# *HIGGS BOSON CHALLENGE*

## *FANTASTIC FOUR*
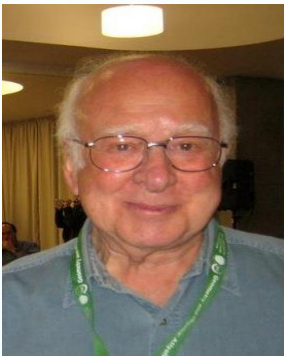
- AKHILA IYENGAR
- LEONARDO CITRARO
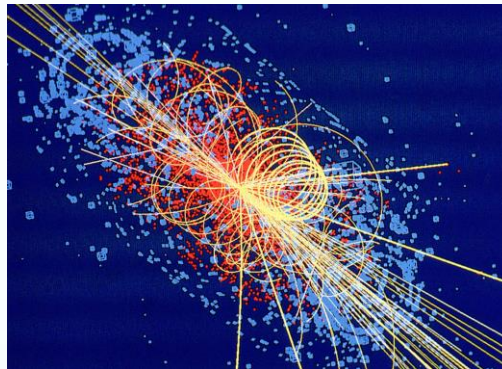- ARTHUR PERODOU
- BAUDOUIN ROULLIER

**LIST OF CONTENTS**

- Introduction
- Data Description
- Evaluation
- Support Vector Machines
- Neural Networks
- Deep Neural Networks
- Results

# Introduction

- Higgs boson is an elementary particle in the standard model of particle physics.

- Peter Higgs and Francois Englert discovered Higgs which was later confirmed by ATLAS and CMS at CERN



**Peter Higgs**

**Francois Englert**

- Small signal buried in huge background noise

- Goal is to classify the events into respective regions (Signal and Background)

# Data Description

PRI_lep_pt

PRI_tau_p

- Dataset comprises of simulated events of which **250,000 are training** set and **550,000 are test** set

- Training data has several attributes like event ID along with **30 features**, labels and weights while test data has event ID and 30 feature columns

- Data has signal (Higgs Boson) and background with associated **weights**

**Binary classification problem**

**Training dataset
30 features
250'000 samples + weights**

**Test dataset
30 features
550'000 samples + weights**

3

# **Evaluation**

- Performance evaluation of predictive model is done by AMS (Approximate Median Significance).

- This makes events with higher weights more significant.

- Let s and b be the true and false positive rates respectively, the AMS is defined as:

$$AMS = \sqrt{2\left(\left(s + b + b_{reg}\right)\ln\left(1 + \frac{s}{b+b_{reg}}\right) - s\right)}$$
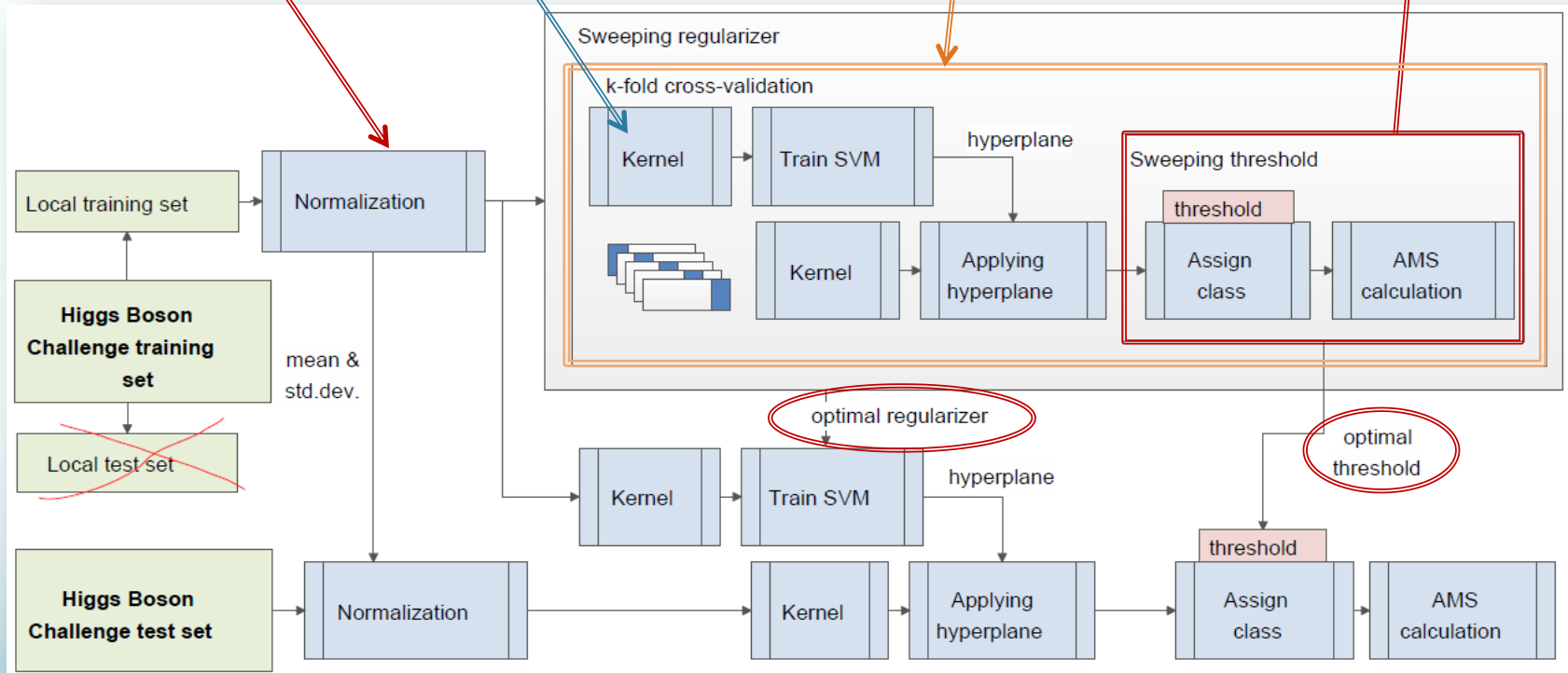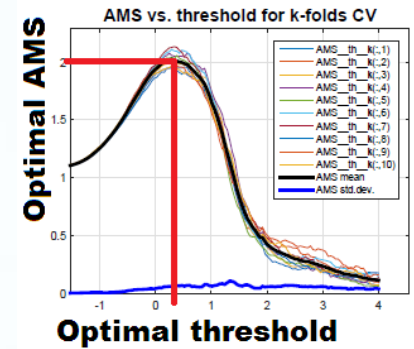
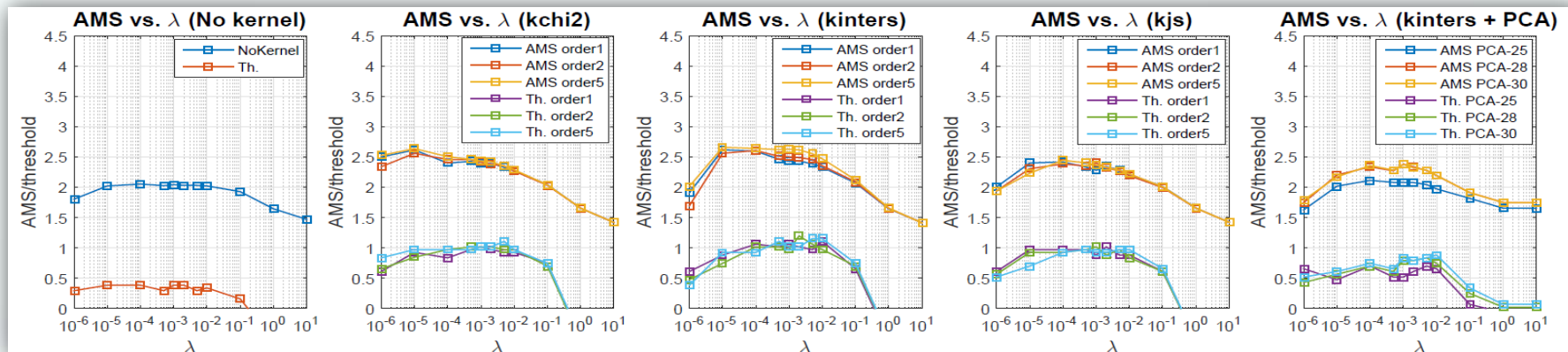$$b_{reg} = 10$$

# Support Vector Machine



**Zero mean unit std.**

**Homogeneous kernel map (linear approximation of non-linear mappings) [VLFeat toolbox]**

**4-folds cross-validation**

AMS vs. threshold for k-folds CV

Optimal AMS

Optimal threshold

Sweeping regularizer

k-fold cross-validation

Kernel → Train SVM → hyperplane

Kernel → Applying hyperplane

Sweeping threshold

threshold

Assign class → AMS calculation

Local training set → Normalization

Higgs Boson Challenge training set

mean & std.dev.

Local test set

optimal regularizer

optimal threshold

Higgs Boson Challenge test set → Normalization → Kernel → Train SVM → hyperplane → Kernel → Applying hyperplane
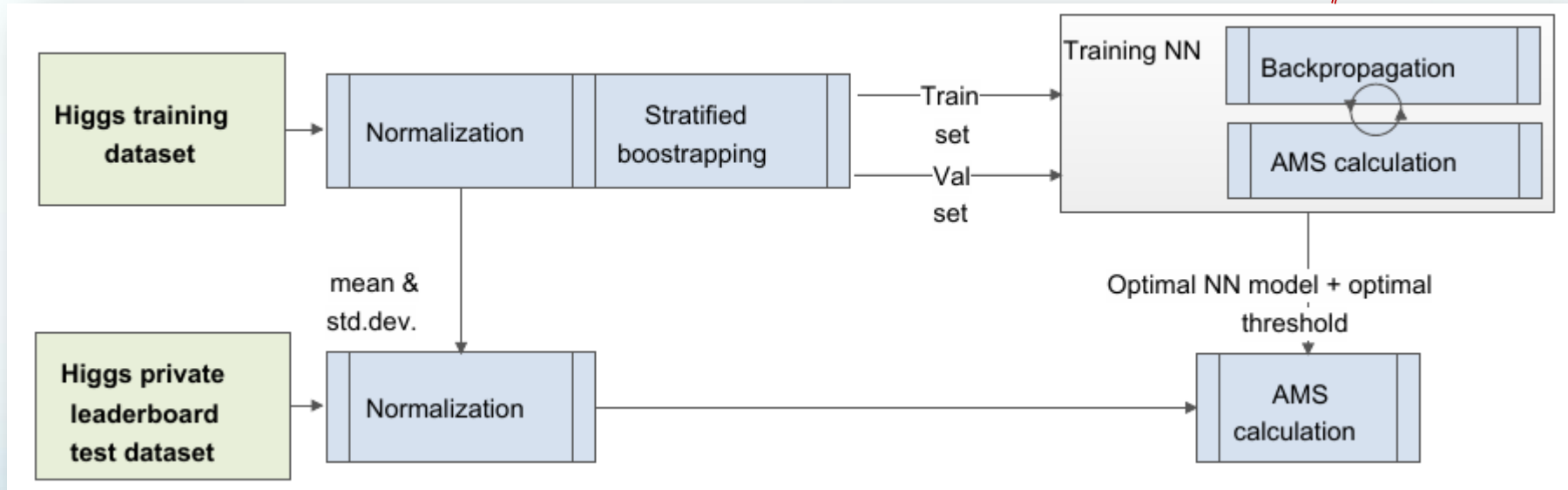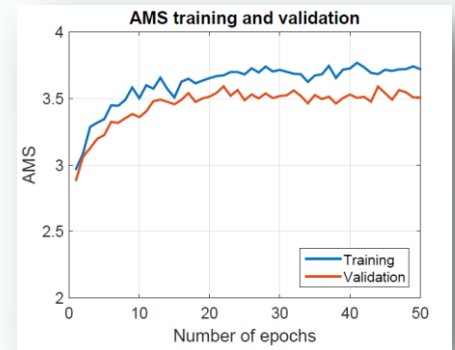
threshold

Assign class → AMS calculation

5

# Results Support Vector Machine
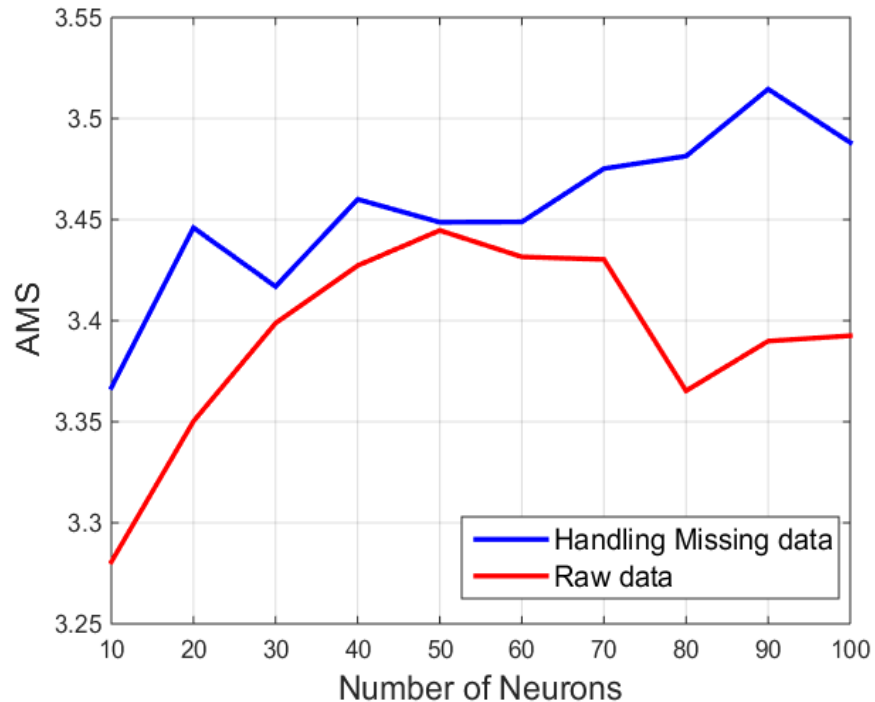
- Maximum AMS vs. Regularizer λ



- Maximum AMS scored in cross-validation: **2.66 (kernel intersect)**

- AMS scored using the test set: **2.44**

- Result far from the best AMS probably due to the linear kernel.

- For big data the SVM is slow and non-linear kernel require a huge amount of memory
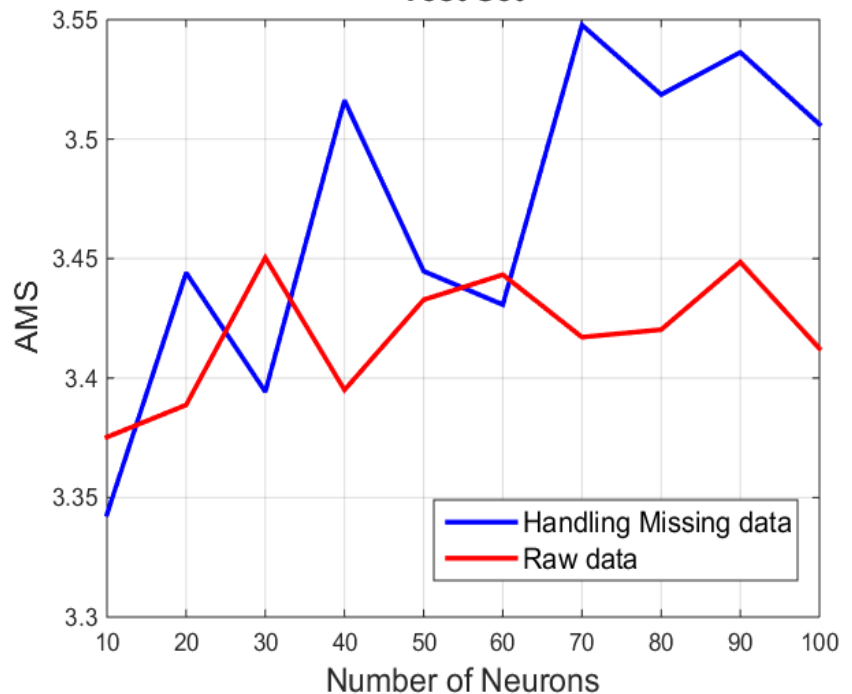
# Single Layer Neural Network
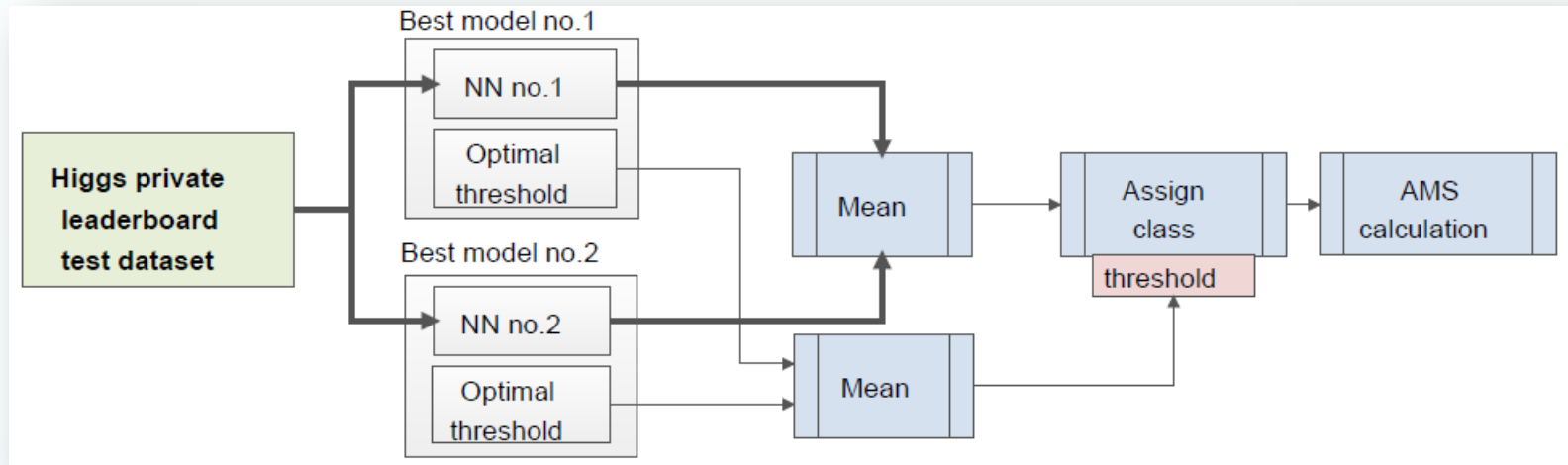
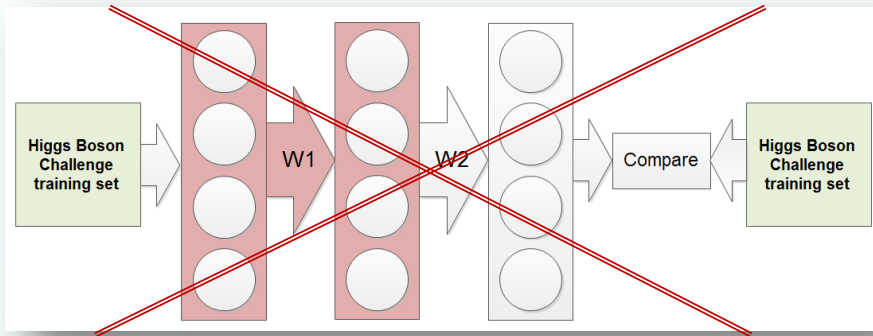# Result Single Layer Neural Network + Missing Data
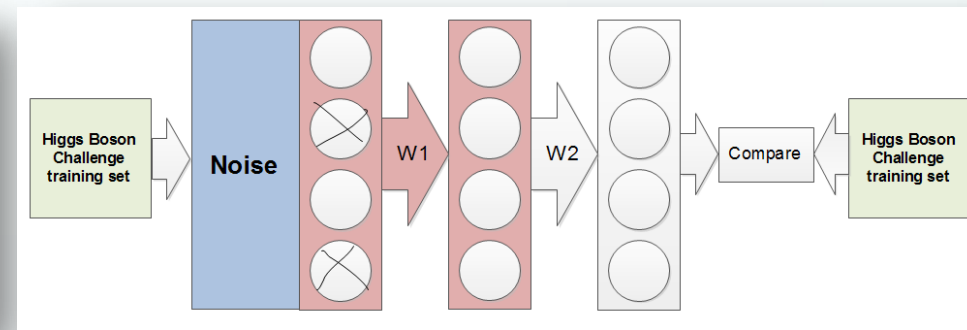
# Averaging multiple Neural Networks



- Using k NN

- Averaging of their prediction

- Averaging of their optimal threshold

# Deep Neural Network

## Auto-encoder

| Higgs Boson Challenge training set | → | W1 | → | W2 | → | Compare | ← | Higgs Boson Challenge training set |

## Denoising Auto-encoder

| Higgs Boson Challenge training set | → | Noise | W1 | → | W2 | → | Compare | ← | Higgs Boson Challenge training set |

**Thousand of ways to set W1 and W2 giving error=0**

Train ↻ Fix layer

**more meaningful mapping of the dataset**

Greedy layer-wise unsupervised pre-trained process

Input layer   Hidden layer   Hidden layer   Hidden layer
W    W    W    W    Output layer

**← Deep network!**

10

# Multiple Averaged Deep Neural Network + Missing Data

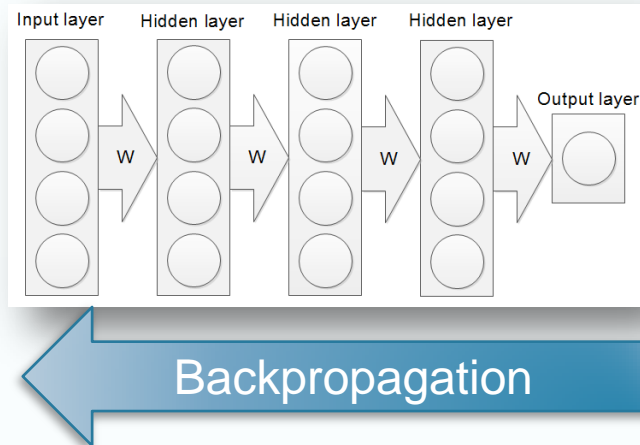| Network type/results | Number of NN averaged | Test AMS | | | Averaged result Test AMS |
|---|---|---|---|---|---|
| Deep NN [30 80 80 80 1] with stacked AE denoising 10% | 5 | 3.54   3.48   3.56<br>3.48   3.54 | | | 3.60 |

**Architecture**

**Noise percentage Auto-Encoders**

**Number NN trained**

**Results of each DeepNN Mean:3.53**

**Result of averaging process**

# Error weights

**Classic error function**



Backpropagation

$$E_n = \frac{1}{2} \sum_k (\widehat{y_{nk}} - t_{nk})^2$$

**Weighted error function**



Backpropagation

$$E_n = \frac{1}{2} \sum_k (\widehat{y_{nk}} - t_{nk})^2 \cdot \boldsymbol{w_n}$$

**Set according to the weights of the training dataset**

**Forces the network to learn where matter**

12

# Multiple Averaged Deep Neural Network + Missing Data + Error weights

| Network type/results | Number of networks averaged | Test AMS | Averaged result Test AMS |
|---|---|---|---|
| Deep NN [30 80 80 80 1] AE denoising 10% | 5 | 3.67 3.67 3.62 3.64 3.72 | 3.75 |

**Same architecture as before**

**Number NN**

**Results of each DeepNN**

**Result of averaging process!!!!**

# Results



Averaged DNN + missing data + error weights — **3.75**

Averaged DNN + missing data — **3.64**

DNN + missing data — **3.53**

Averaged Single layer NN + missing data — **3.49**

Single layer NN — **3.39**

SVM + homogeneous kernel — **2.44**

Completed • $13,000 • 1,785 teams

## Higgs Boson Machine Learning Challenge

Higgs challenge

kaggle

Mon 12 May 2014 – Mon 15 Sep 2014 (6 months ago)

| # | Δrank | Team Name ‡ model uploaded * in the money | Score | Entries | Last Submission UTC (Best – Last Submission) |
|---|---|---|---|---|---|
| 1 | ↑1 | Gábor Melis ‡ * | 3.80581 | 110 | Sun, 14 Sep 2014 09:10:04 (-0h) |
| 2 | ↑1 | Tim Salimans ‡ * | 3.78913 | 57 | Mon, 15 Sep 2014 23:49:02 (-40.6d) |
| 3 | ↑1 | nhlx5haze ‡ * | 3.78682 | 254 | Mon, 15 Sep 2014 16:50:01 (-76.3d) |

**Our rank: 14th over ~1800 teams** ➡ **top 8% in 2 month**

# Questions?