# Introduction to Machine Learning in Astroparticle Physics

Tim Ruhe, TU Dortmund University

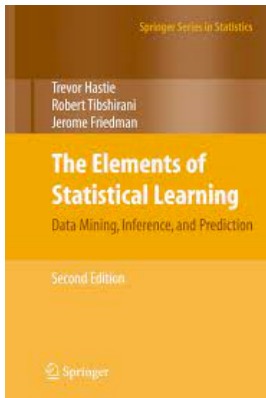PHYSTAT-Gamma 09/27/22

tim.ruhe@tu-dortmund.de

IceCube

## Motivation

Machine Learning provides tools to accomplish an analysis task faster and more accurately (when used correctly).

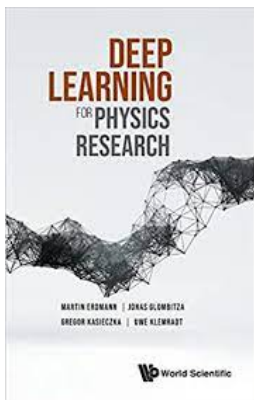# Further Reading

- General Introduction to Statistical Learning
- Good start to get an overview
- A lot of extra material: https://hastie.su.domains/ElemStatLearn/
- (I believe you can also download the pdf there....)

Source: https://hastie.su.domains/ElemStatLearn/

- Focus on Deep Learning and Neural Networks
- Nice pedagogic approach

Source: amazon

# Further Reading

**DE GRUYTER**

**MACHINE LEARNING UNDER RESOURCE CONSTRAINTS**

**DISCOVERY IN PHYSICS**

*Edited by Katharina Morik and Wolfgang Rhode*

- Focus on astroparticle and particle physics
- Contains a lot of topics also covered in this talk
- Open access
- To be published by the end of 2022
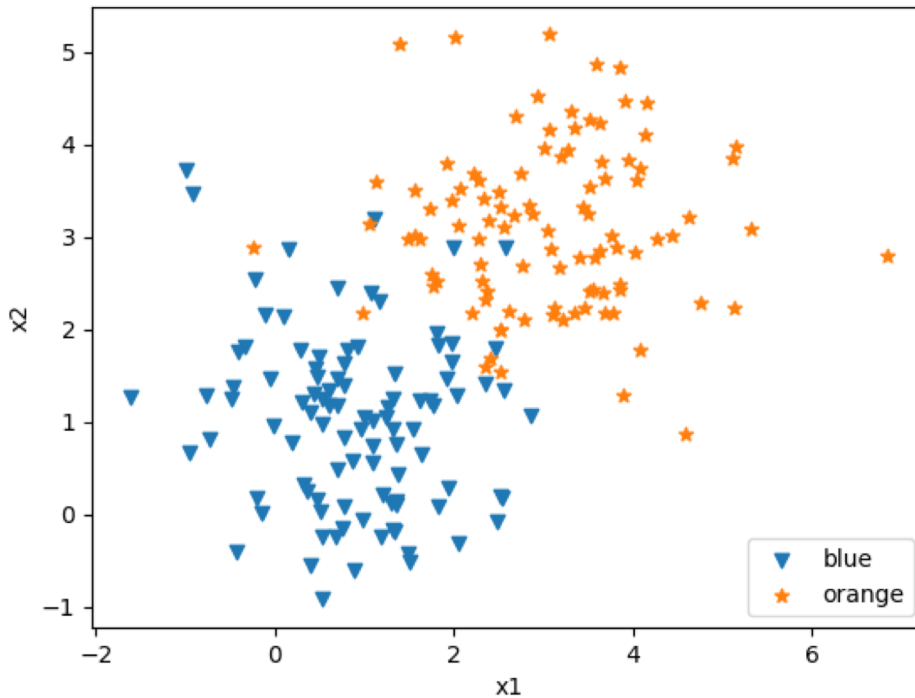
Source: Von The scikit-learn developers - github.com/scikit-learn/scikit-learn/blob/master/doc/logos/scikit-learn-logo.svg, BSD, https://commons.wikimedia.org/w/index.php?curid=71445288

# Outline

- Nomenclature and stuff
- Feature Selection
- Selected Algorithms
- Neural Networks and Deep Learning

## Nomenclature



N $(\vec{X}, y)$ pairs are referred to as training set
Or annotated data

Events (Examples) are characterized by a feature vector:

$$\vec{X} = (x_1 \ \dots \ x_n)$$

In this example

$$\vec{X} = (x_1, x_2)$$

And a class variable

$$y \ \in \ [y_1 \ \dots \ y_n]$$

In this example

$$y \ \in [blue, orange]$$

## Nomenclature

Task: Build a model to separate blue and orange examples.



Events (Examples) are characterized by a feature vector:

$$\vec{X} = (x_1 \ \dots \ x_n)$$

In this example

$$\vec{X} = (x_1, x_2)$$

And a class variable

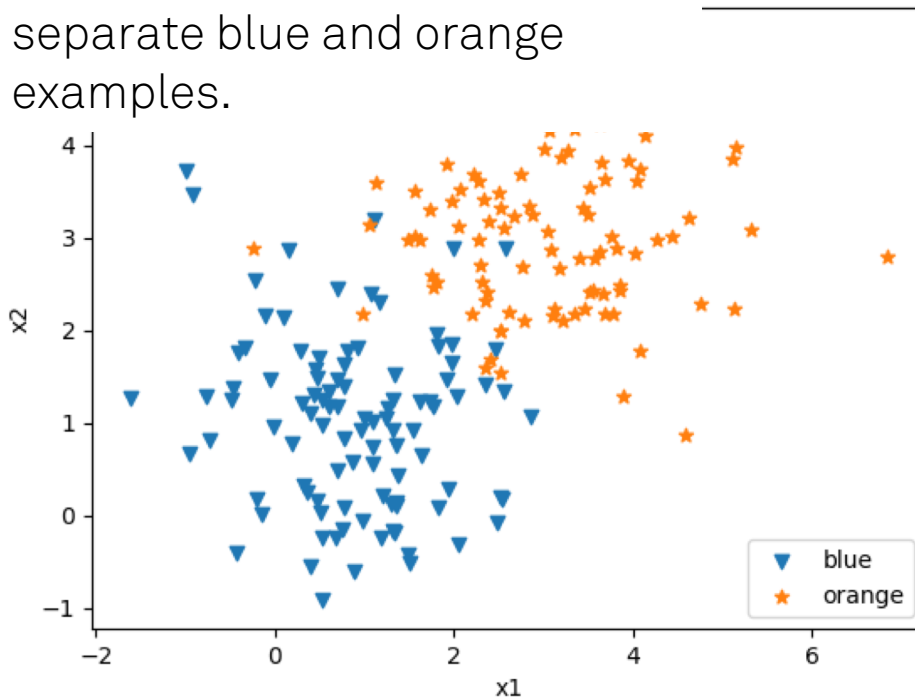$$y \in [y_1 \ \dots \ y_n]$$

In this example

$$y \in [blue, orange]$$

# The Linear Model



$$\hat{y} = \beta_0 + \sum_{i=1}^{p} x_i \beta_i$$

$$\hat{y} = \begin{cases} \text{orange: } 0 \\ \text{blue: } 1 \end{cases}$$

Solve e.g. by least squares fit

# The Linear Model: Graphical Representation of the Model

The model is not able to perfectly describe the training data.

Above line: Classify as orange

Below line: Classify as blue

# Application to Unseen Data



Confusion Matrix

|  | Blue (+) | Orange (-) |
|---|---|---|
| **Blue (+)** | 22 | 2 |
| **Orange (-)** | 3 | 22 |

*I defined that border will be part of the orange class.

# True and False Negatives and Postives

True Positives (TP)

False Positives (FP)

Confusion Matrix

|  | Blue (+) | Orange (-) |
|---|---|---|
| Blue (+) | 22 | 2 |
| Orange (-) | 3 | 22 |

False Negatives (FN)

True Negatives (TN)

blue: 1
orange: 0

*I defined that border will be part of the orange class.

# Quality Measures

Accuracy:

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + FP + TN + FN}$$

Precision:

$$PREC = \frac{TP}{TP + FP}$$

Recall:

$$REC = \frac{TP}{TP + FP}$$

\* These measures can sometimes have different names

True Positives (TP)

False Positives (FP)

Confusion Matrix

|  | Blue (+) | Orange (-) |
|---|---|---|
| **Blue (+)** | 22 | 2 |
| **Orange (-)** | 3 | 22 |

False Negatives (FN)

True Negatives (TN)

# Area Under Curve

Graphics: M. Linhoff [Learning Under Resource Constraints – Discovery in Physics] (in preparation)



Mean area under curve: $0.8263 \pm 0.0022$

Single CV ROC Curve
Mean ROC curve

Graphic: M. Linhoff

Source: By cmglee, MartinThoma - Roc-draft-xkcd-style.svg, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=109730045

ROC characteristic for the FACT Open Crab data set

# Exemplary Data Analysis Pipeline

| Variable Selection | → | Classifier Training | → | Cut on Classifier Output |
|---|---|---|---|---|







Picture: CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=14260

Source: https://www.pinterest.com/pin/550354016946043419/

# Variable Selection: Try possible combinations



Yes, but...

Source: By Thore Husfeldt at English Wikipedia, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=31823619

# Feature Selection

| | |
|---|---|
| Initially | 1219 |
| blacklisted | 1129 |
| constant & useless | 855 |
| Correlation cut | 323 |
| Data/MC Clf | 311 |
| mRMR | 60 |

M. Börner, PhD thesis (2018)

Exclude features that either bias the selection or are only present in simulation.



After correlation cut:
323 Features

Before $mRMR$:
311 Features

Final Sample:
61 Features

Data
Simulations

Frequence

Classification Score

True Positive Rate

# Feature Selection

| Initially |
|:---:|
| 1219 |

| blacklisted |
|:---:|
| 1129 |

| constant & useless |
|:---:|
| 855 |

| Correlation cut |
|:---:|
| 323 |

| Data/MC Clf |
|:---:|
| 311 |

| mRMR |
|:---:|
| 60 |

M. Börner, PhD thesis (2018)

T. Ruhe, PHYSTAT-Gamma

Exclude features that either bias the
selection or are only present in simulation.

Constant features do not carry information.



After correlation cut:
323 Features

Data
Simulations

Before $mRMR$:
311 Features

Final Sample:
61 Features

Classification Score

True Positive Rate

Fa

# Feature Selection

| Initially |
|---|
| 1219 |

| blacklisted |
|---|
| 1129 |

| constant & useless |
|---|
| 855 |

| Correlation cut |
|---|
| 323 |

| Data/MC Clf |
|---|
| 311 |

| mRMR |
|---|
| 60 |

M. Börner, PhD thesis (2018)

Exclude features that either bias the selection or are only present in simulation.

Constant features do not carry information.

Strongly correlated features do not contain new information (or only very little)



After correlation cut: 323 Features

Before *mRMR*: 311 Features

Final Sample: 61 Features

## Feature Selection

| Initially |
|:---:|
| 1219 |

| blacklisted |
|:---:|
| 1129 |

| constant & useless |
|:---:|
| 855 |

| Correlation cut |
|:---:|
| 323 |

| Data/MC Clf |
|:---:|
| 311 |

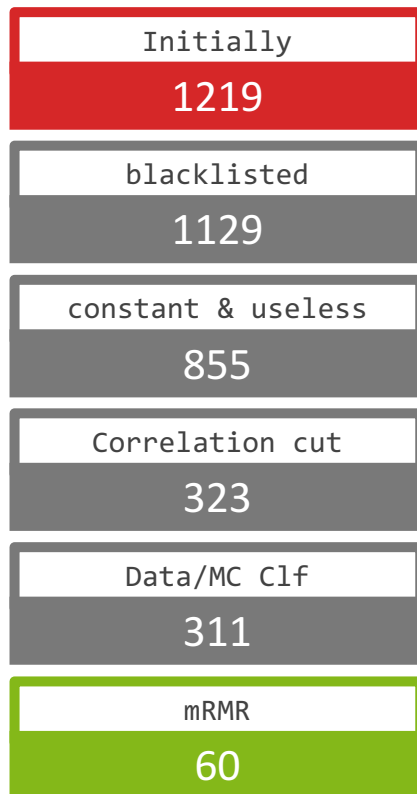| mRMR |
|:---:|
| 60 |

M. Börner, PhD thesis (2018)

T. Ruhe, PHYSTAT-Gamma

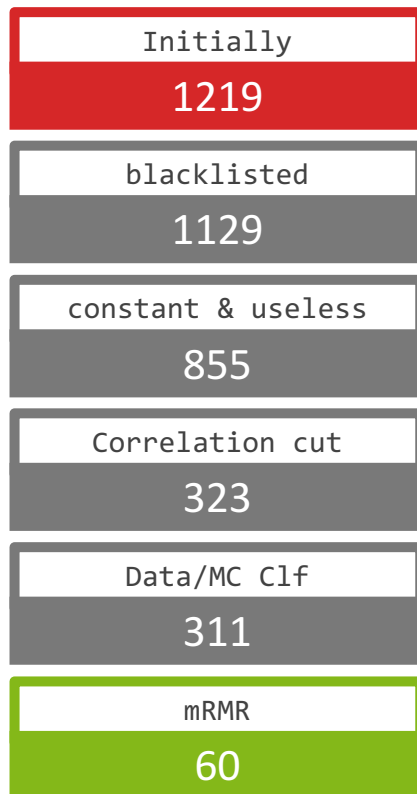Exclude features that either bias the selection or are only present in simulation.

Constant features do not carry information.

Strongly correlated features do not contain new information (or only very little)

Simulated and experimental data should agree to not bias the result.



After correlation cut: 323 Features

Before *mRMR*: 311 Features

Final Sample: 61 Features

Data
Simulations

Frequence

Classification Score

True Positive Rate

# Feature Selection

| Initially |
|:---:|
| **1219** |

| blacklisted |
|:---:|
| **1129** |

| constant & useless |
|:---:|
| **855** |

| Correlation cut |
|:---:|
| **323** |

| Data/MC Clf |
|:---:|
| **311** |

| mRMR |
|:---:|
| **60** |

M. Börner, PhD thesis (2018)

T. Ruhe, PHYSTAT-Gamma

Exclude features that either bias the selection or are only present in simulation.
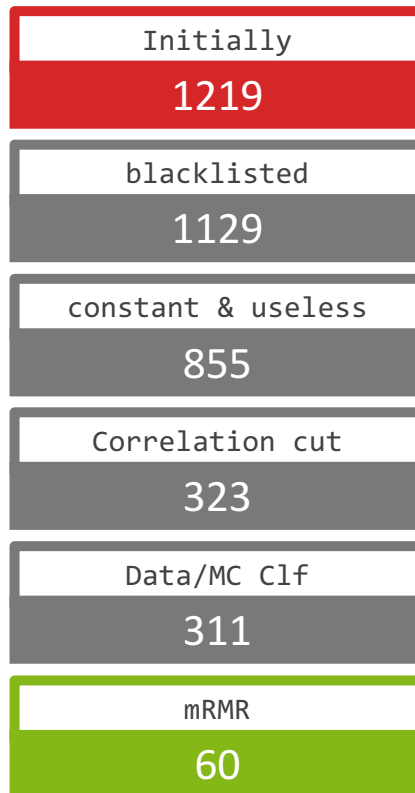
Constant features do not carry information.

Strongly correlated features do not contain new information (or only very little)

Simulated and experimental data should agree to not bias the result.

Automated selection by a feature selection algorithm.



After correlation cut: 323 Features

Before *mRMR*: 311 Features

Final Sample: 61 Features

Data
Simulations

Frequence

Classification Score

True Positive Rate

technische universität
dortmund

technische universität
dortmund

lehrstuhl
physik e5

lehrstuhl
physik e5

# Minimum Redundancy Maximum Relevance

| Initially |
|---|
| 1219 |

| blacklisted |
|---|
| 1129 |

| constant & useless |
|---|
| 855 |

| Correlation cut |
|---|
| 323 |

| Data/MC Clf |
|---|
| 311 |

| mRMR |
|---|
| 60 |

M. Börner, PhD thesis (2018)

- Select features according to relevance and redundancy

- Feature set is built by iteratively adding features that fulfill the following criterion

$$\max_{x_j \in X - S_{m-1}} \left[ I(x_j, c) - \frac{1}{m-1} \sum_{x_i \in S_{m-1}} I(x_i, x_j) \right]$$

Peng, H.C., Long, F., and Ding, C., IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 8, pp. 1226–1238, 2005.
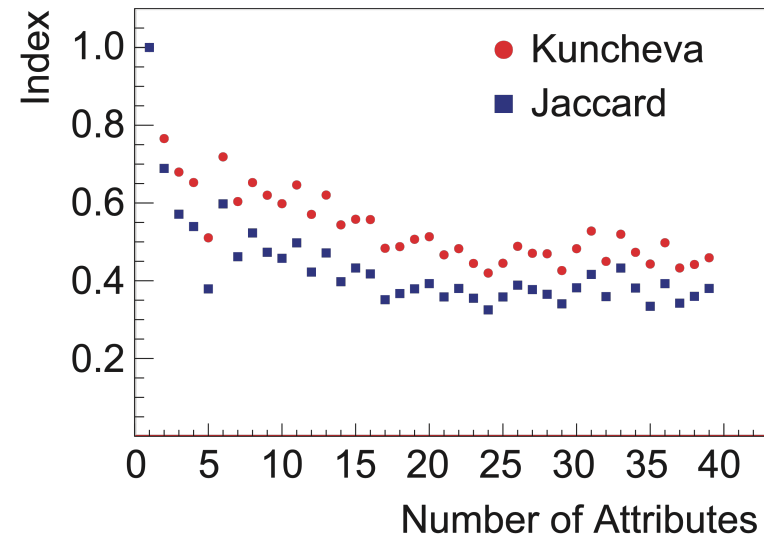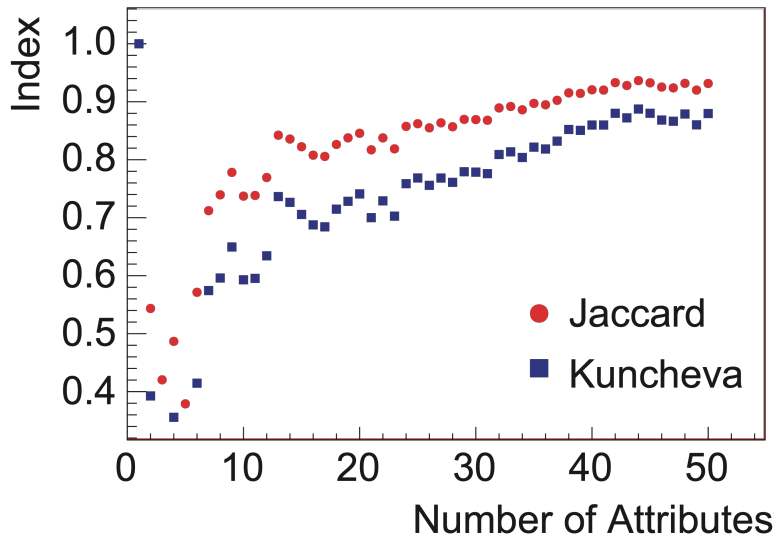
Ding, C., & Peng, H. *Journal of bioinformatics and computational biology*, 3(02), 185-205. (2005)

# Feature Selection Stability

Ludmila I. Kuncheva: A STABILITY INDEX FOR FEATURE SELECTION



$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

$$I_C(A,B) = \frac{rn - k^2}{k(n-k)}$$
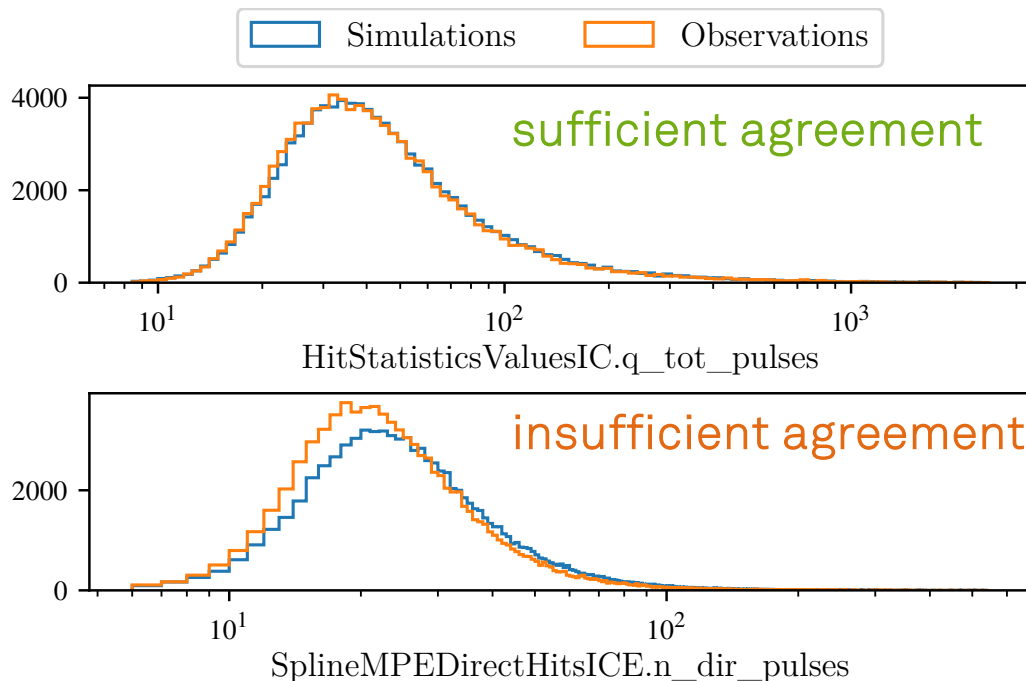
$$k = |A| = |B|$$

$$r = |A \cap B|$$

$n$: number of features

# Detection of Data/Simulation Disagreements



Legend: Simulations, Observations

sufficient agreement

HitStatisticsValuesIC.q__tot__pulses

insufficient agreement

SplineMPEDirectHitsICE.n__dir__pulses

Graphics: M. Linhoff [Learning Under Resource Constraints – Discovery in Physics] (in preparation)

Challenges when inspecting distributions by eye:

- only looking at one-dimensional distributions

- Systematic errors in simulation will also affect correlations between features

- Which metric ???

- Which threshold ???

# Detection of Data/Simulation Disagreements

Random Forest Feature Importance

| Feature | |
|---|---|
| SPEFit2Bayesian.time | |
| SplineMPEDirectHitsICE.n__dir__pulses | |
| HitStatisticsValuesIC.q__tot__pulses | |
| LineFitTimeSplit2Params.n__hits | |
| SplineMPE__MillipedeHighEnergyMIEFitParams.logl | |
| SplineMPEMuEXDifferential.energy | |
| SPEFitSingle__TWHVFitParams.rlogl | |
| SplineMPE__SegementFitParams.rlogl | |
| LineFitGeoSplit1Params.n__hits | |
| SplineMPE__MillipedeHighEnergyMIEFitParams.chi__squared | |
| MuEXAngular4__rllt.value | |
| SplineMPEDirectHitsD.dir__track__length | |
| SplineMPECramerRaoParams.cramer__rao__theta | |
| ProjectedQ.ratio | |
| SplineMPECramerRaoParams.variance__theta | |
| BestTrackCramerRaoParams.variance__theta | |
| SPEFit2TimeSplit1Bayesian.z | |
| SplineMPETruncatedEnergy__SPICEMie__DOMS__Muon.energy | |
| SplineMPEDirectHitsICC.dir__track__hit__distribution__smoothness | |
| SplineMPECramerRaoParams.variance__x | |

0.00    0.05

Graphics: M. Linhoff [Learning Under Resource
Constraints – Discovery in Physics] (in preparation)
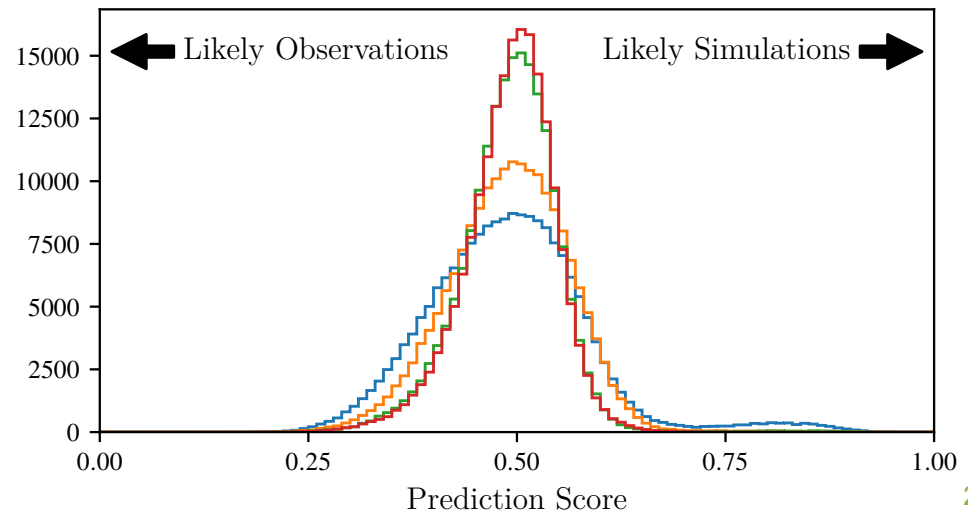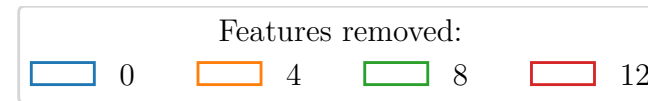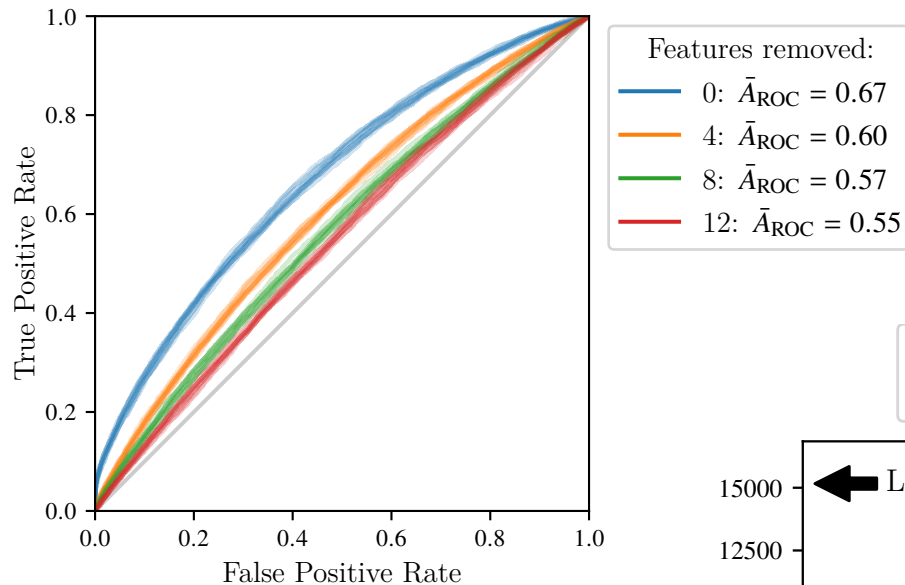
General Idea:

- Train classifier to distinguish simulated and experimental data
- Hard to impossible for a perfect agreement
- Sort features according to their importance
- Discard top n features
- Advantage: Extent to which mismatches can be tolerated is set by the classifier

# Detection of Data/Simulation Disagreements

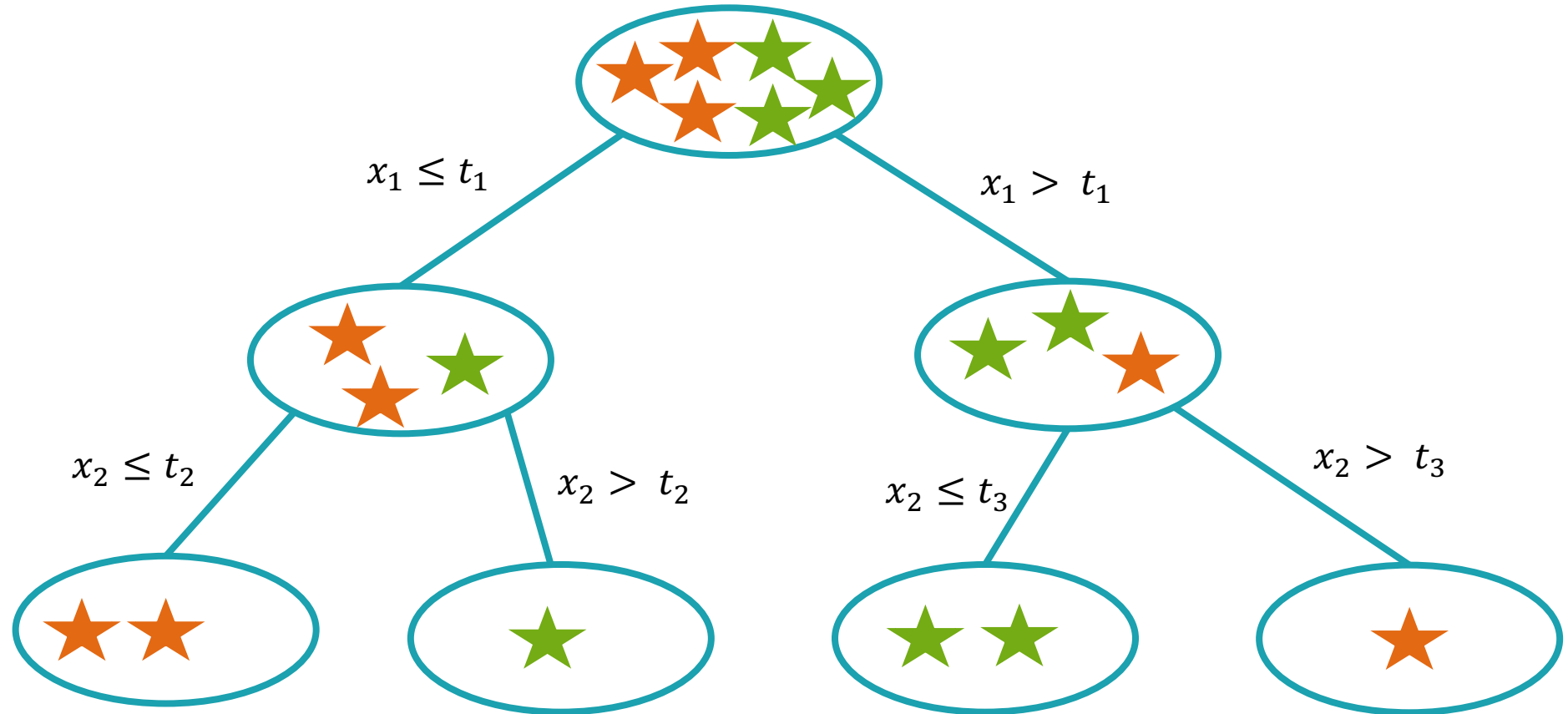Graphics: M. Linhoff [Learning Under Resource Constraints – Discovery in Physics] (in preparation)

**Features removed:**
- 0: $\bar{A}_{\text{ROC}} = 0.67$
- 4: $\bar{A}_{\text{ROC}} = 0.60$
- 8: $\bar{A}_{\text{ROC}} = 0.57$
- 12: $\bar{A}_{\text{ROC}} = 0.55$

Area under Curve is close to 0.5 (close to random guess).

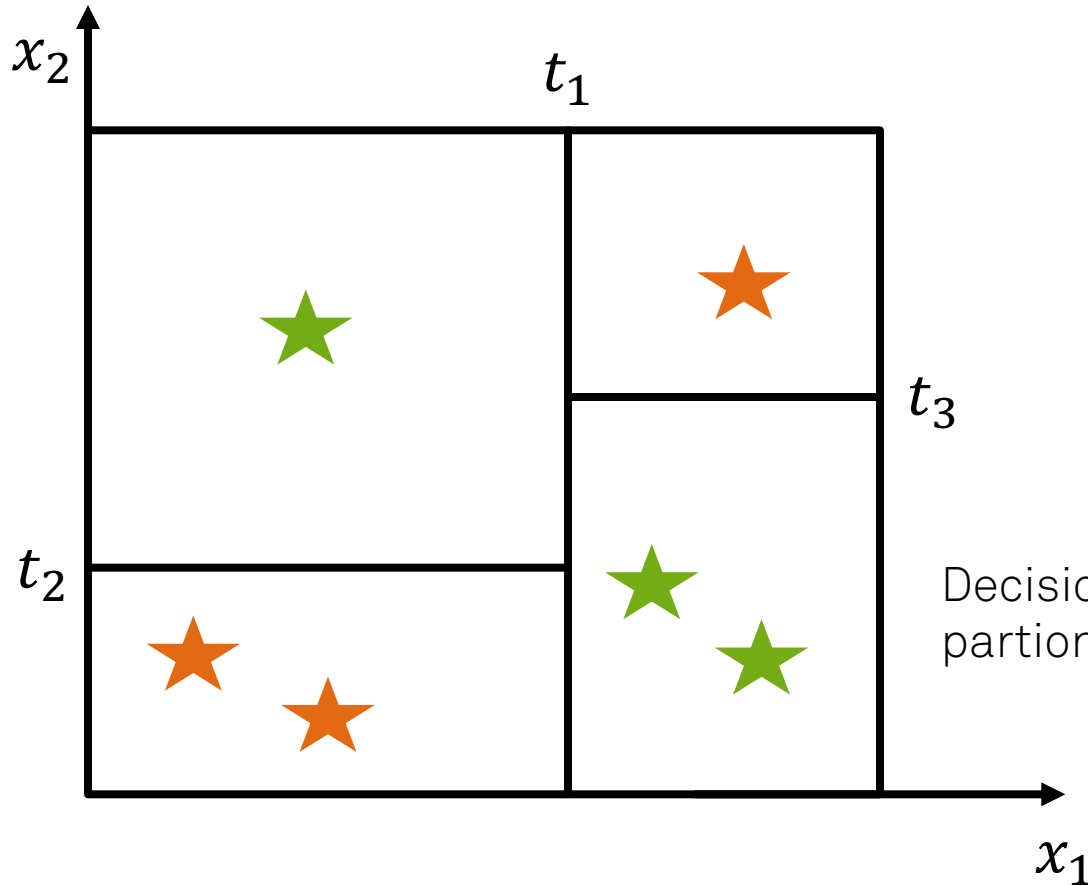Prediction score centered around 0.5 (close to random guess).

**Features removed:** 0  4  8  12

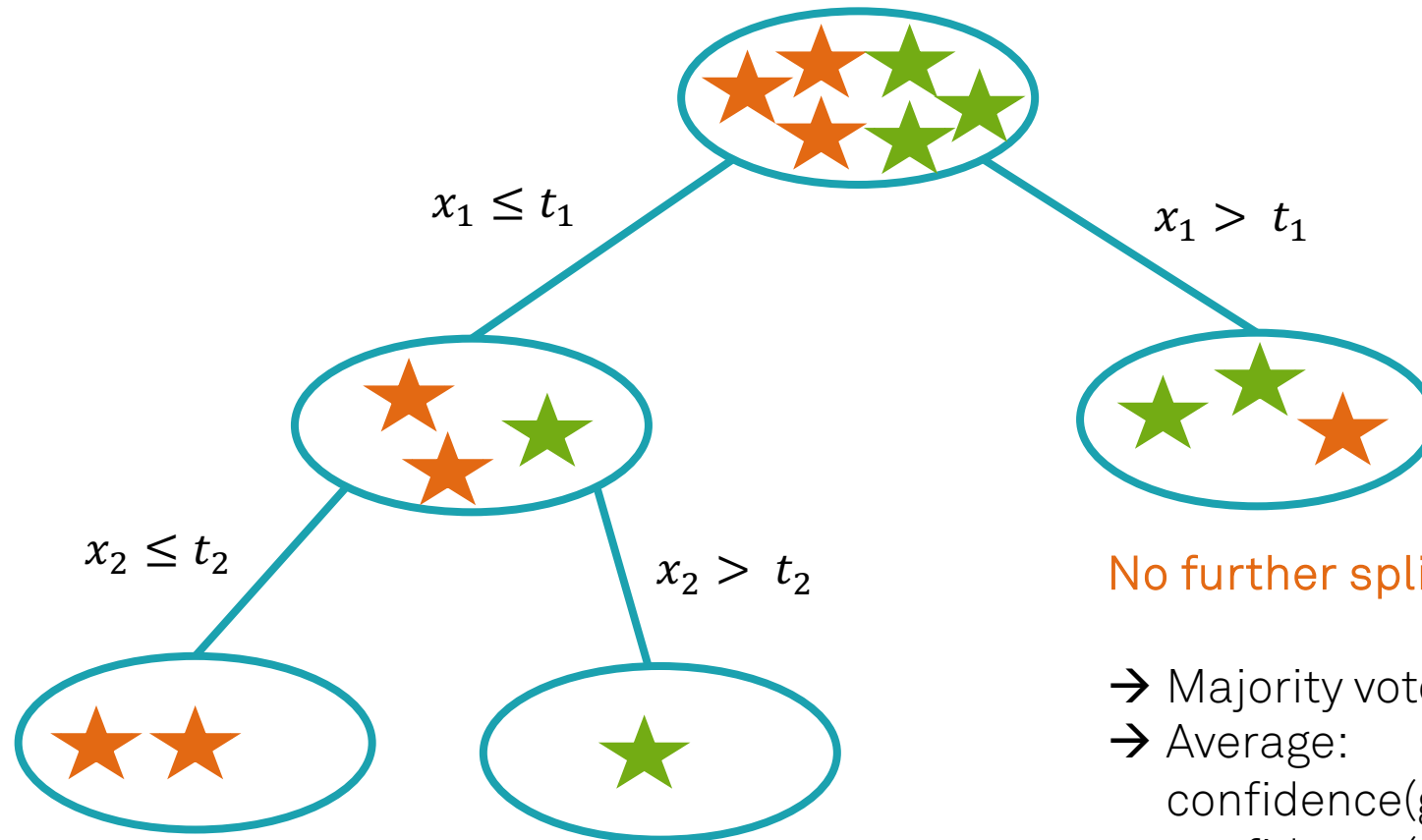Likely Observations ← → Likely Simulations

# Replacing Missing Values



- Some algorithms can only handle numerical values
- Missing values (nans, infs, …) can be replaced
  - Average
  - Median
  - Constant
  - …
- Features with too many missing values can be excluded
- Missing value can actually provide valuable information, e.g. reconstruction algorithm failed because this is an event with poor information

# Decision Trees



The decision tree structure:

- Root node with condition $x_1 \leq t_1$ (left) and $x_1 > t_1$ (right)
- Left node splits with $x_2 \leq t_2$ and $x_2 > t_2$
- Right node splits with $x_2 \leq t_3$ and $x_2 > t_3$

# Decision Trees Nomenclature

root node

$$x_1 \leq t_1$$

Depth of the tree: 3

$$x_2 \leq t_2$$

$$x_2 > t_2$$

terminal nodes/leaves

## Decision Trees



Decision Trees achieve a binary partioning of the feature space.

# Decision Trees

$x_1 \leq t_1$     $x_1 > t_1$

$x_2 \leq t_2$     $x_2 > t_2$

No further split possible

→ Majority vote: green
→ Average:
   confidence(green)= 2/3,
   confidence(orange) = 1/3

# Decision Trees for Regression



Root node: $1$ $1$ $2$ $2$ $1$ $2$ $3$

$x_1 \leq t_1$     $x_1 > t_1$

Left child: $1$ $1$ $3$

Right child: $2$ $2$ $1$

$x_2 \leq t_2$     $x_2 > t_2$

Leaf: $1$ $1$

Leaf: $3$

No further split possible

$$\frac{1}{3}(2 + 2 + 1) = \frac{5}{3}$$

# How to decide where to split?

Impurity Measures:

$P(\omega_j)$: Fraction of patterns at node $N$ in class $\omega_j$

Cross entropy:

$$i(N) = -\sum_{j=1}^{K} P(\omega_j) log_2(\omega_j)$$

Gini Impurity:

$$i(N) = \sum_{i \neq j} P(\omega_i)P(\omega_j) = \frac{1}{2}\left[1 - \sum_{j} P^2(\omega_j)\right]$$

# How to decide where to split?

Misclassification impurity:

$$i(N) = 1 - max_j P(\omega_j)$$

To determine the optimal split consider the decrease in impurity:

$$\Delta i(N) = i(N) - P_L \cdot i(N_L) - (1 - P_L) \cdot i(N_R)$$

# Boosted Decsion Trees

FINAL CLASSIFIER

$$G(x) = \text{sign}\left[\sum_{m=1}^{M} \alpha_m G_m(x)\right]$$

Weighted Sample $\dashrightarrow$ $G_M(x)$

Weighted Sample $\dashrightarrow$ $G_3(x)$

Weighted Sample $\dashrightarrow$ $G_2(x)$

Training Sample $\dashrightarrow$ $G_1(x)$

- Classifiers are weighted by
  $$\alpha_m = \log((1 - err_m)/err_m)$$
- Better classifiers obtain higher weights
- Example weights are updated in every iteration
  $$w_i \leftarrow w_i \cdot \exp(\alpha_m \cdot I(y_i \neq G(x_i))$$
- Falsely classified examples obtain higher weights in the next iteration

Source: Elements of Statistical Learning, Figure 10.1

# Random Forest



$x_1 \leq t_1$

Random subset of examples to build each tree.

Random subset features to determine the optimal split



Random forests utilize an ensemble of independent weak classifiers (decision trees) to obtain a better classification.

Final classification is achieved via:

$$c_j = \frac{1}{n_{tress}} \sum_{i=1}^{n_{trees}} c_{ij}$$

$c_{ij}$: Classification for example j by tree i

## Validation



- In order to not optimize on statistical fluctuations in the test set it is advisable to use cross validation for parameter optimization
- It can also be useful to have an additional validation set to validate the performance of the optimized parameter settings

 Use for training

 Use for testing

# Neural Networks



Input

Hidden

Output



Source: By Alvesgaspar - Top left: File:Cat August 2010-4.jpg by AlvesgasparTop middle: File:Gustav chocolate.jpg by Martin BahmannTop right: File:Orange tabby cat sitting on fallen leaves-Hisashi-01A.jpg by HisashiBottom left: File:Siam lilacpoint.jpg by Martin BahmannBottom middle: File:Felis catus-cat on snow.jpg by Von.grzankaBottom right: File:Sheba1.JPG by Dovenetel, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=17960205

Source: By en:User:Cburnett - This W3C-unspecified vector image was created with Inkscape., CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=1496812

- Highly successful, e.g in image classification
- Hyped
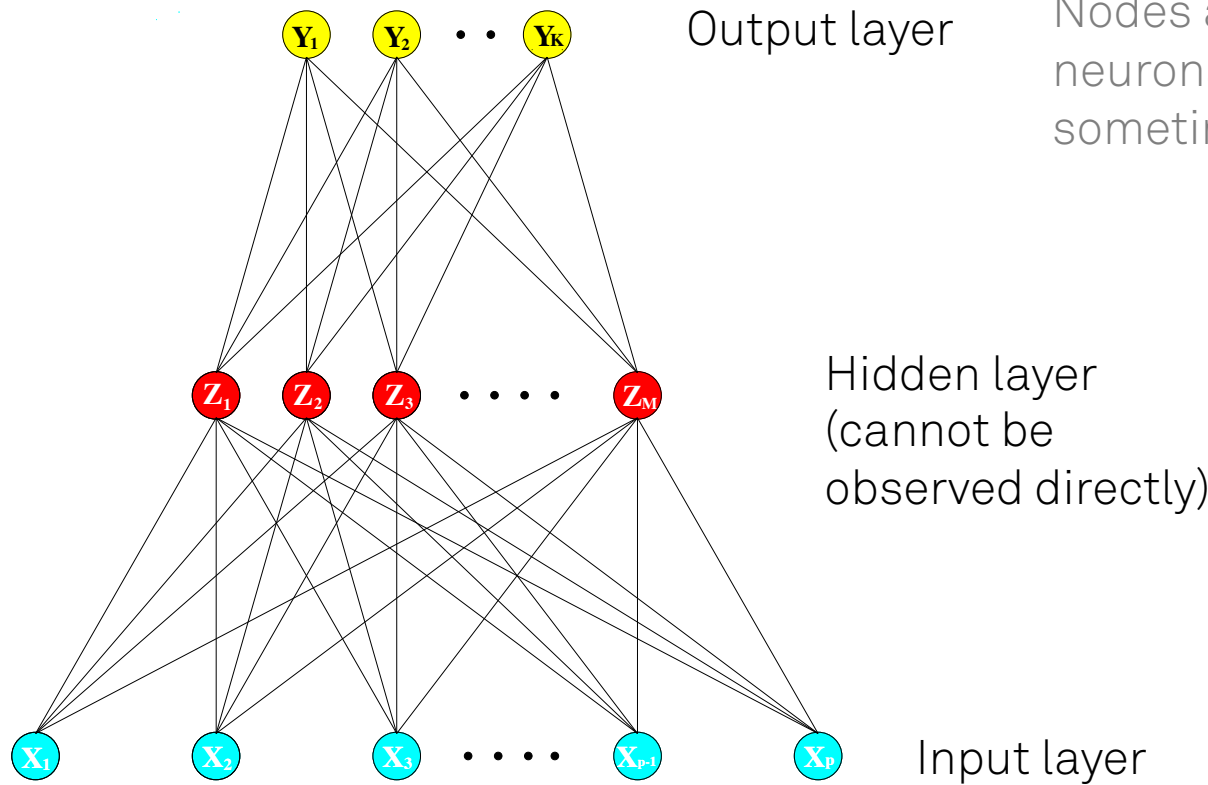
# Reasons for Using Neural Networks (in IceCube)

- Improved reconstruction methods will lead to increased sensitivity for the detection of sources
- Hardware limitations at the South Pole
- Events need to be processed in a given time frame to prevent pileup
- Limitations call for robust method that can handle raw data in constant time
- Neural networks are compuattionally inexpensive once the network is trained
- Fixed amount of operations, runtime is (largely) independent of the input
- Translational invariance (position of the classified object does not impact the class)
- Physics of neutrino interaction is invariant in time and space
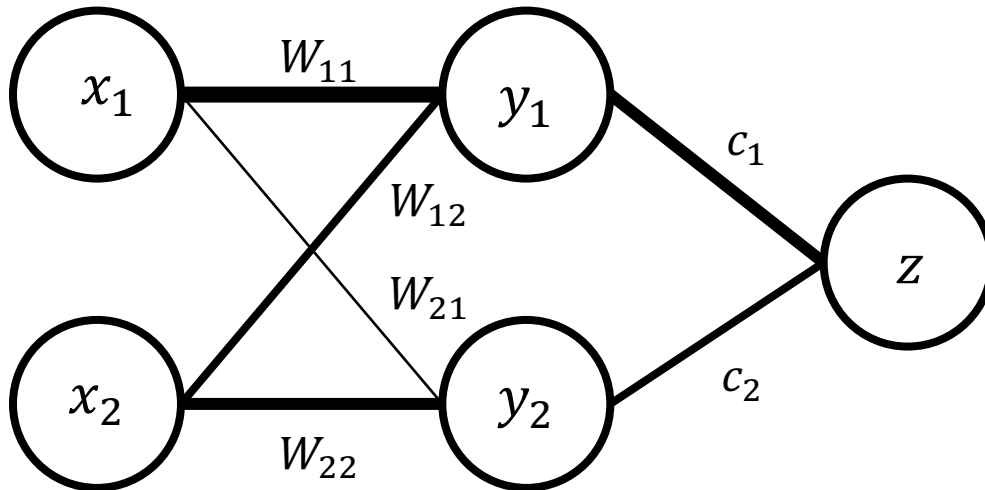
# General Idea of Neural Networks

Originally developed to mimic the human brain.

Nodes are sometimes called neurons, and connections are sometimes called synapses.

Output layer

Hidden layer
(cannot be
observed directly)

Input layer

## General Idea of Neural Networks

Nodes are linear combinations of nodes from previous layers.

The task is to optimize the weights such that the estimated label $z$ matches the true label $\hat{z}$.
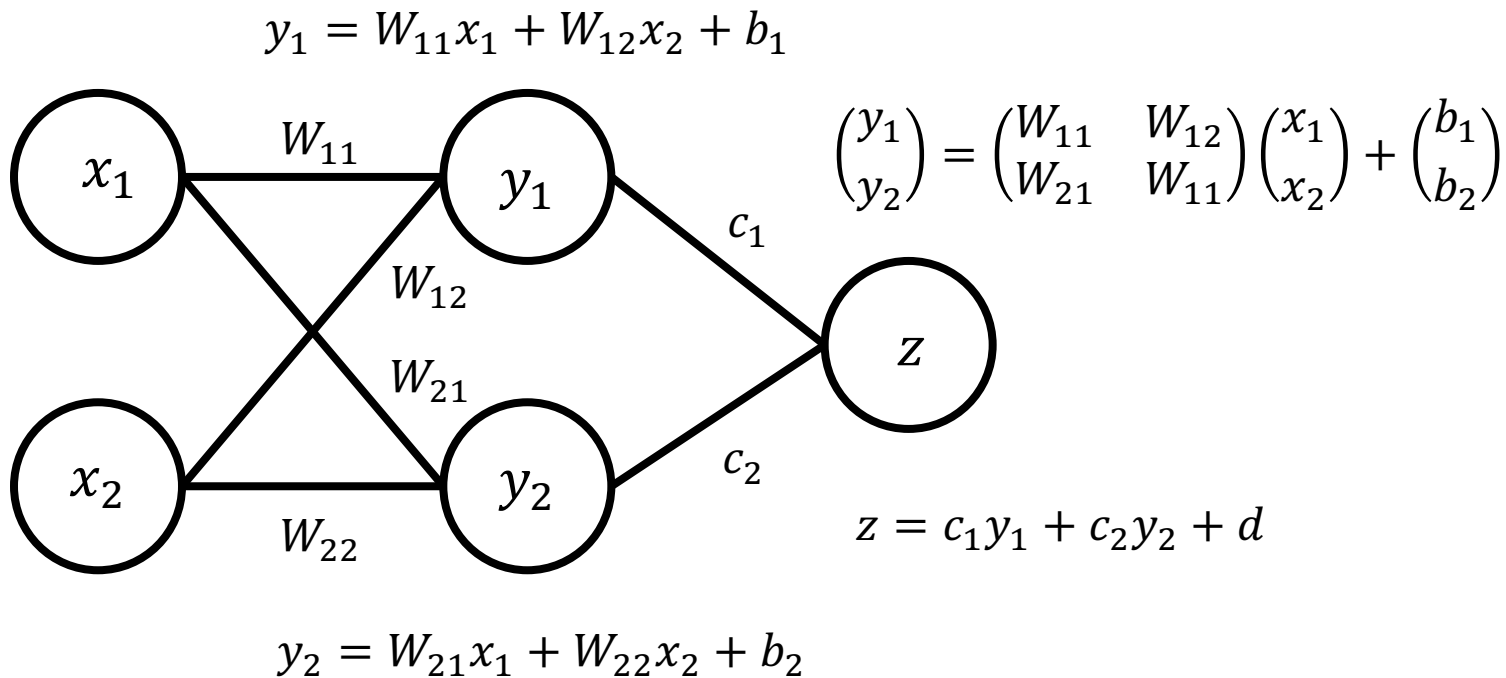


A feed-forward neural network with linear output and at least one hidden layer with a finite number of nodes can approximate any of the above* functions with arbitrary precision**.

*Continuous functions on closed bounded subsets of the Eucledian space $\mathbb{R}^n$.
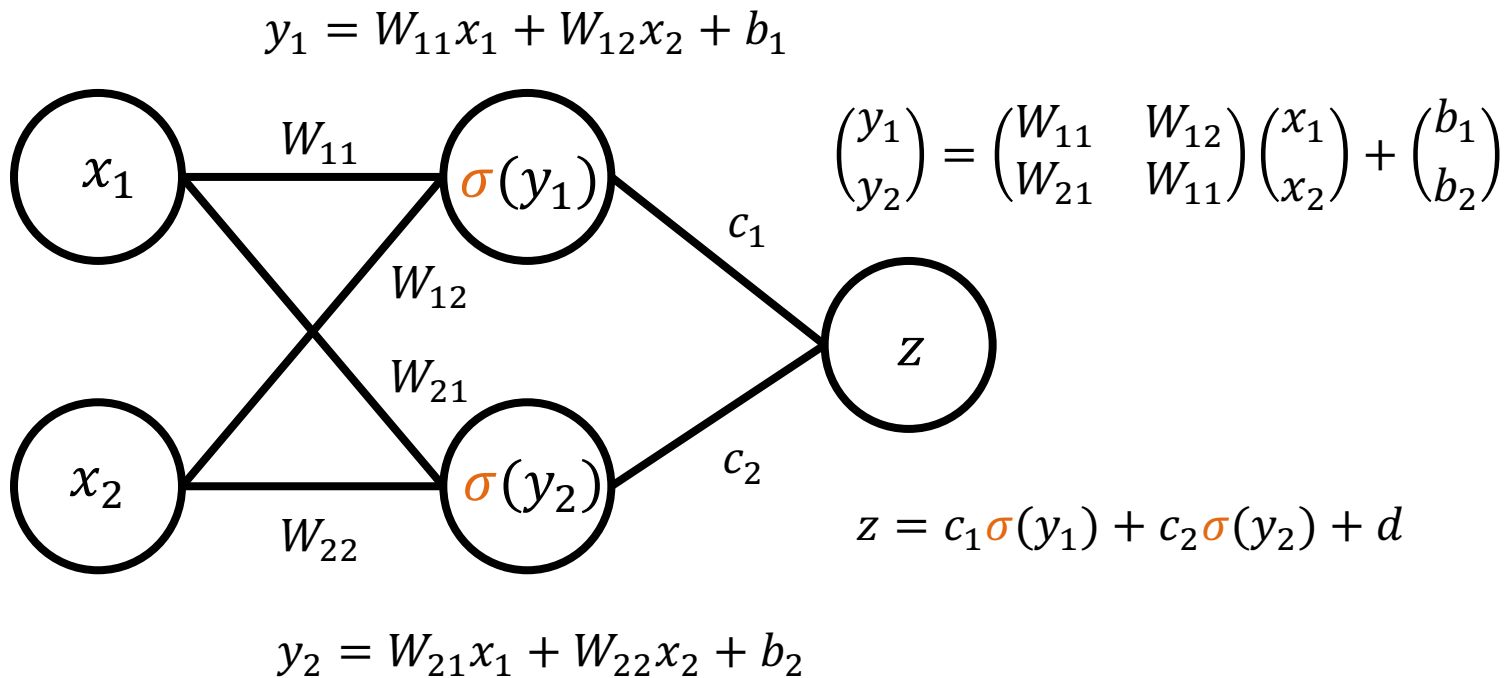
**Figure and definition adopted from Erdmann et al.

# More mathematically speaking

$$y_1 = W_{11}x_1 + W_{12}x_2 + b_1$$



$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} W_{11} & W_{12} \\ W_{21} & W_{11} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$z = c_1 y_1 + c_2 y_2 + d$$

$$y_2 = W_{21}x_1 + W_{22}x_2 + b_2$$

**Figure adopted from Erdmann et al.

## Adding non-linearity

$$y_1 = W_{11}x_1 + W_{12}x_2 + b_1$$



$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} W_{11} & W_{12} \\ W_{21} & W_{11} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$z = c_1\sigma(y_1) + c_2\sigma(y_2) + d$$

$$y_2 = W_{21}x_1 + W_{22}x_2 + b_2$$

**Figure adopted from Erdmann et al.

$\sigma$ is generally referred to as the activation function.

# Popular choices for $\sigma$

# Input Preparation

- **Zero-centered:** ReLU changes drastically around $0$, $x_i - \langle x_i \rangle$ to include positive and negative values
- **Order of magnitude:** Large variables could be preferred in the network training $x'_i = \frac{x_i - \langle x_i \rangle}{\sigma_i}$
- **Logarithm** to achieve more evenly distributed data
- **Decorrelation:** highly correlated variables should be decorrelated
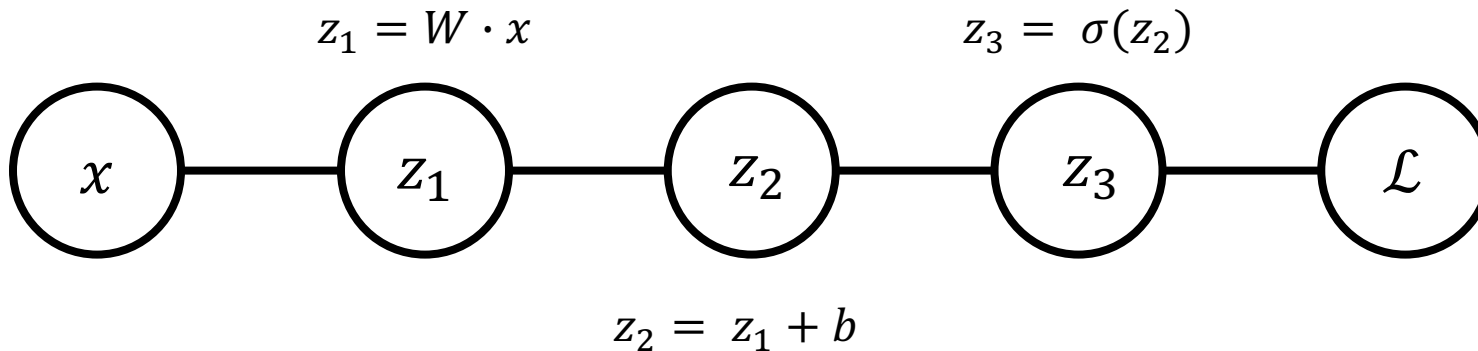
# Epoch and (mini)batch

- **Minibatch, Batch:** Using all examples can be infeasible in case many parameters need to be optimized, instead random subsets (batches) of examples are used. The optimal size of the batch depends on the problem to be solved. Popular choices are $2^k$

- **Epoch:** Complete use of all examples.

## Weight Updates

$$z_1 = W \cdot x$$

$$z_3 = \sigma(z_2)$$



$$z_2 = z_1 + b$$

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial z_3} \cdot \frac{\partial z_3}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial W}$$
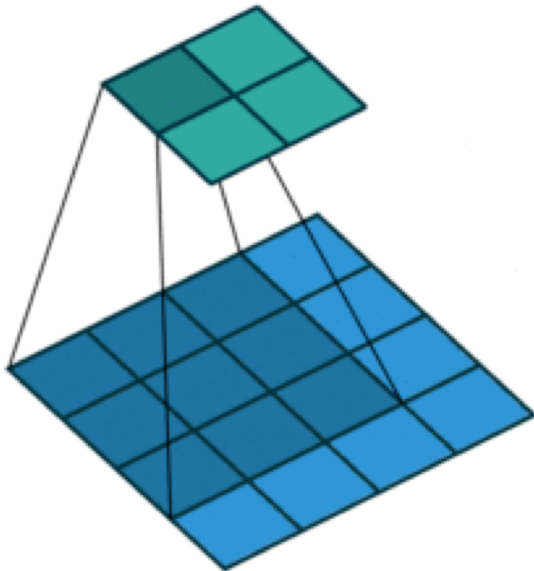
$$W_{t+1} = W_t - \alpha \mathbb{E}\left[\frac{\partial \mathcal{L}}{\partial W}\right]_t$$

$$\mathbb{E}\left[\frac{\partial \mathcal{L}}{\partial W}\right] = \frac{1}{k}\sum_{i=1}^{k}\left(\frac{\partial \mathcal{L}}{\partial W}\right)_i$$

This is the basic idea, this will most likely be handled by an optimizer.
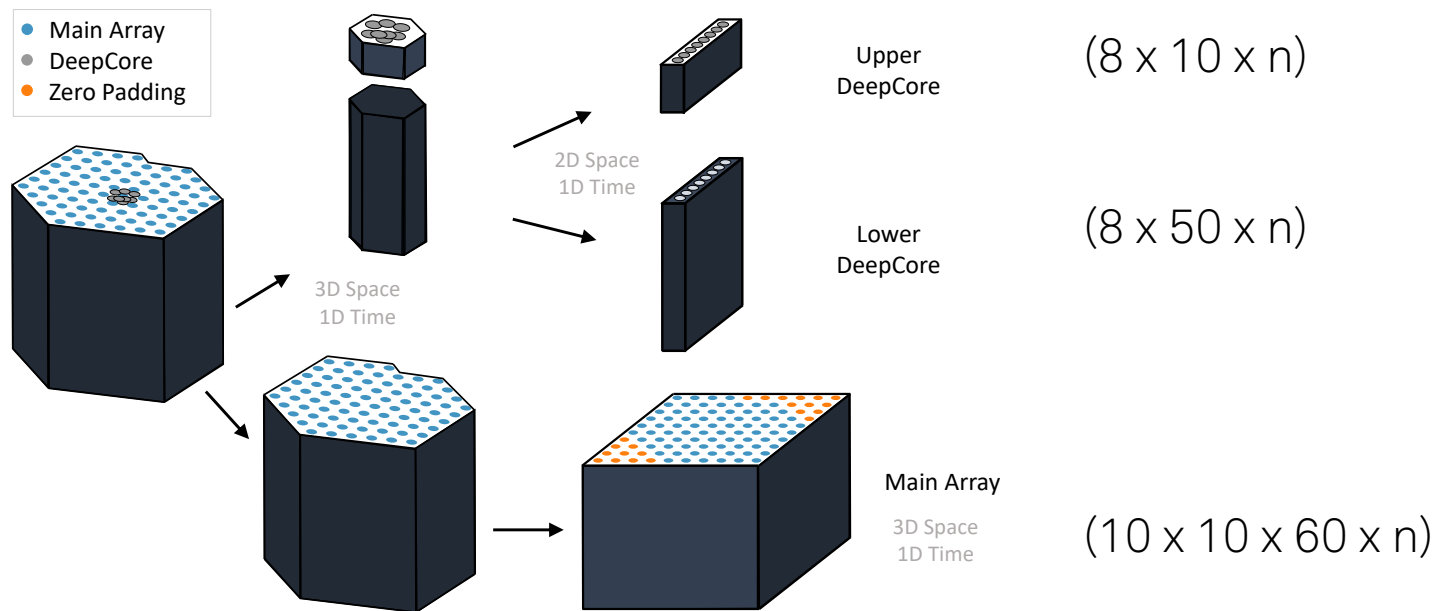
# Convolutional Layers



- Considering all pixels in an image in a fully connected network, results in too many parameters to be optimized
- The position of an object in an image should not alter the prediction (translational invariance)
- The convolutional operation exploits the neighbourhood of each pixel

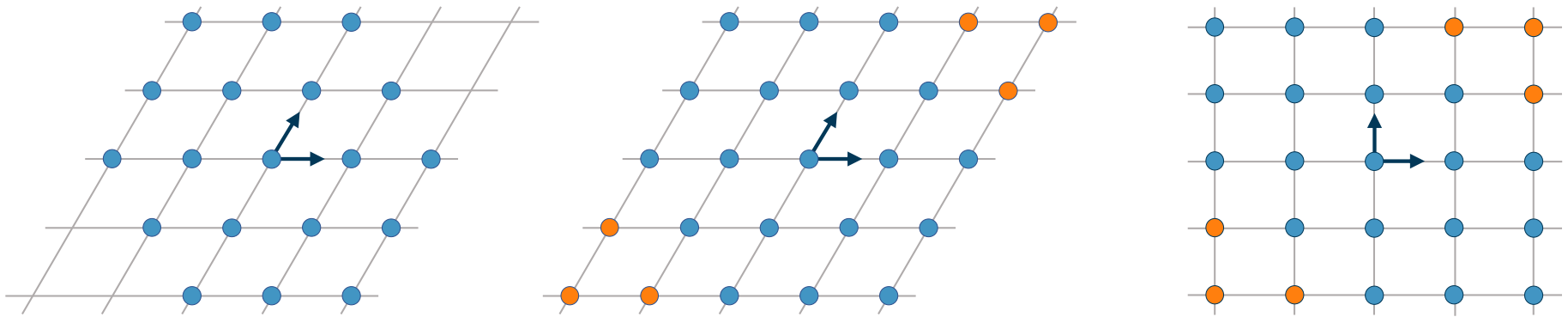Source: By Vincent Dumoulin, Francesco Visin - https://github.com/vdumoulin/conv_arithmetic, MIT, https://commons.wikimedia.org/w/index.php?curid=78003423

# Hexagonal Input Data



- Main Array
- DeepCore
- Zero Padding

Upper DeepCore

$(8 \times 10 \times n)$

2D Space
1D Time

Lower DeepCore

$(8 \times 50 \times n)$

3D Space
1D Time

Main Array

3D Space
1D Time

$(10 \times 10 \times 60 \times n)$

Abbasi et al., JINST, 16 (7) (2004).

# Hexagonal Kernels



Abbasi et al., JINST, 16 (7) (2004).

# Take-Away Messages

- Machine Learning and esp. Deep Learning is not magic
- Machine Learning and Deep Learning are tools that will help you to accomplish an analysis task faster and more accurately (when used correctly)
- The preprocessing of data is part of machine learning (and very important)
- Not every classifier is suited for every problem (consider runtime)
- If something fast and simple does the job: use it
- Make sure simulated and experimental data agree
- ...