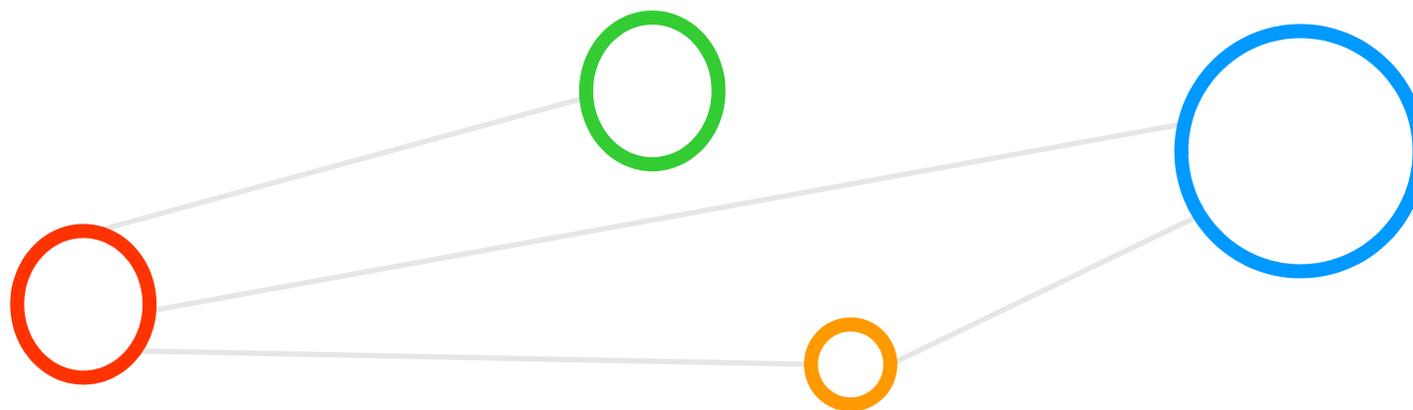
Lecture Bibliography



Lecture Bibliography

- [1] Multilayer Perceptron, Online:
https://www.eecs.yorku.ca/course_archive/2012-13/F/4404-5327/lectures/10%20Multilayer%20Perceptrons.pdf
- [2] MIT 6.S191: Introduction to Deep Learning, Online:
<http://www.introtodeeplearning.com>
- [3] Species Iris Group of North America Database, Online:
<http://www.signa.org>
- [4] Cheng, A.C, Lin, C.H., Juan, D.C., InstaNAS: Instance-aware Neural Architecture Search, Online:
<https://arxiv.org/abs/1811.10201>
- [5] Evgeniou T., Pontil M., Poggio T., ‘Statistical Learning Theory: A Primer’, International Journal of Computer Vision 38(1), 9–13, 2000, Online:
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.92.9504&rep=rep1&type=pdf>
- [6] YouTube Video, Udacity, ‘Overfitting’, Online:
<https://www.youtube.com/watch?v=CxAxRCv9WoA>
- [7] YouTube Video, ‘Rethinking Learning Theory’, Online:
<https://www.youtube.com/watch?v=LEYglxKclo&t=2921s>

Acknowledgements



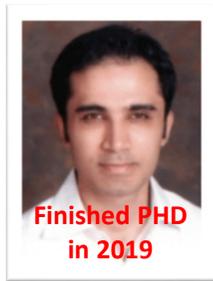
Acknowledgements – High Productivity Data Processing Research Group



PD Dr.
G. Cavallaro



Senior PhD
Student A.S. Memon



Senior PhD
Student M.S. Memon



PhD Student
E. Erlingsson



PhD Student
S. Bakarat



PhD Student
R. Sedona



Dr. M. Goetz
(now KIT)



MSc M.
Richerzhagen
(now other division)



MSc
P. Glock
(now INM-1)



MSc
C. Bodenstein
(now Soccerwatch.tv)



MSc Student
G.S. Guðmundsson
(Landsverkjun)



This research group has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 763558 (DEEP-EST EU Project)

