

# Machine learning methods applied to the analysis of central exclusive production events in ALICE

Sebastian Ratzenböck<sup>1</sup>

<sup>1</sup>Stefan Meyer Institut  
Österreichische Akademie der Wissenschaften

26. April 2018

# Outline

1 Central exclusive production at ALICE

2 ML: an overview

3 Rectangular cuts

- Decision Trees
- Example
- Improvements

4 Linear cuts

- Single layer perceptron
- Activation functions
- Kernel trick

5 Non-linear cuts

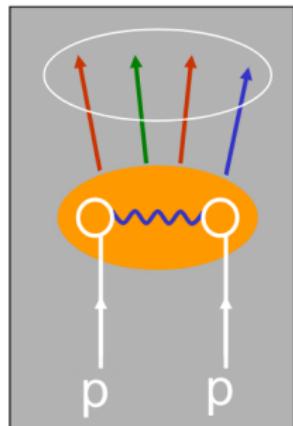
6 Neural networks

7 Results

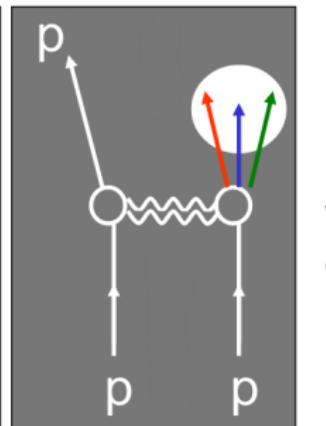
# Diffractive events

Diffractive events are reactions where no quantum numbers are exchanged

Incident hadrons  
acquire color  
 $\rightarrow$  break apart



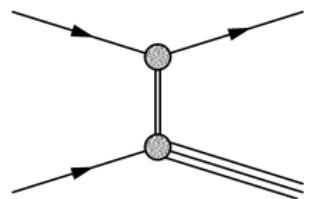
$\eta$  gap exponentially suppressed



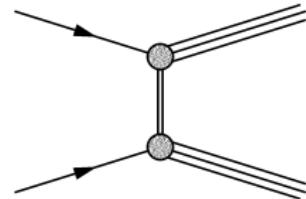
$\eta$  gap not suppressed

# Diffractive processes

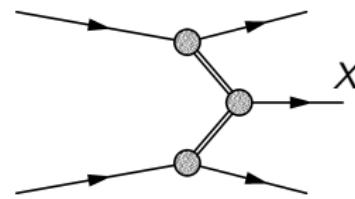
Different types of diffractive events are distinguished



Single diffractive



Double diffractive



CEP

Described by *Regge theory*

# Regge theory - short overview

- Reactions characterized by a color neutral  $t$ -channel exchange carrying vacuum quantum numbers
- $t$ -channel process governed by exchange of resonances sitting on *Regge trajectory*  $\alpha(t) = \alpha + \alpha't$
- Resonance of mass  $m$  sitting on the trajectory at  $t = m^2$  has angular momentum of value  $\alpha(m^2)$
- Trajectory makes contribution to total inelastic cross section  $\sigma_{tot}(s) \propto s^{\alpha-1}$

# Regge theory - short overview

- Reactions characterized by a color neutral  $t$ -channel exchange carrying vacuum quantum numbers
- $t$ -channel process governed by exchange of resonances sitting on *Regge trajectory*  $\alpha(t) = \alpha + \alpha' t$
- Resonance of mass  $m$  sitting on the trajectory at  $t = m^2$  has angular momentum of value  $\alpha(m^2)$
- Trajectory makes contribution to total inelastic cross section  $\sigma_{tot}(s) \propto s^{\alpha-1}$

# Regge theory - short overview

- Reactions characterized by a color neutral  $t$ -channel exchange carrying vacuum quantum numbers
- $t$ -channel process governed by exchange of resonances sitting on *Regge trajectory*  $\alpha(t) = \alpha + \alpha' t$
- Resonance of mass  $m$  sitting on the trajectory at  $t = m^2$  has angular momentum of value  $\alpha(m^2)$
- Trajectory makes contribution to total inelastic cross section  $\sigma_{tot}(s) \propto s^{\alpha-1}$

# Regge theory - short overview

- Reactions characterized by a color neutral  $t$ -channel exchange carrying vacuum quantum numbers
- $t$ -channel process governed by exchange of resonances sitting on *Regge trajectory*  $\alpha(t) = \alpha + \alpha't$
- Resonance of mass  $m$  sitting on the trajectory at  $t = m^2$  has angular momentum of value  $\alpha(m^2)$
- Trajectory makes contribution to total inelastic cross section  $\sigma_{tot}(s) \propto s^{\alpha-1}$

# Regge theory - short overview

- Pomeron  $\mathbb{P}$  represents special Regge trajectory & is a color singlet state of gluons & has vacuum quantum numbers
- Contribution to  $\sigma_{tot}$  grows with increasing energy
- Double Pomeron exchange = production mechanism to study glueballs in the mass region below 2 MeV
- Pomeron exchange leads to restriction of quantum numbers of  $X$   
 $\rightarrow J^{PC} = (\text{even})^{++}$

# Regge theory - short overview

- Pomeron  $\mathbb{P}$  represents special Regge trajectory & is a color singlet state of gluons & has vacuum quantum numbers
- Contribution to  $\sigma_{tot}$  grows with increasing energy
- Double Pomeron exchange = production mechanism to study glueballs in the mass region below 2 MeV
- Pomeron exchange leads to restriction of quantum numbers of  $X$   
 $\rightarrow J^{PC} = (\text{even})^{++}$

# Regge theory - short overview

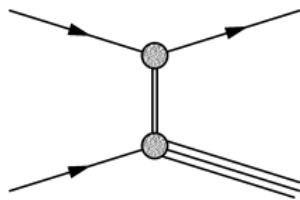
- Pomeron  $\mathbb{P}$  represents special Regge trajectory & is a color singlet state of gluons & has vacuum quantum numbers
- Contribution to  $\sigma_{tot}$  grows with increasing energy
- Double Pomeron exchange = production mechanism to study glueballs in the mass region below 2 MeV
- Pomeron exchange leads to restriction of quantum numbers of  $X$   
 $\rightarrow J^{PC} = (\text{even})^{++}$

# Regge theory - short overview

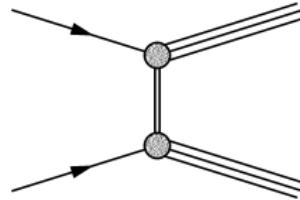
- Pomeron  $\mathbb{P}$  represents special Regge trajectory & is a color singlet state of gluons & has vacuum quantum numbers
- Contribution to  $\sigma_{tot}$  grows with increasing energy
- Double Pomeron exchange = production mechanism to study glueballs in the mass region below 2 MeV
- Pomeron exchange leads to restriction of quantum numbers of  $X$   
 $\rightarrow J^{PC} = (\text{even})^{++}$

# Event topology

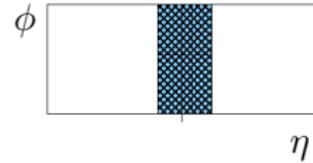
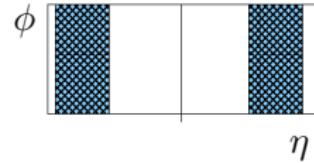
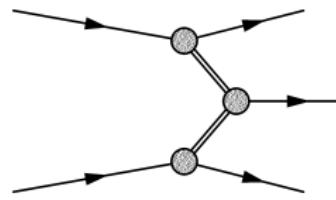
Single diffractive



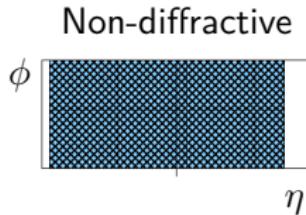
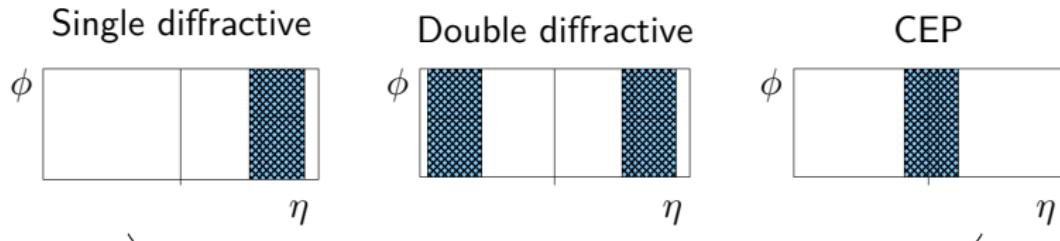
Double diffractive



CEP

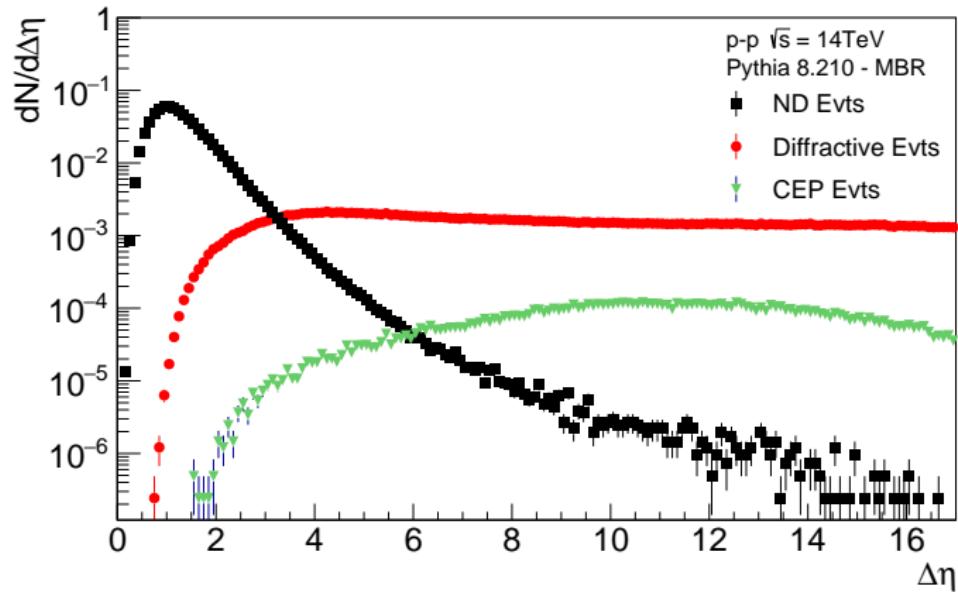


# Event topology



# Event selection

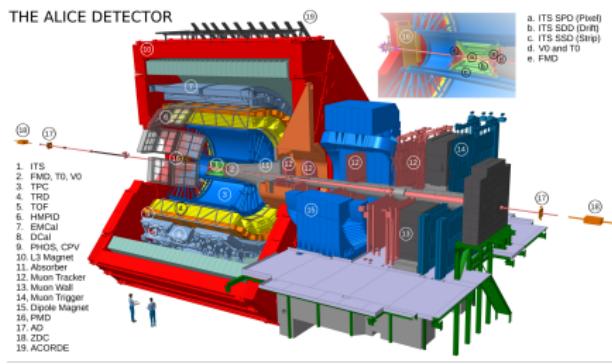
Large rapidity gaps in *non-diffractive* events are exponentially suppressed



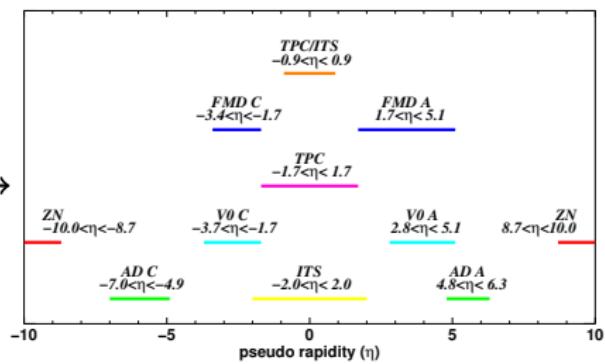
# The ALICE detector system

ALICE detector

THE ALICE DETECTOR



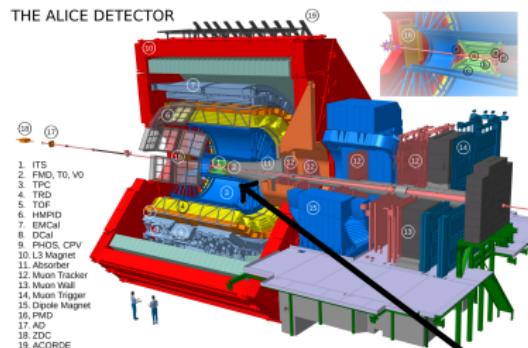
$\eta$  coverage



# The ALICE detector system

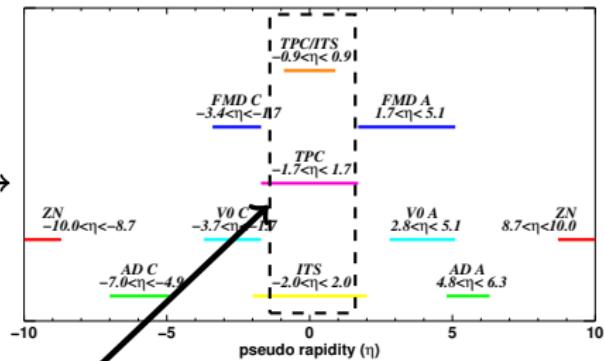
ALICE detector

THE ALICE DETECTOR



a.  
b.  
c.  
d.  
e.

$\eta$  coverage

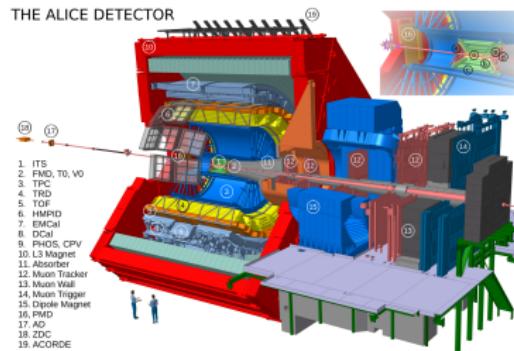


Central barrel  $\rightarrow$  determine  $p^\mu$

# The ALICE detector system

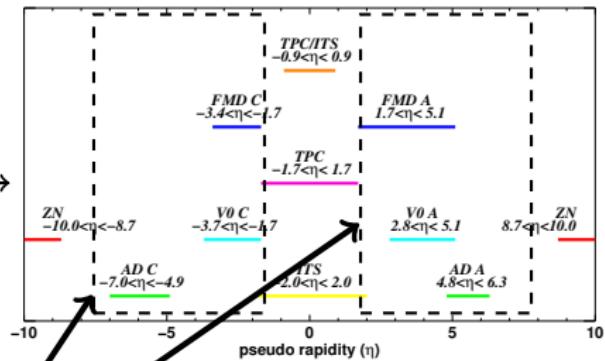
## ALICE detector

THE ALICE DETECTOR



- a. ITS SPD (Pixel)  
b. ITS TPC (Drift)  
c. ITS SSD (Stripl)  
d. V0 and T0  
e. FMD

## $\eta$ coverage

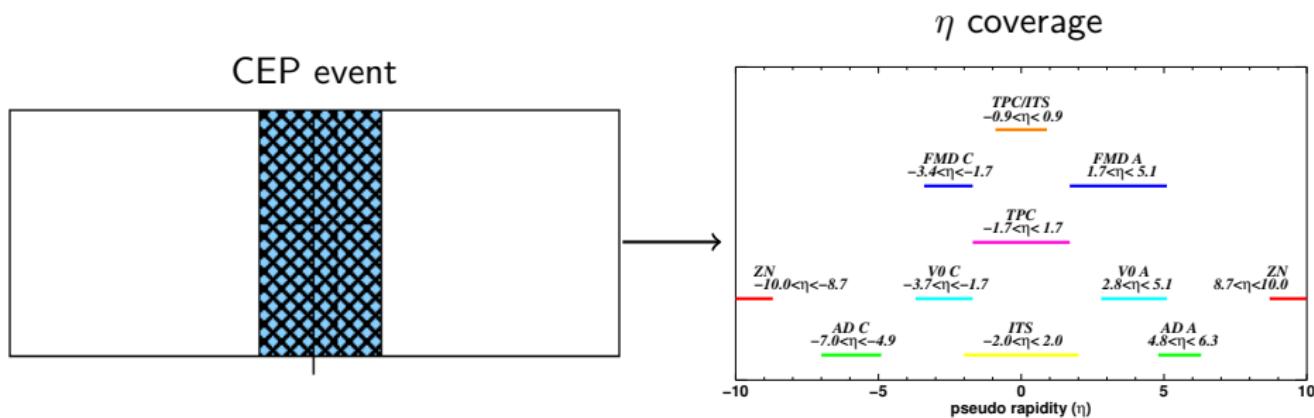


Small detectors for global event characterisation

# Central exclusive production at ALICE

To study the CEP events a  $\eta$  gap condition is used

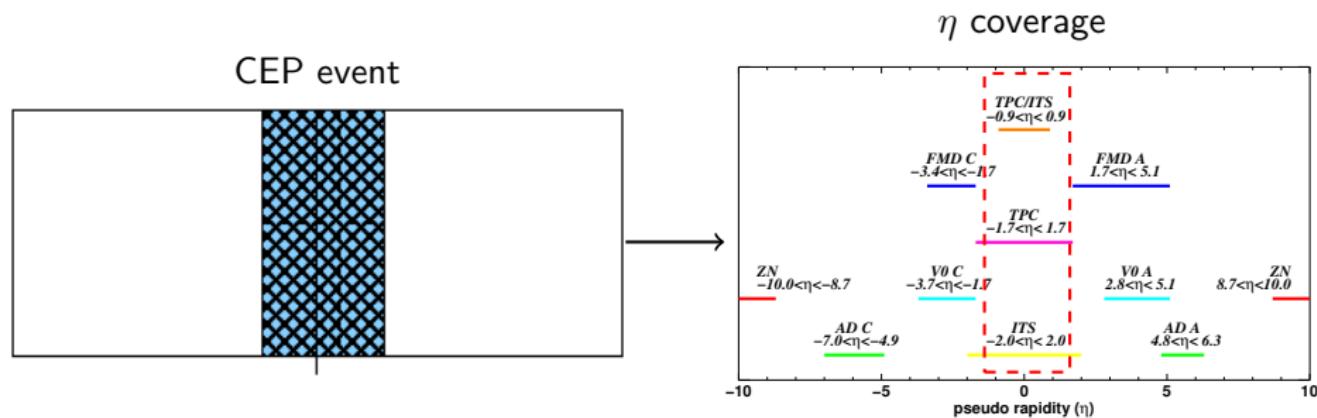
- ① Signal in the central barrel
- ② No activity in veto detectors



# Central exclusive production at ALICE

To study the CEP events a  $\eta$  gap condition is used

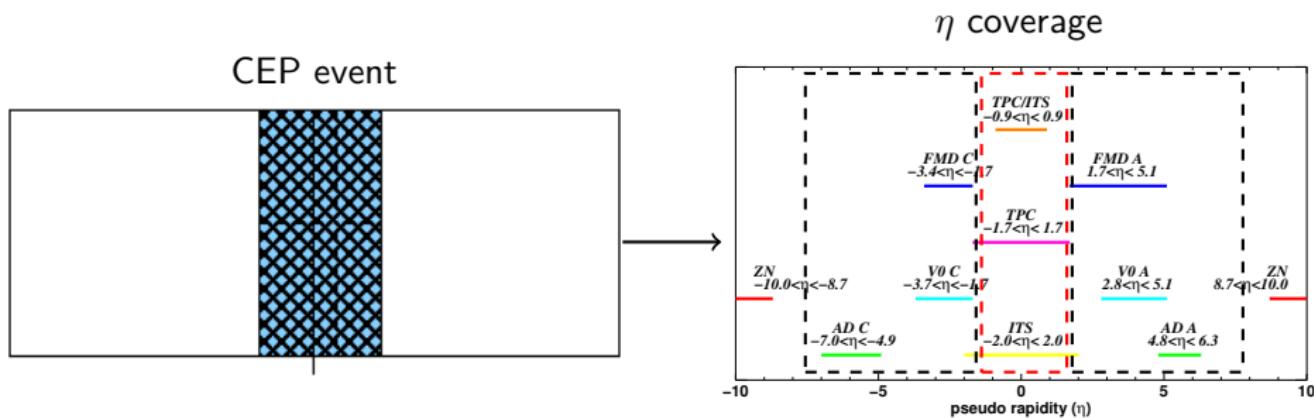
- ➊ Signal in the central barrel
- ➋ No activity in veto detectors



# Central exclusive production at ALICE

To study the CEP events a  $\eta$  gap condition is used

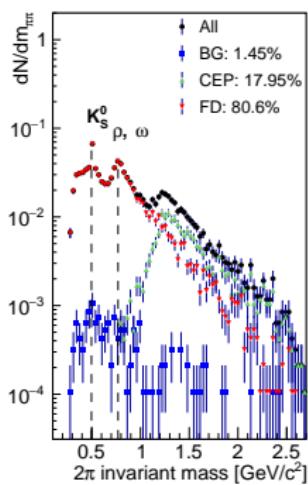
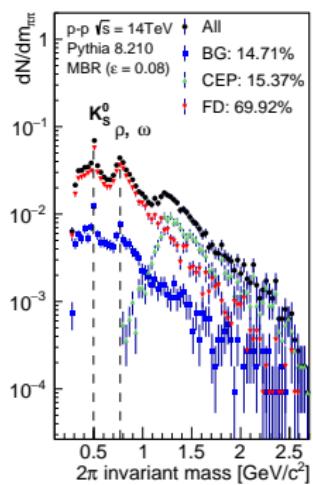
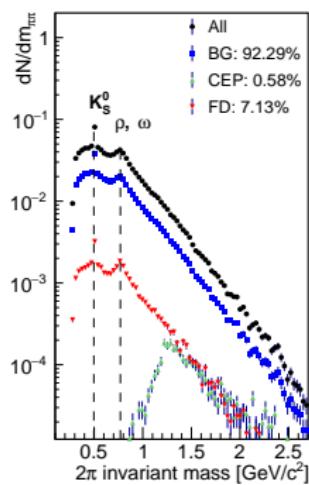
- ① Signal in the central barrel
- ② No activity in veto detectors



# Invariant mass spectrum

Studying *Pythia-8* simulations yield

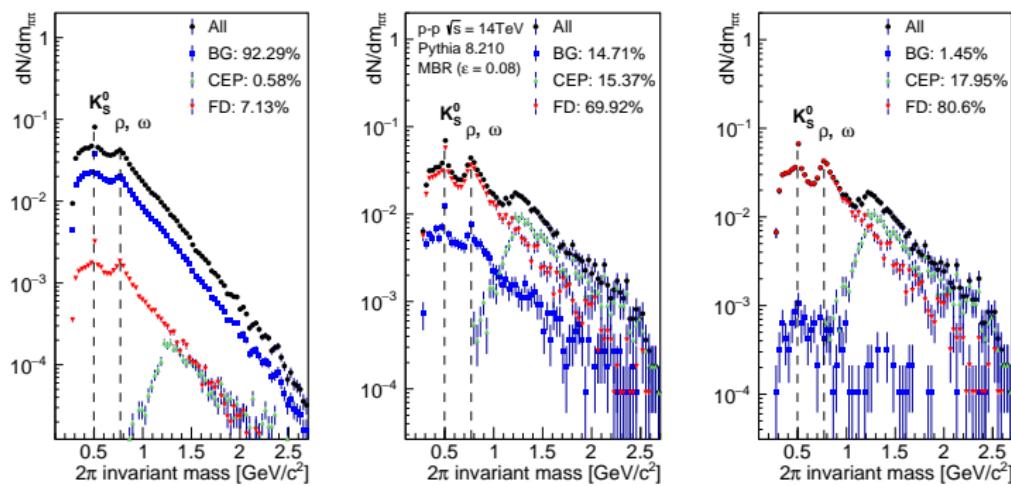
- Enforcing  $\eta$  gap cut reduces non-diffractive almost entirely
- Remaining background are partially reconstructed CEP events - feed down



# Invariant mass spectrum

Studying *Pythia-8* simulations yield

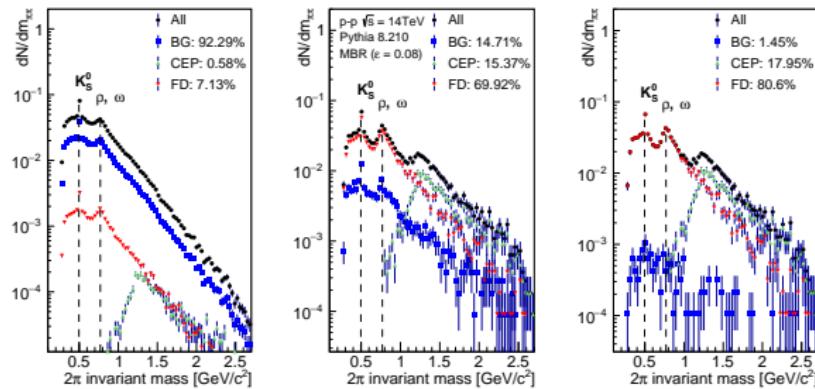
- Enforcing  $\eta$  gap cut reduces non-diffractive almost entirely
- Remaining background are partially reconstructed CEP events - **feed down**



# Invariant mass spectrum

Studying *Pythia-8* simulations yield

- Enforcing  $\eta$  gap cut reduces non-diffractive almost entirely
- Remaining background are partially reconstructed CEP events - **feed down**



# ML: an overview

In general ML represents a contrast to a *rule based systems*

## Rule-based system

System that uses rules to make deductions or choices

- Domain-specific expert system
- Knowledge base: facts & rules (if → then statement)
- Rules manually specified (by expert) → expensive, incomplete

# ML: an overview

In general ML represents a contrast to a *rule based systems*

## Rule-based system

System that uses rules to make deductions or choices

- Domain-specific expert system
- Knowledge base: facts & rules (if → then statement)
- Rules manually specified (by expert) → expensive, incomplete

# ML: an overview

In general ML represents a contrast to a *rule based systems*

## Rule-based system

System that uses rules to make deductions or choices

- Domain-specific expert system
- Knowledge base: facts & rules (if → then statement)
- Rules manually specified (by expert) → expensive, incomplete

# ML: an overview

In general ML represents a contrast to a *rule based systems*

## Machine learning

- Algorithms that learn from *data* & make predictions on *data*
- Automatic methods → no human needed
- Human work required for defining problem & assessing the data

# ML: an overview

In general ML represents a contrast to a *rule based systems*

## Machine learning

- Algorithms that learn from *data* & make predictions on *data*
- Automatic methods → no human needed
- Human work required for defining problem & assessing the data

# ML: an overview

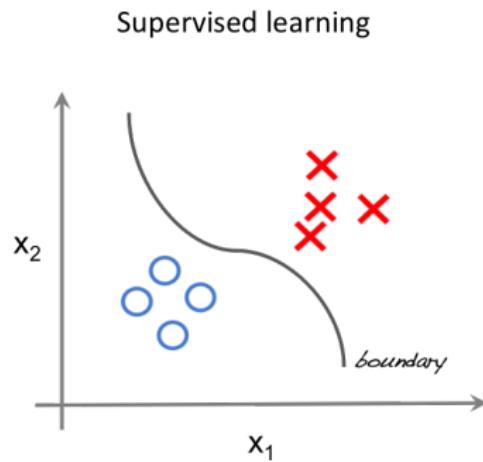
In general ML represents a contrast to a *rule based systems*

## Machine learning

- Algorithms that learn from *data* & make predictions on *data*
- Automatic methods → no human needed
- Human work required for defining problem & assessing the data

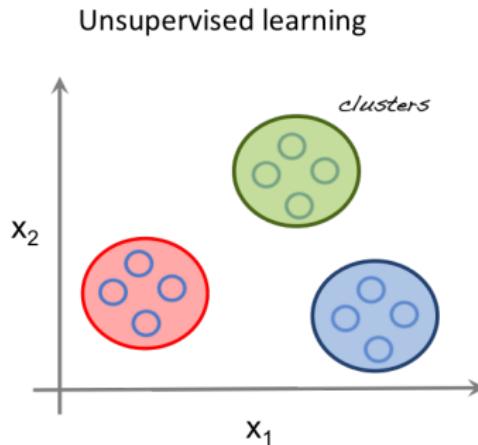
# Types of ML

- Supervised
  - ▶ Classification
  - ▶ Regression
- Unsupervised



# Types of ML

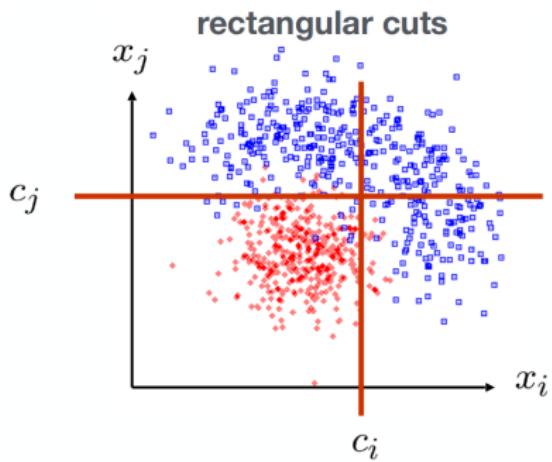
- Supervised
  - ▶ Classification
  - ▶ Regression
- Unsupervised



# Rectangular cuts

Standard cut in one variable

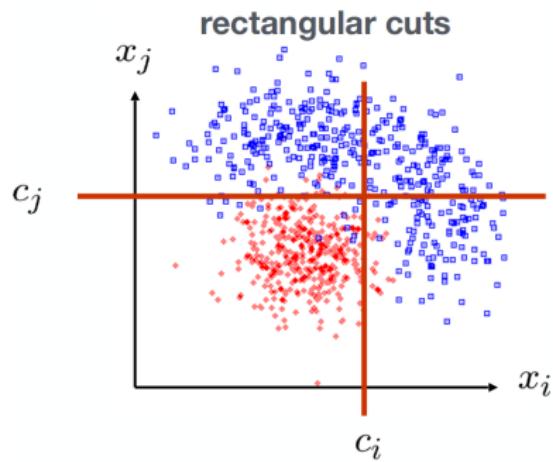
- Cuts only in lower-dimensional subspaces
- Ignores possible dependencies between the input variables
- Signal might behave like BG in several observables  
→ misclassification



# Rectangular cuts

Standard cut in one variable

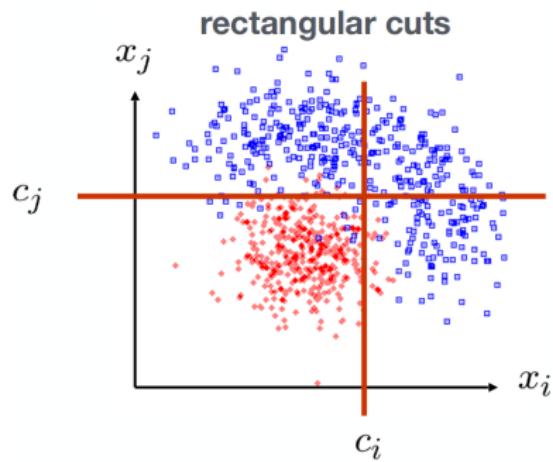
- Cuts only in lower-dimensional subspaces
- Ignores possible dependencies between the input variables
- Signal might behave like BG in several observables  
→ misclassification



# Rectangular cuts

Standard cut in one variable

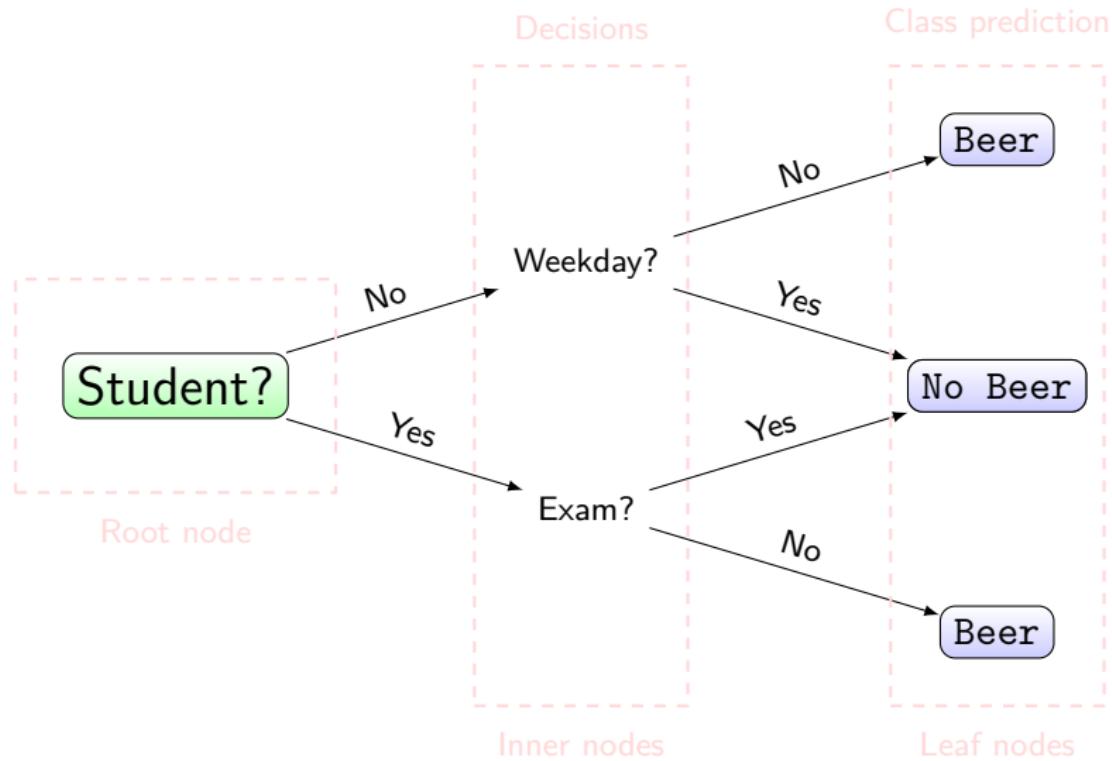
- Cuts only in lower-dimensional subspaces
- Ignores possible dependencies between the input variables
- Signal might behave like BG in several observables  
→ misclassification



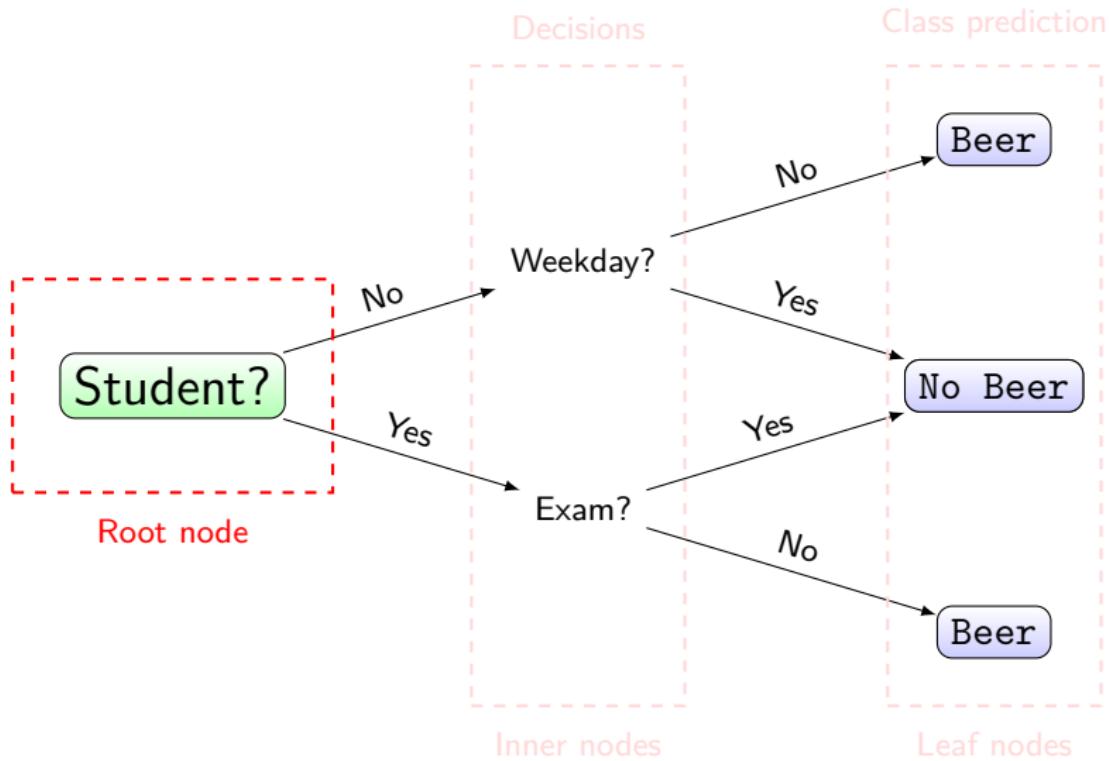
# Rectangular cuts with *decision trees*

- Tree-like graph → flowchart
- Easy to understand
- Either be manually modelled by experts or learned from training data

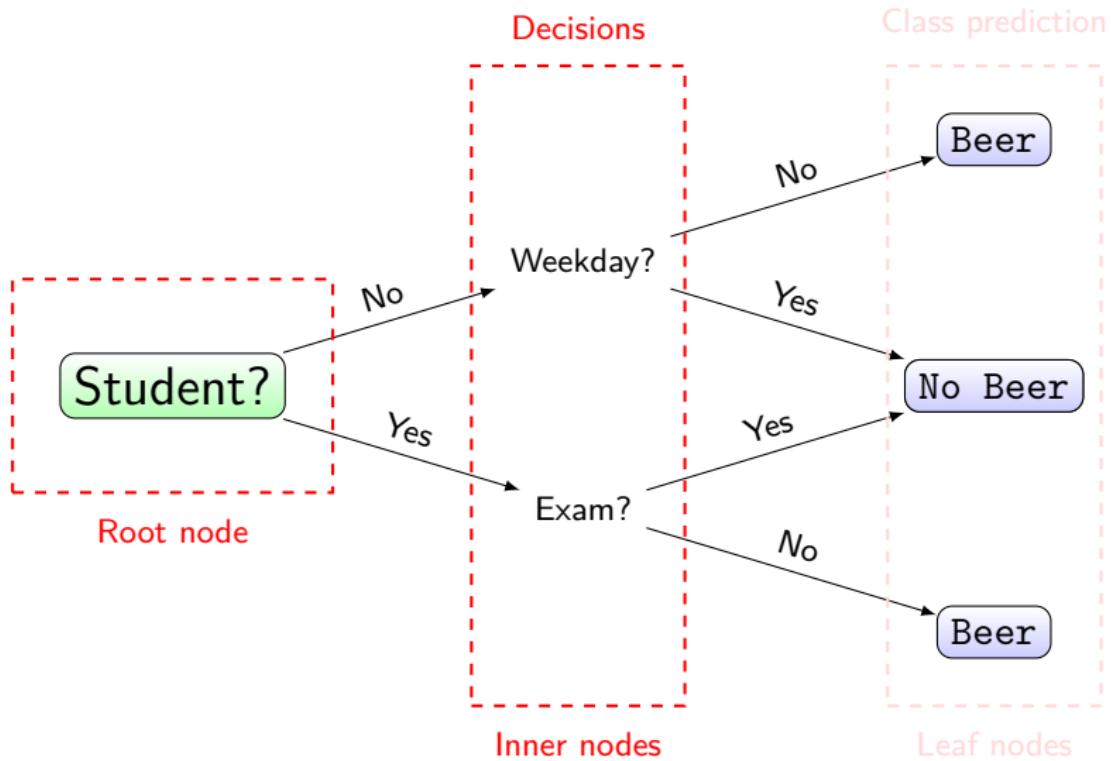
# Rectangular cuts with *decision trees*



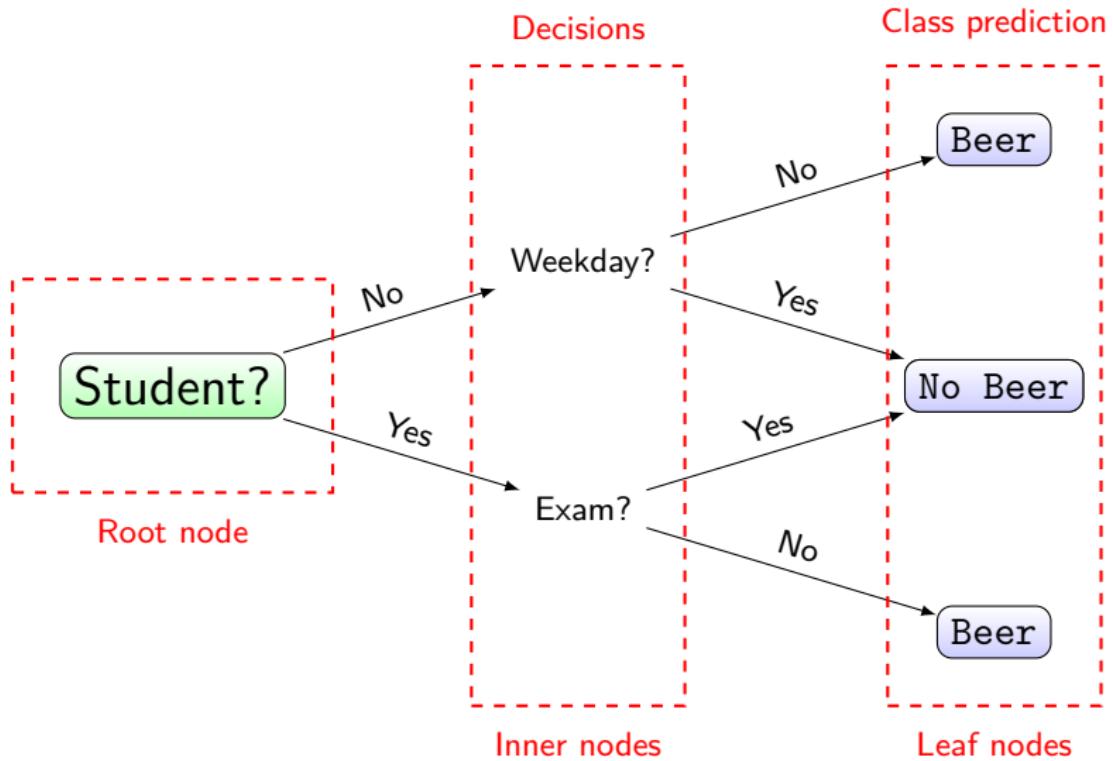
# Rectangular cuts with *decision trees*



# Rectangular cuts with *decision trees*



# Rectangular cuts with *decision trees*



# Decision tree learning

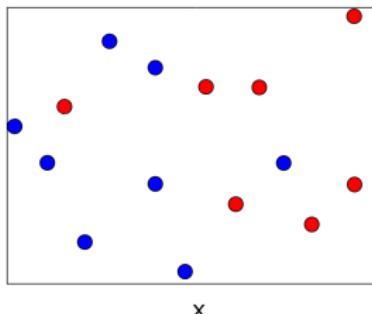
## Training

*Recursively split feature space into sub-spaces at each step*  
→ Measures to evaluate split

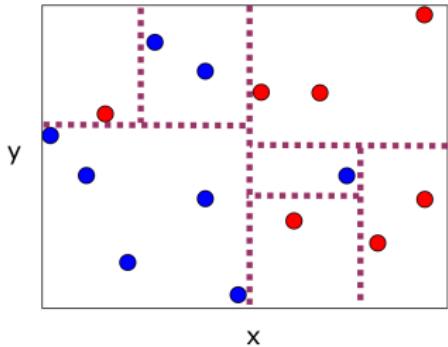
- Error rate
- Information gain
- Gini index

# Decision tree learning

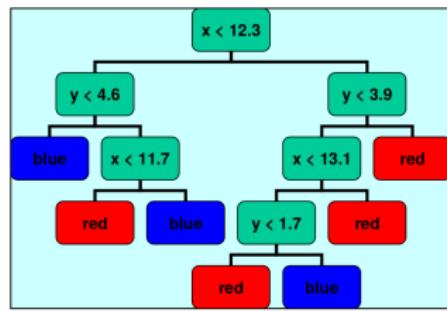
Feature space



Classification



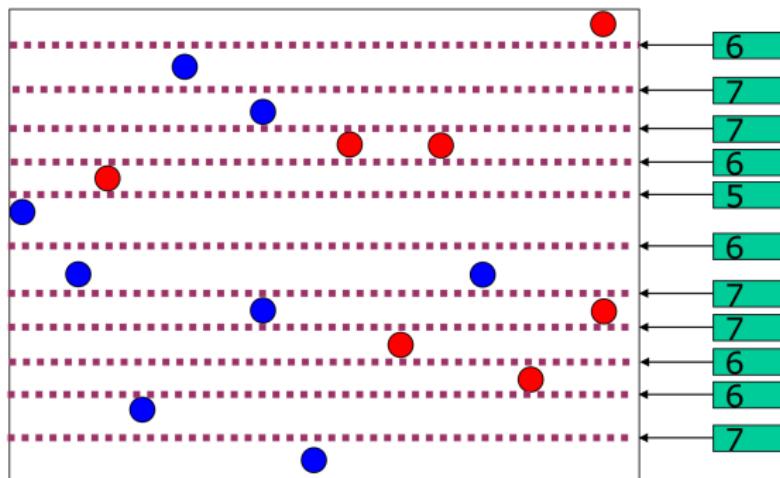
Decision tree



# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)

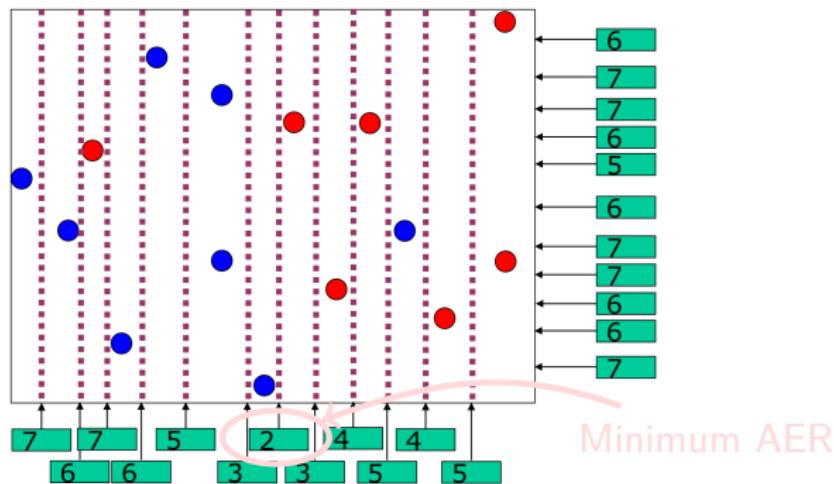
$y$  split AER



# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)

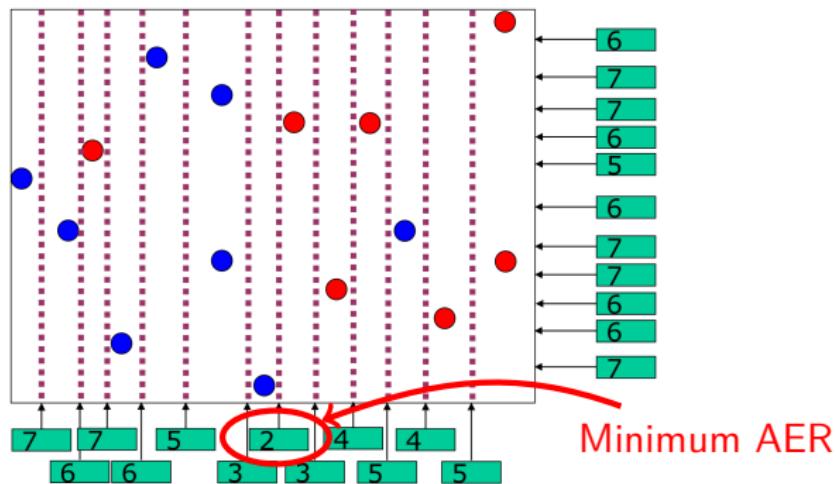
x-y split AER



# Decision tree learning

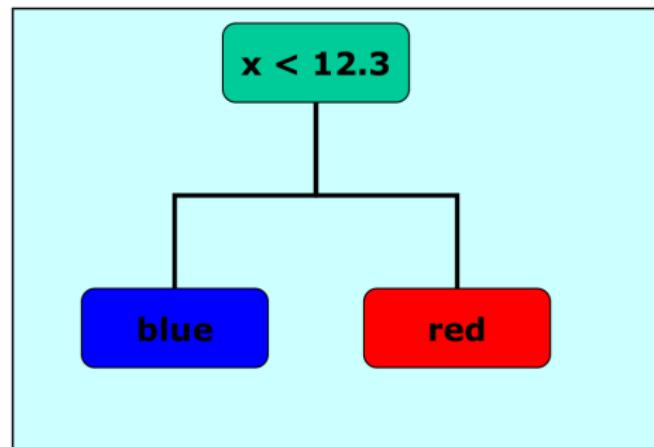
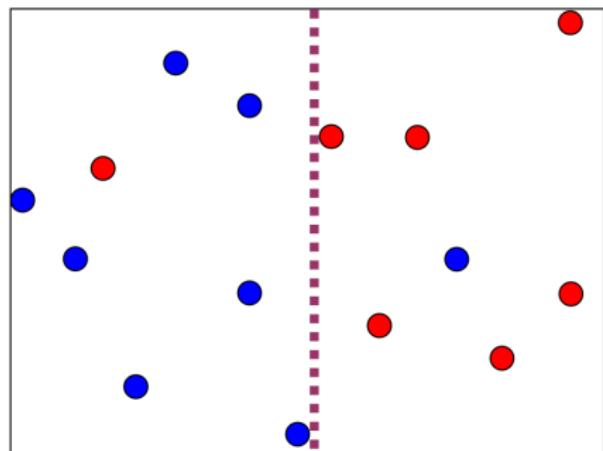
- 1) We compute a measure for *each possible split* in each feature  
 → here **absolute error rate** (AER)

x-y split AER



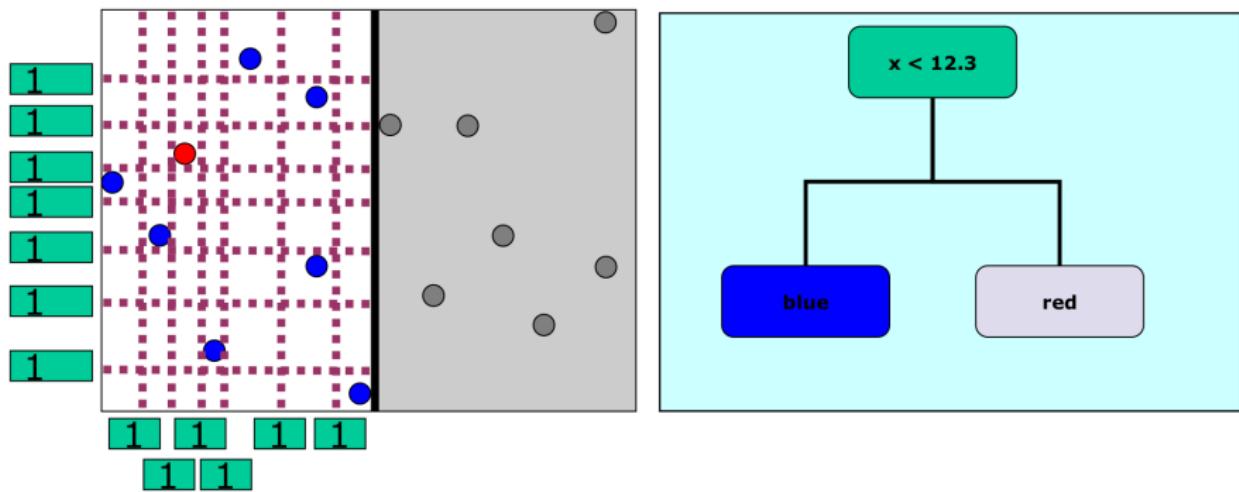
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until AER → 0



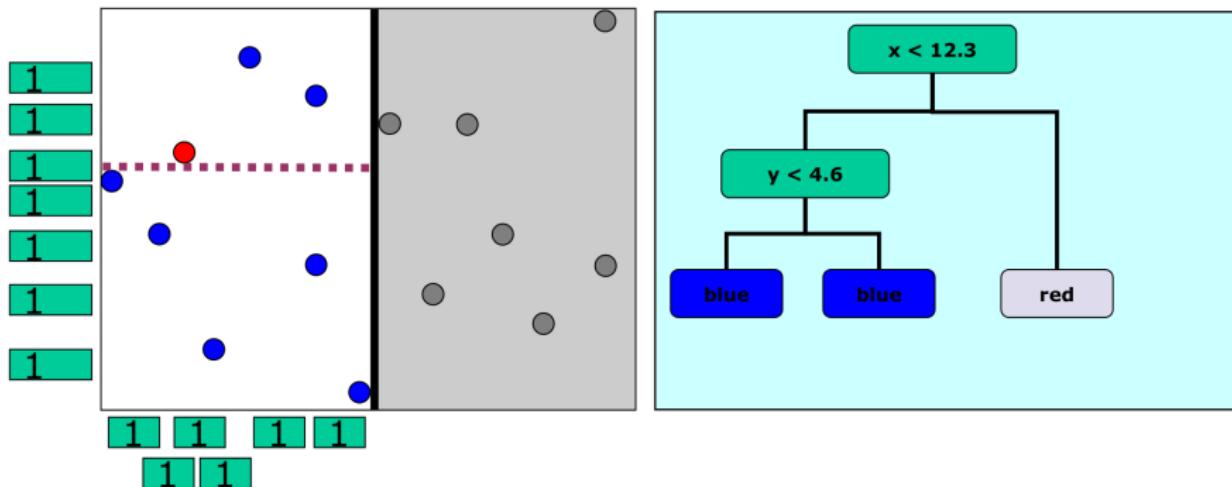
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until AER → 0



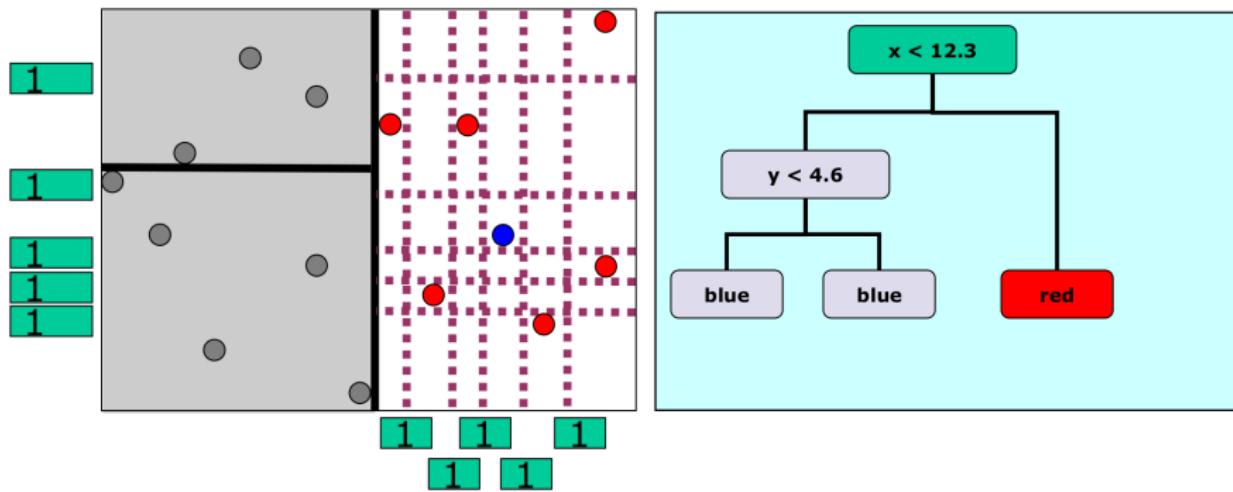
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until AER → 0



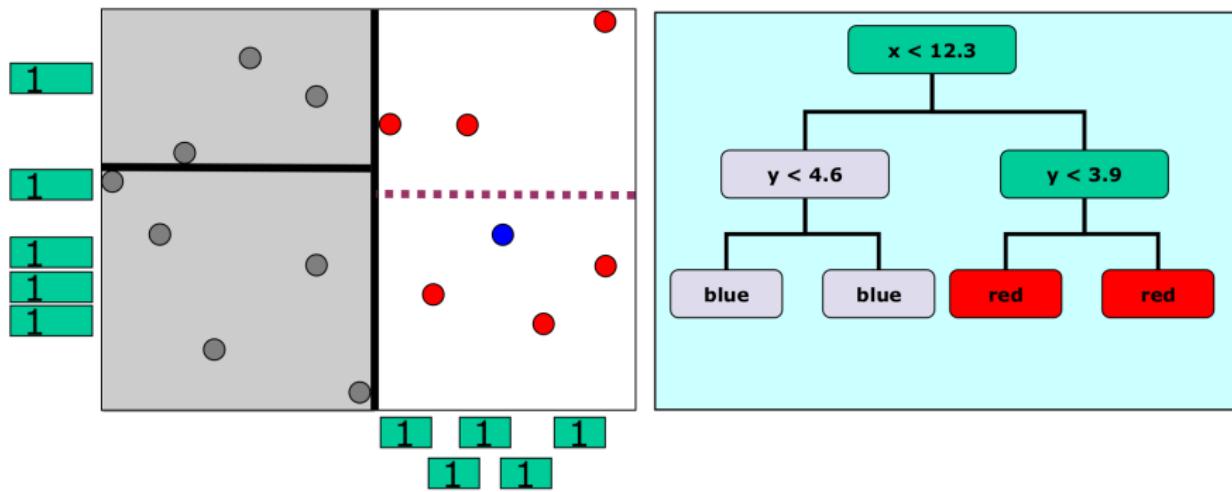
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until AER → 0



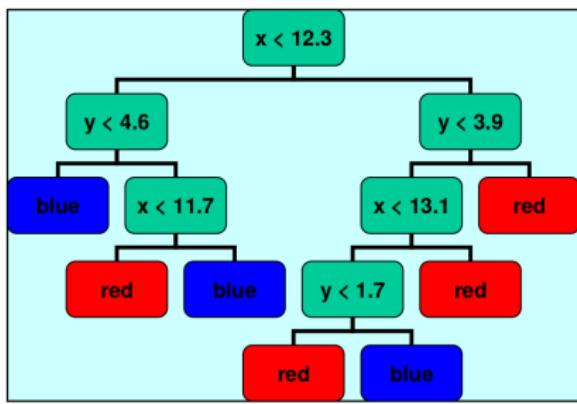
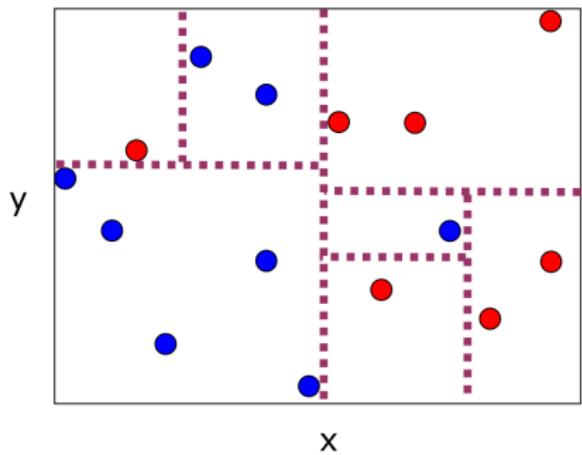
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until  $\text{AER} \rightarrow 0$



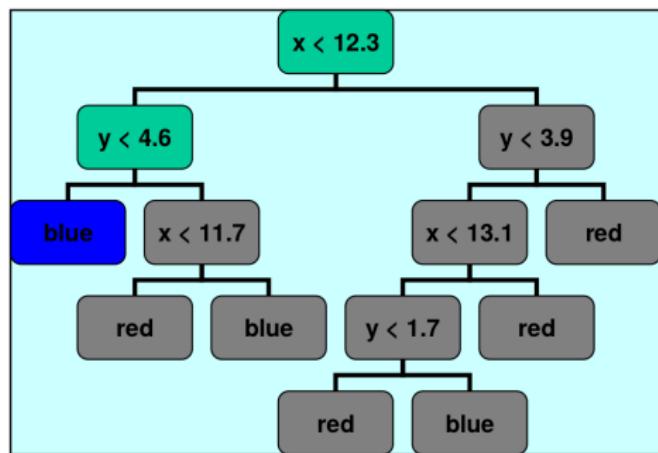
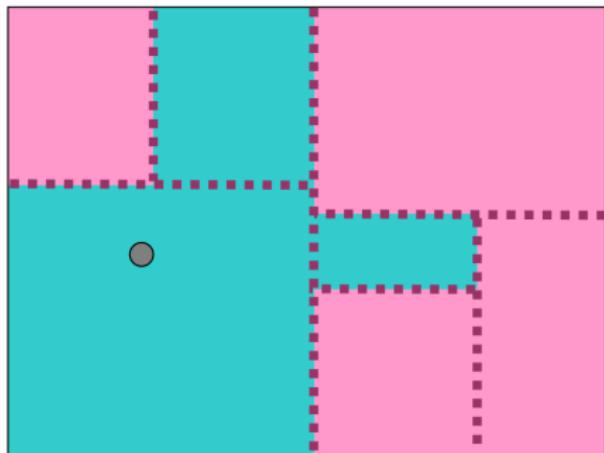
# Decision tree learning

- 1) We compute a measure for *each possible split* in each feature  
→ here **absolute error rate** (AER)
- 2) Recursively repeat step (1) for each subspace until  $\text{AER} \rightarrow 0$



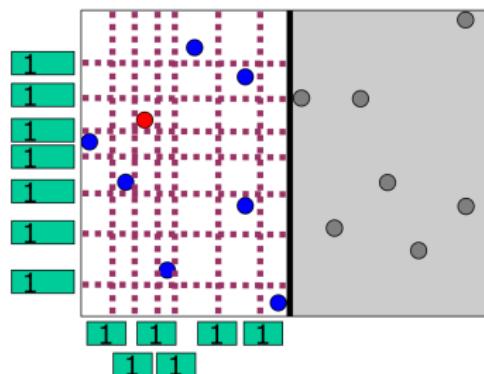
# Decision tree classification

## 3) Classification

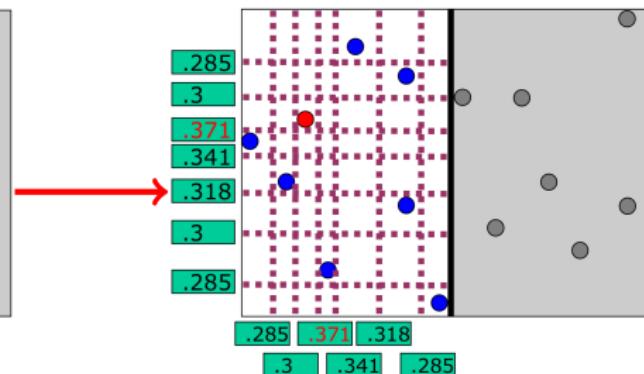


# Decision tree improvements I

- Use more sophisticated split measures
  - ▶ *Information gain*  $\leftrightarrow$  (im-)purity of splitted sub-sets
  - ▶ Gini index
- Pruning



Absolute error rate



Information gain

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample (*—bootstrapping*)
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample (=bootstrapping)
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample  
*(=bootstrapping)*
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample  
*(=bootstrapping)*
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample  
*(=bootstrapping)*
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample  
*(=bootstrapping)*
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample (=bootstrapping)
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample (=bootstrapping)
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Decision tree improvements II

## Random forest

- Ensemble of DTs
- For each tree use:
  - ▶ Random sub-sample (=bootstrapping)
  - ▶ Random number of the original features  
→ large number of rather shallow trees
- Classify data by majority voting of individual trees

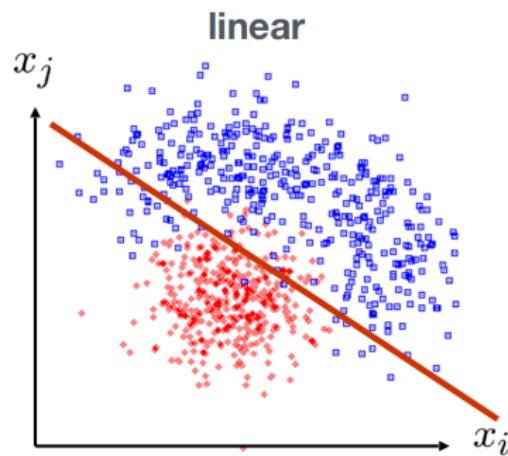
## Boosted DT

- Sequential ensemble of evolving DTs
- Assign weights to training samples → initially equal weights
- Adjust weights & give more importance to misclassified samples
- Subsequent predictors thus should focus on those

# Linear cuts

More degrees of freedom than rectangular cut

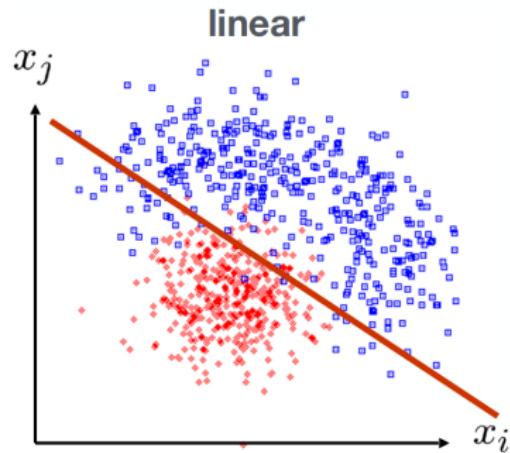
- Simple white box methods
- Very fast classification
- Can become very powerful by using *kernel trick*



# Linear cuts

More degrees of freedom than rectangular cut

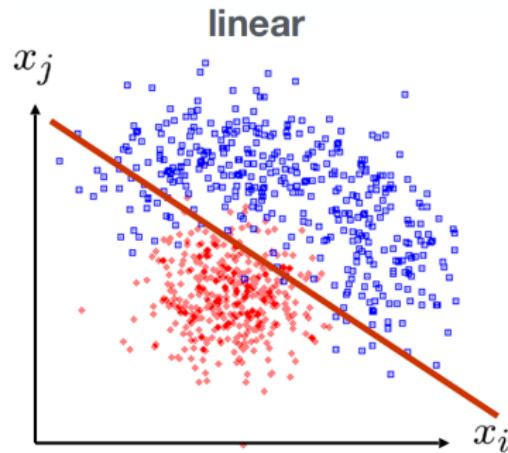
- Simple white box methods
- Very fast classification
- Can become very powerful by using *kernel trick*



# Linear cuts

More degrees of freedom than rectangular cut

- Simple white box methods
- Very fast classification
- Can become very powerful by using *kernel trick*



# Linear models

Takes a *linear function* of its inputs  $\mathbf{x} = (x_1, \dots, x_n)$  to base its decision on.

$$y : \mathbb{R}^n \rightarrow \mathbb{R} \mid \mathbf{x} \mapsto y = f(\mathbf{w} \cdot \mathbf{x}) = f\left(\sum_j w_j x_j\right)$$

$\mathbf{w}$  ... weight vector

Simplest case

$$y = f(x) = \Theta(x) = \begin{cases} 0 & x < 0 & \text{signal} \\ 1 & x \geq 0 & \text{background} \end{cases}$$

→ Function can be approximated by **single layer perceptron**

# Linear models

Takes a *linear function* of its inputs  $\mathbf{x} = (x_1, \dots, x_n)$  to base its decision on.

$$y : \mathbb{R}^n \rightarrow \mathbb{R} \mid \mathbf{x} \mapsto y = f(\mathbf{w} \cdot \mathbf{x}) = f\left(\sum_j w_j x_j\right)$$

$\mathbf{w}$  ... weight vector

Simplest case

$$y = f(x) = \Theta(x) = \begin{cases} 0 & x < 0 & \text{signal} \\ 1 & x \geq 0 & \text{background} \end{cases}$$

→ Function can be approximated by **single layer perceptron**

# Linear models

Takes a *linear function* of its inputs  $\mathbf{x} = (x_1, \dots, x_n)$  to base its decision on.

$$y : \mathbb{R}^n \rightarrow \mathbb{R} \mid \mathbf{x} \mapsto y = f(\mathbf{w} \cdot \mathbf{x}) = f\left(\sum_j w_j x_j\right)$$

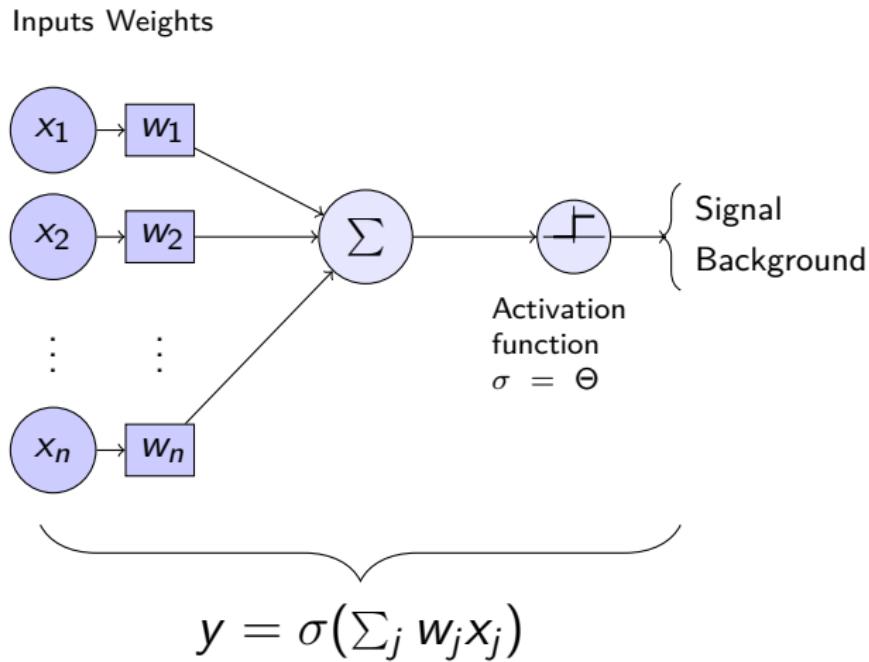
$\mathbf{w}$  ... weight vector

Simplest case

$$y = f(x) = \Theta(x) = \begin{cases} 0 & x < 0 & \text{signal} \\ 1 & x \geq 0 & \text{background} \end{cases}$$

→ Function can be approximated by **single layer perceptron**

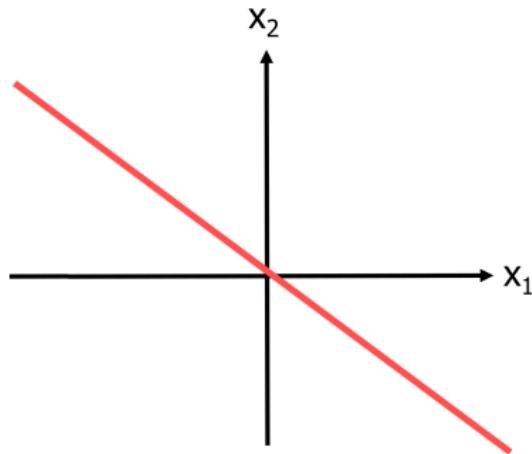
# Single layer perceptron (SLP)



# SLP training

## Algorithm

- ➊ Initialize weights  $w_{init}$
- ➋ Repeat until  $y = y_{target}$ :
  - ➌ Present training sample  $x$
  - ➍ Predict sample label  
 $y = \Theta(xw_{init})$  and compute error  $\Delta = y_{target} - y$
  - ➎ If  $\Delta \neq 0 \rightarrow$  update weights  
 $w' = w + \alpha \Delta x$

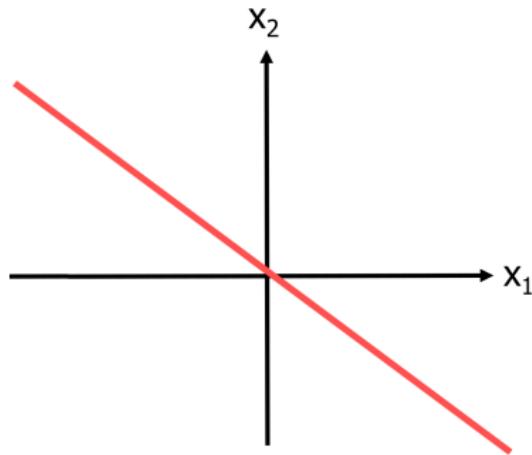


$$y = w_1 x_1 + w_2 x_2$$

# SLP training

## Algorithm

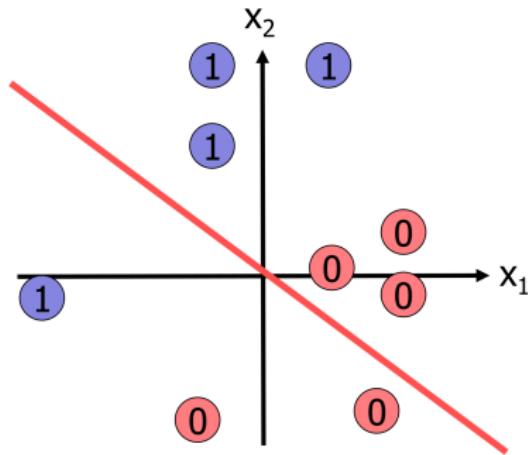
- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $\mathbf{x}$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$



# SLP training

## Algorithm

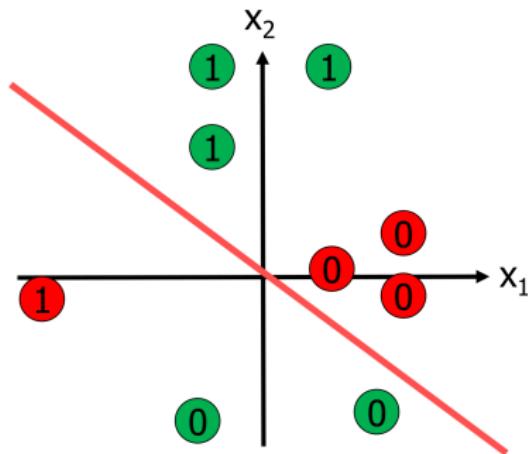
- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $\mathbf{x}$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$



# SLP training

## Algorithm

- ➊ Initialize weights  $\mathbf{w}_{init}$
- ➋ Repeat until  $y = y_{target}$ :
  - ➌ Present training sample  $\mathbf{x}$
  - ➍ Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ➎ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

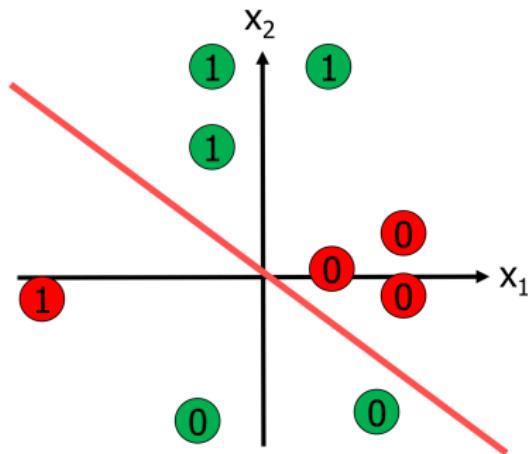


# SLP training

## Algorithm

- ➊ Initialize weights  $\mathbf{w}_{init}$
- ➋ Repeat until  $y = y_{target}$ :
  - ➌ Present training sample  $\mathbf{x}$
  - ➍ Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ➎ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

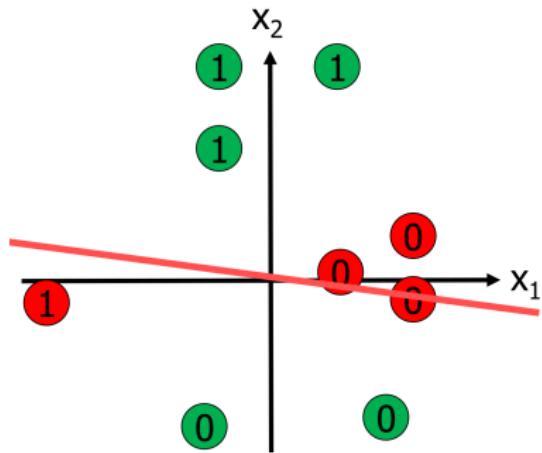


# SLP training

## Algorithm

- ➊ Initialize weights  $\mathbf{w}_{init}$
- ➋ Repeat until  $y = y_{target}$ :
  - ➌ Present training sample  $x$
  - ➍ Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ➎ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

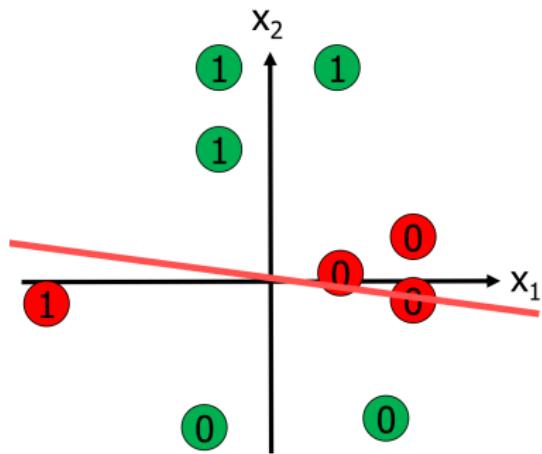


# SLP training

## Algorithm

- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $\mathbf{x}$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

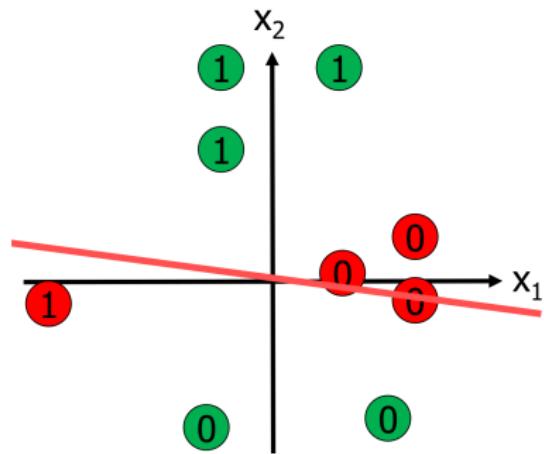


# SLP training

## Algorithm

- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $x$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute  
 error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

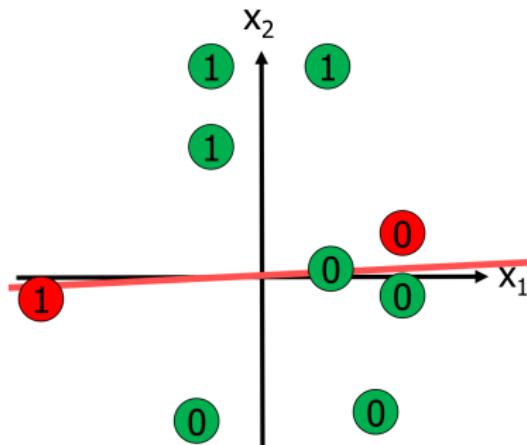


# SLP training

## Algorithm

- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $\mathbf{x}$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

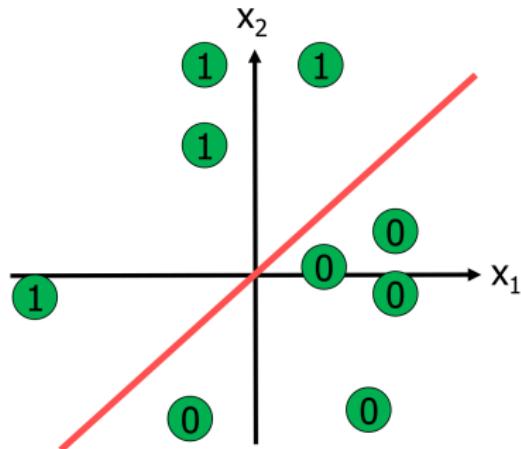


# SLP training

## Algorithm

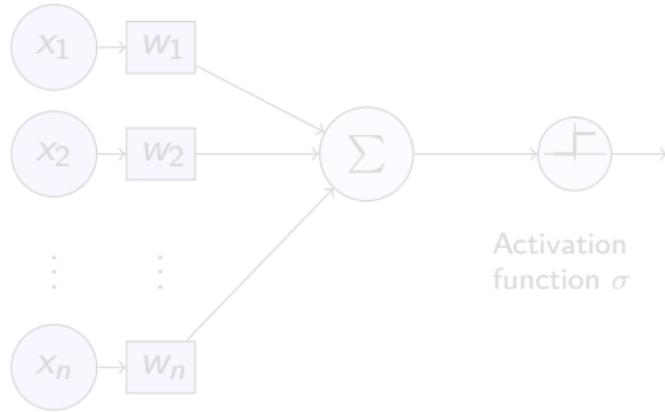
- ① Initialize weights  $\mathbf{w}_{init}$
- ② Repeat until  $y = y_{target}$ :
  - ① Present training sample  $\mathbf{x}$
  - ② Predict sample label  
 $y = \Theta(\mathbf{x}\mathbf{w}_{init})$  and compute error  $\Delta = y_{target} - y$
  - ③ If  $\Delta \neq 0 \rightarrow$  update weights  
 $\mathbf{w}' = \mathbf{w} + \alpha \Delta \mathbf{x}$

$\alpha \dots$  learning rate

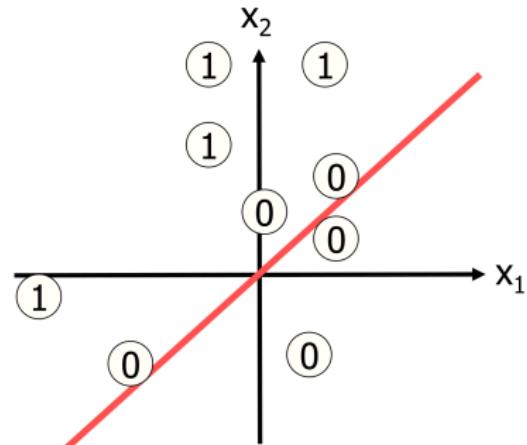


# SLP bias term

Inputs Weights

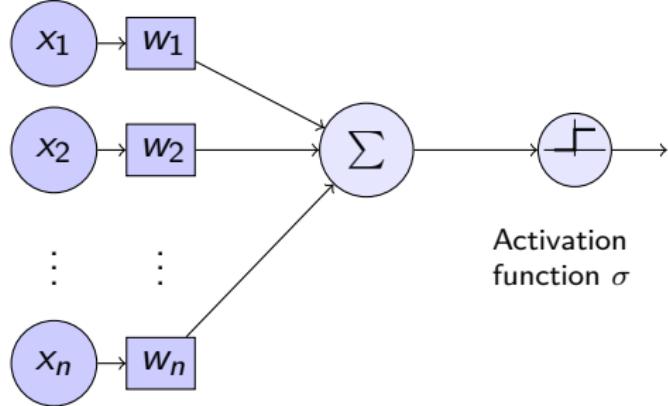


Activation  
function  $\sigma$

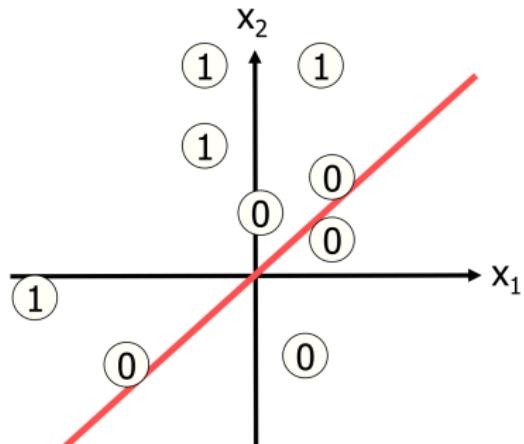


# SLP bias term

Inputs Weights

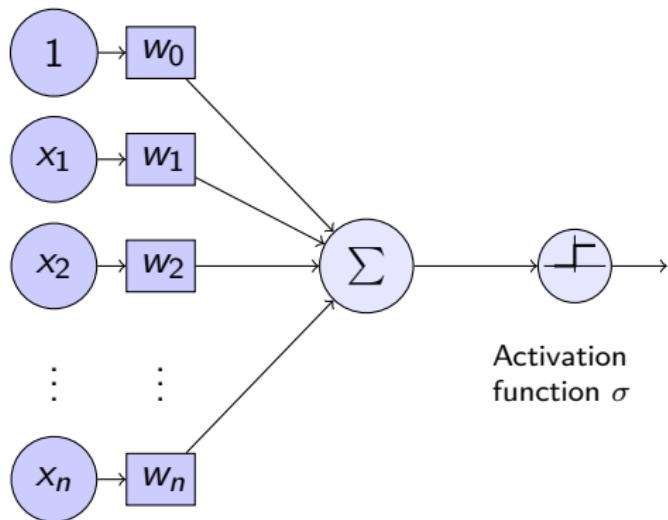


Activation  
function  $\sigma$

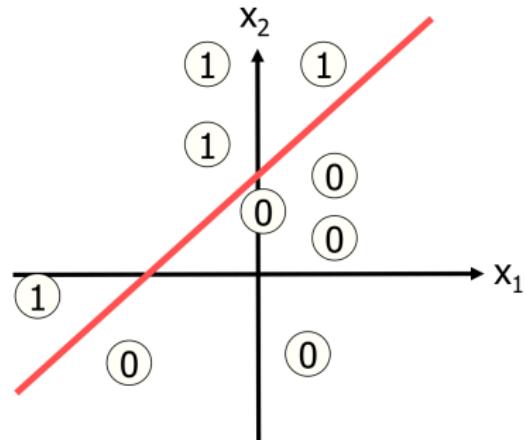


# SLP bias term

Inputs Weights

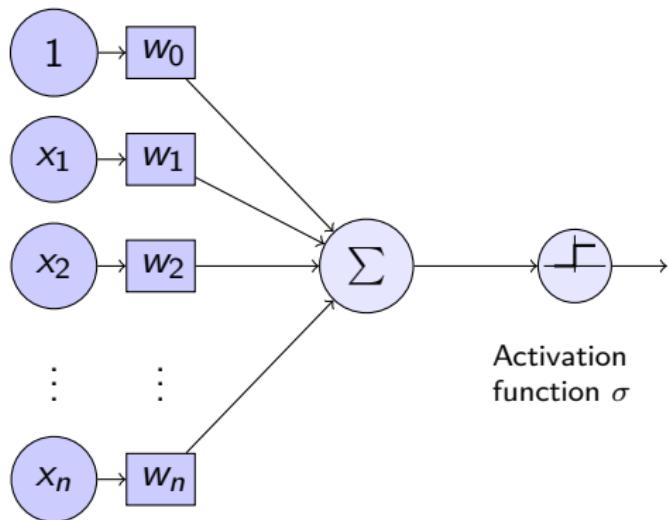


Activation  
function  $\sigma$

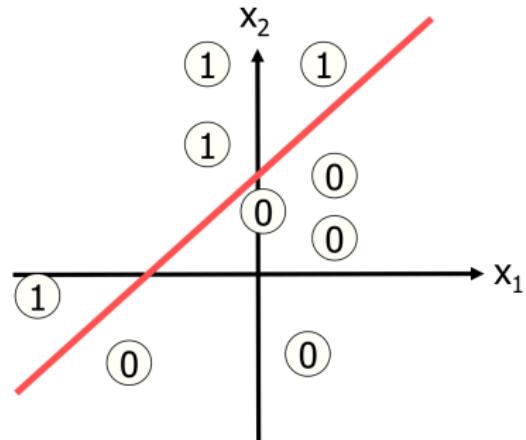


# SLP bias term

Inputs Weights

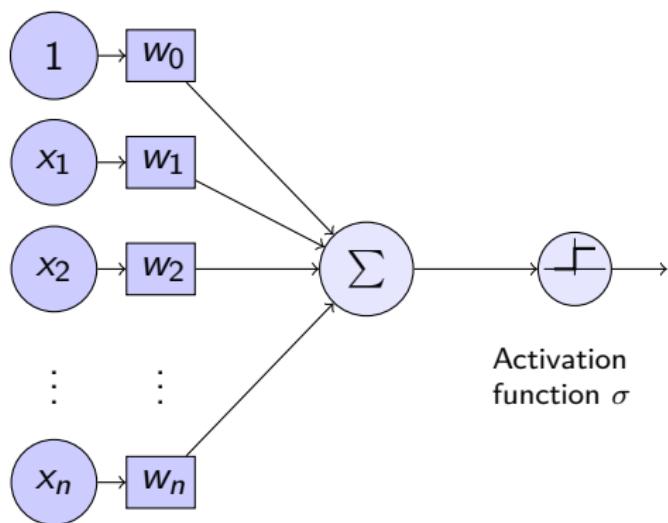


Activation  
function  $\sigma$

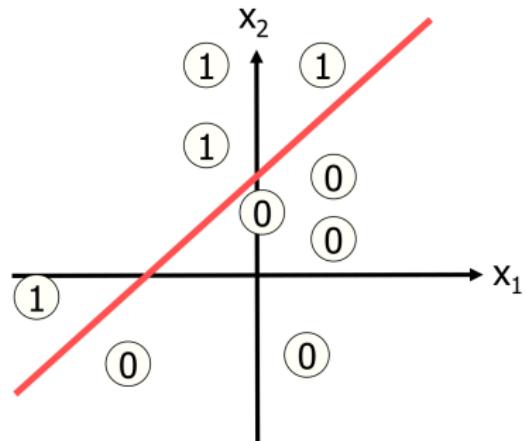


# SLP bias term

Inputs Weights



Activation  
function  $\sigma$

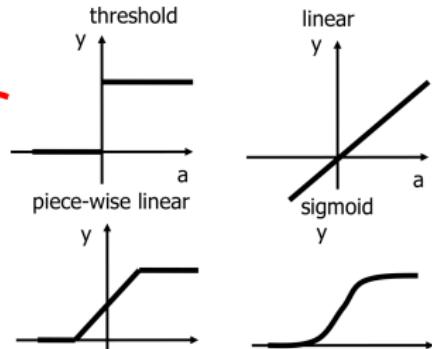
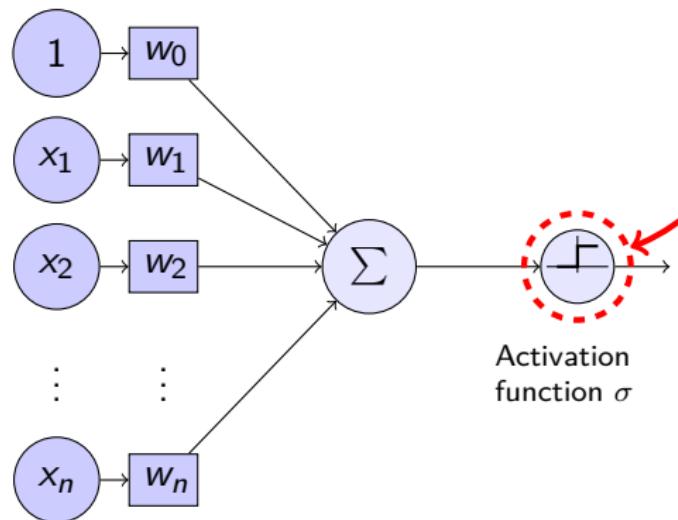


Bias weight  $w_0$  learned just as the other weights

$$y = \sigma(w_0 + \sum_{j=1}^n w_j x_j)$$

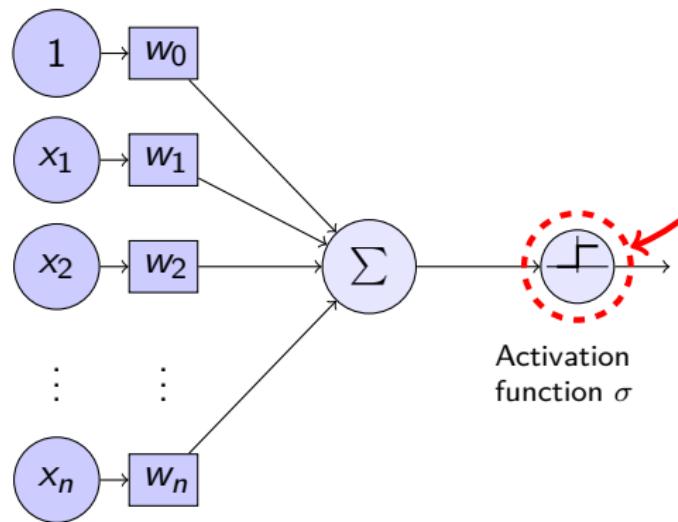
# SLP activation functions

Inputs Weights

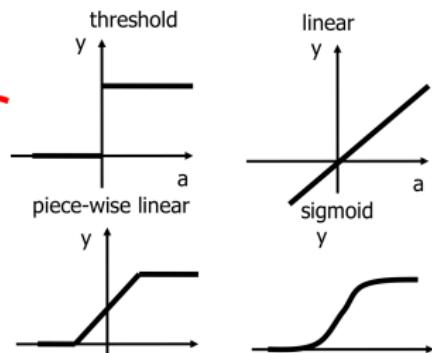


# SLP activation functions

Inputs Weights

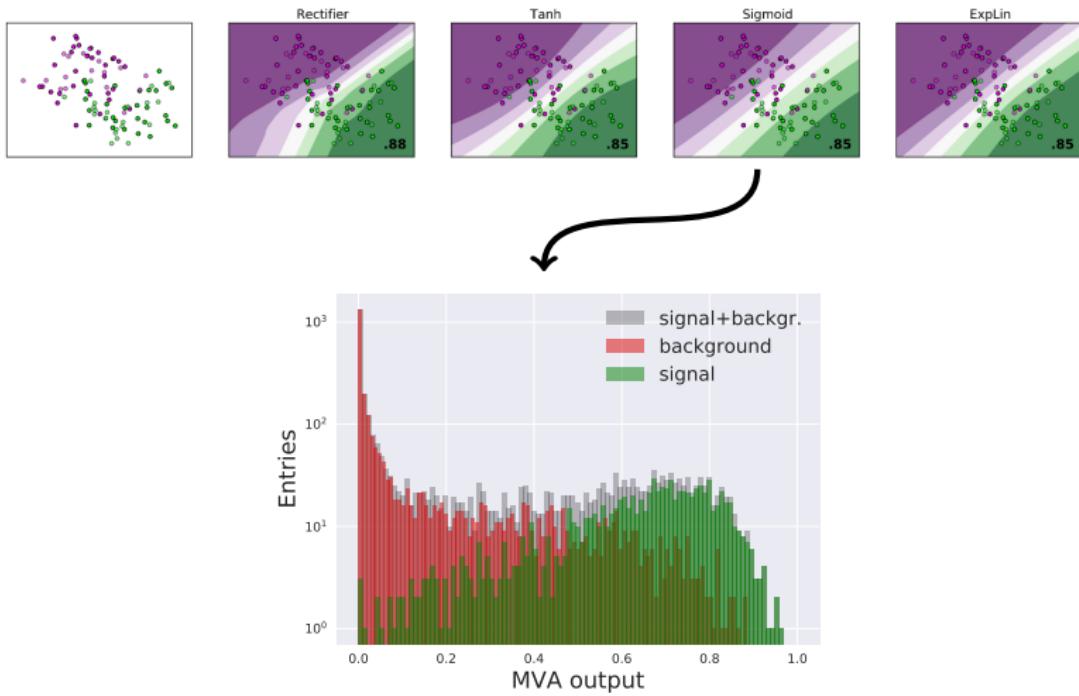


Activation  
function  $\sigma$



Often used:  $\text{sigmoid} \rightarrow \text{output} \in (0, 1)$

# SLP improvement I Non-linear activation functions output



# Improving linear methods

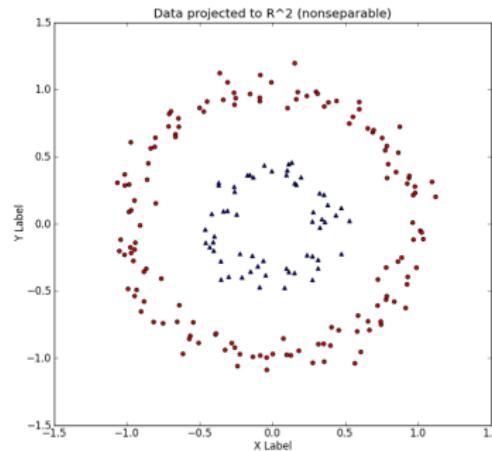
## Kernel trick

Map data to a higher dimensional space where linear hyperplane can again be found

# Improving linear methods

## Kernel trick

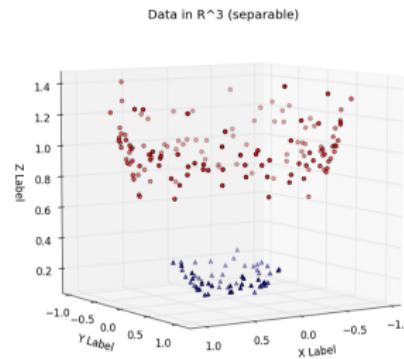
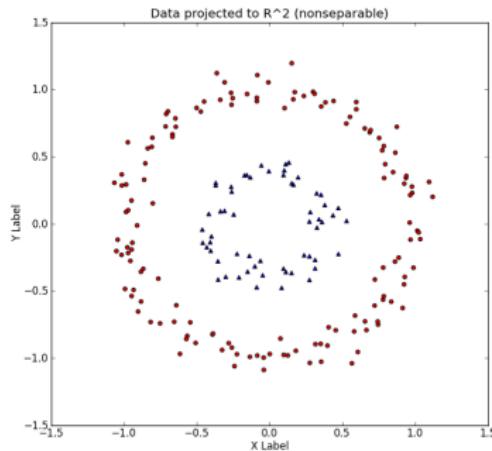
Map data to a higher dimensional space where linear hyperplane can again be found



# Improving linear methods

## Kernel trick

Map data to a higher dimensional space where linear hyperplane can again be found



# Non-linear models

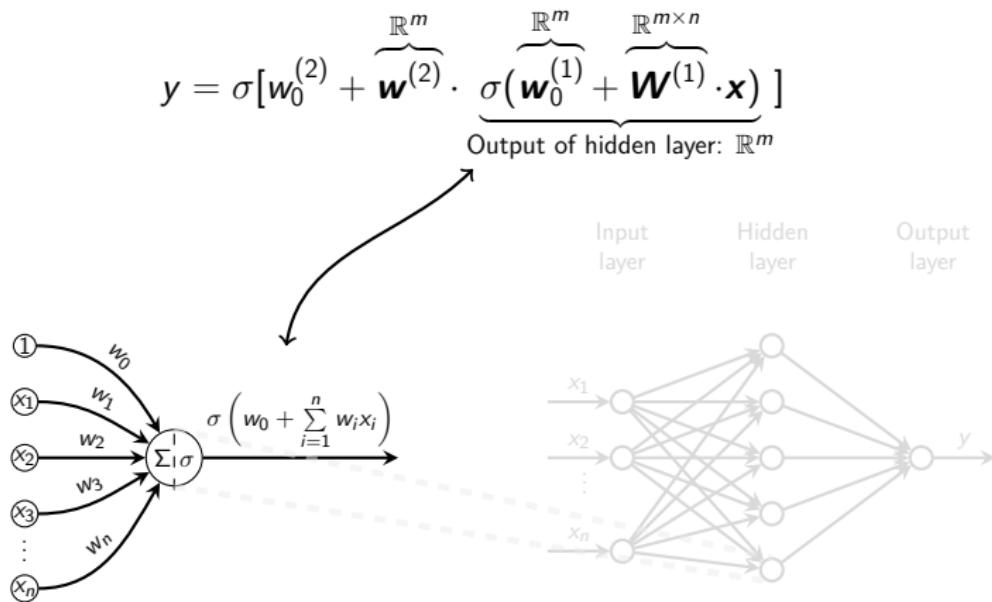
## Assumption

Take *kernel trick* idea: map inputs  $\mathbf{x} = (x_1, \dots, x_n)$  to higher dimensional space  $\mathbb{R}^m$  with  $m > n$  where linear separation is possible:

$$\begin{aligned} y : \mathbb{R}^n &\rightarrow \mathbb{R}^m \rightarrow \mathbb{R} = \sigma \left( w_0^2 + \sum_{i=1}^m \left[ w_i^{(2)} + \sigma \left( w_{0,i}^{(1)} + \sum_{k=1}^n w_{k,i}^{(1)} x_k \right) \right] \right) \\ &= \sigma \left[ w_0^{(2)} + \mathbf{w}^{(2)} \cdot \sigma(\mathbf{w}_0^{(1)} + \mathbf{W}^{(1)} \cdot \mathbf{x}) \right] \end{aligned}$$

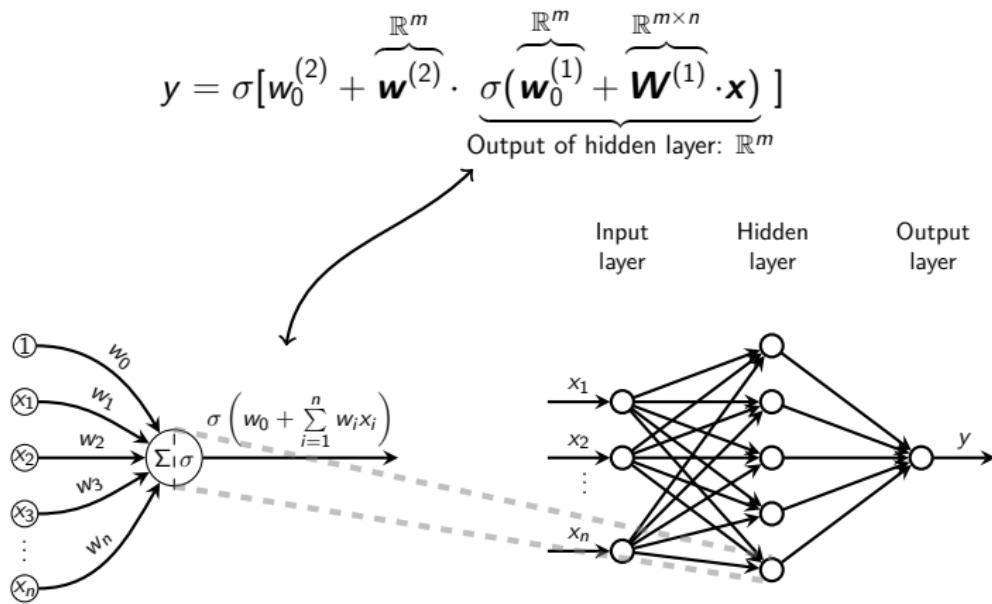
→  $m$  single layer perceptrons

# Visualisation



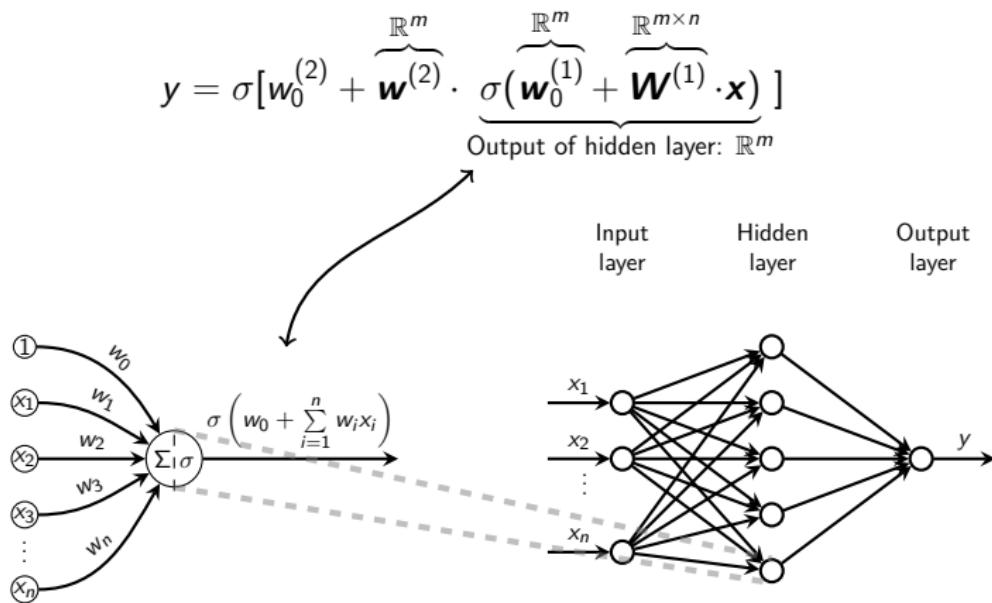
→ The hidden layer *learns* a representation so that the data is linearly separable

# Visualisation



→ The hidden layer *learns* a representation so that the data is linearly separable

# Visualisation



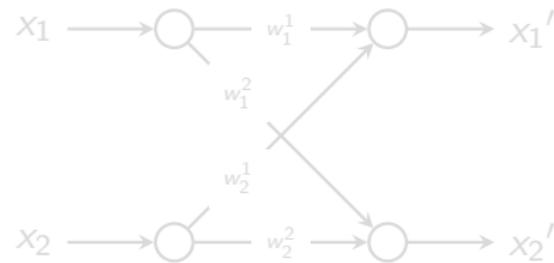
→ The hidden layer *learns* a representation so that the data is linearly separable

# Visualising the hidden layer

## Effects of adding a hidden layer

Input vector  $\mathbf{x}$  undergoes *linear transformation* by multiplication with weight matrix  $\mathbf{W}$ . The bias term  $\mathbf{w}_0$  causes a *translation*.

$$y = \sigma[w_0^{(2)} + \underbrace{\mathbf{w}^{(2)} \cdot \sigma(w_0^{(1)} + \mathbf{W}^{(1)} \cdot \mathbf{x})}_{\tilde{\mathbf{x}}}]$$

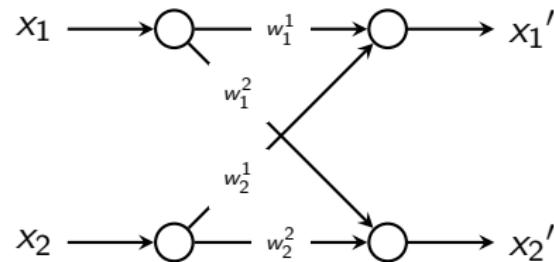


# Visualising the hidden layer

## Effects of adding a hidden layer

Input vector  $\mathbf{x}$  undergoes *linear transformation* by multiplication with weight matrix  $\mathbf{W}$ . The bias term  $\mathbf{w}_0$  causes a *translation*.

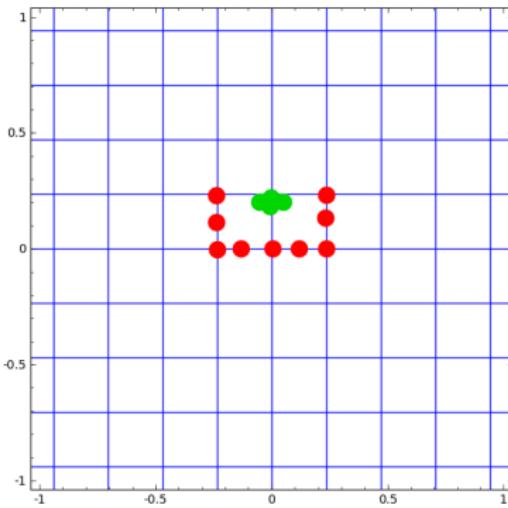
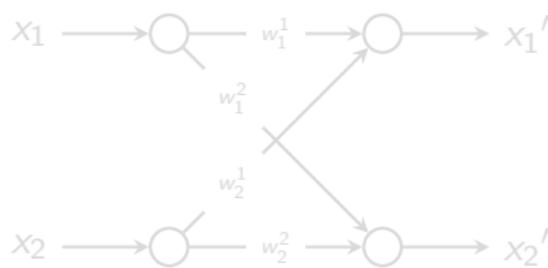
$$\tilde{\mathbf{x}} = \sigma(\mathbf{w}_0^{(1)} + \mathbf{W}^{(1)} \cdot \mathbf{x}) \text{ with e.g. } \mathbf{W} = \begin{pmatrix} w_1^1 & w_1^2 \\ w_2^1 & w_2^2 \end{pmatrix}$$



# Visualising the hidden layer

- ① Linear transformation:  $Wx$
- ② Translation:  $w_0$
- ③ Non-linear activation:  $\sigma$

$$\tilde{x} = x$$



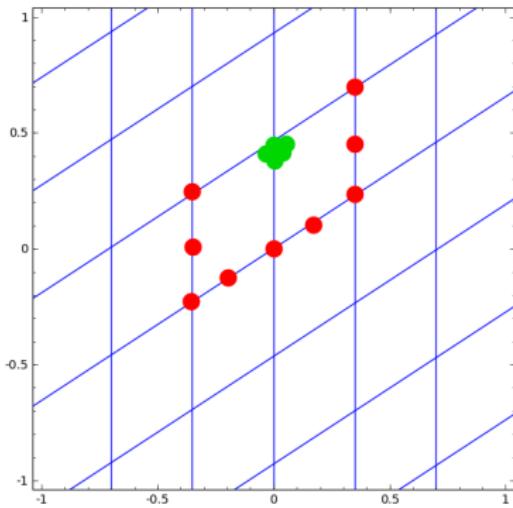
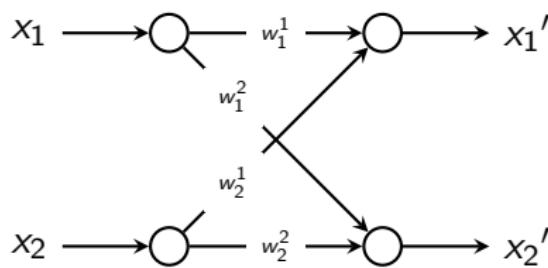
# Visualising the hidden layer

① Linear transformation:  $\mathbf{Wx}$

② Translation:  $\mathbf{w}_0$

③ Non-linear activation:  $\sigma$

$$\tilde{\mathbf{x}} = (\mathbf{W}^{(1)} \cdot \mathbf{x})$$



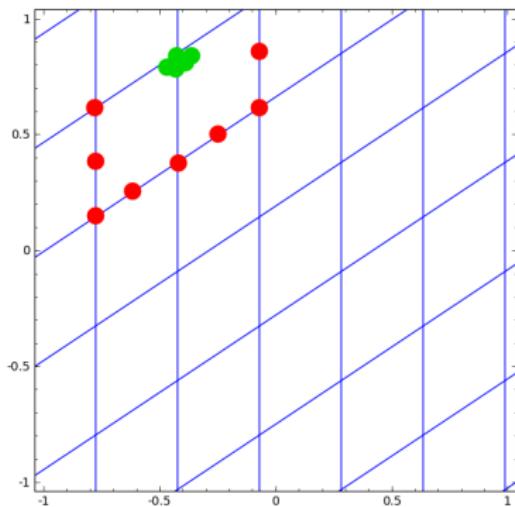
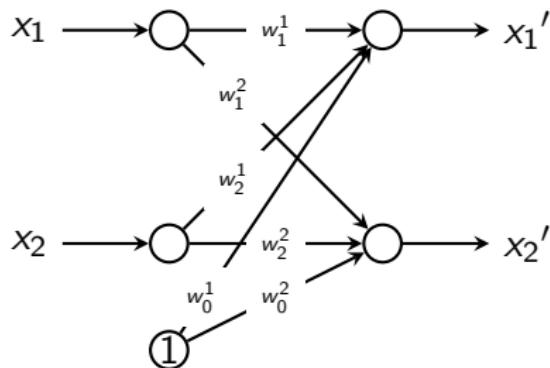
# Visualising the hidden layer

① Linear transformation:  $\mathbf{Wx}$

② Translation:  $\mathbf{w}_0$

③ Non-linear activation:  $\sigma$

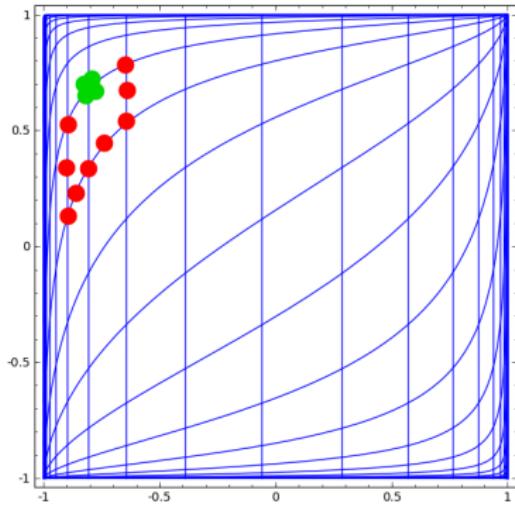
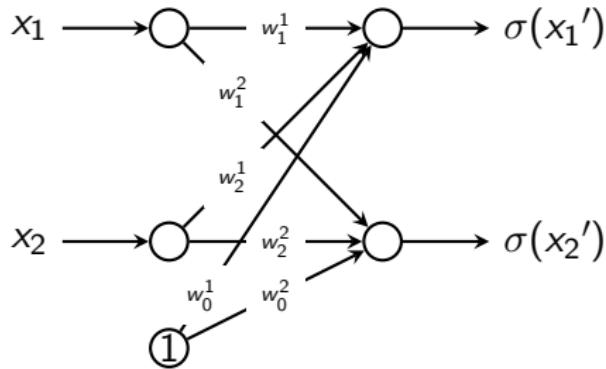
$$\tilde{\mathbf{x}} = (\mathbf{w}_0^{(1)} + \mathbf{W}^{(1)} \cdot \mathbf{x})$$



# Visualising the hidden layer

- ➊ Linear transformation:  $\mathbf{W}\mathbf{x}$
- ➋ Translation:  $\mathbf{w}_0$
- ➌ Non-linear activation:  $\sigma$

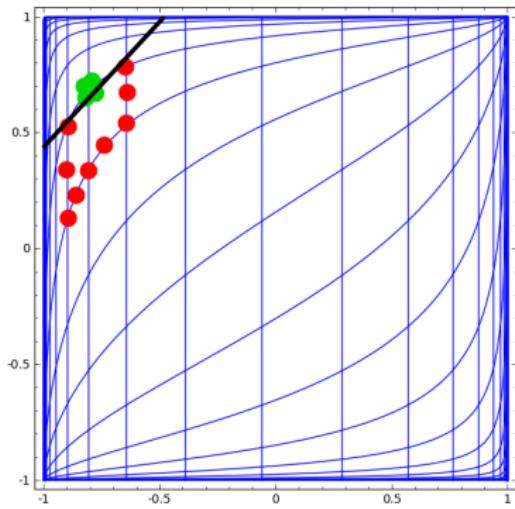
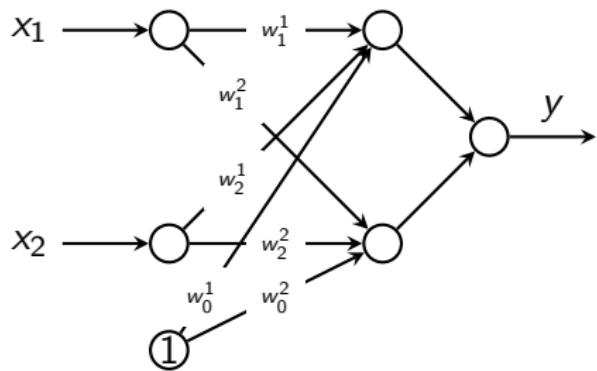
$$\tilde{\mathbf{x}} = \sigma(\mathbf{w}_0^{(1)} + \mathbf{W}^{(1)} \cdot \mathbf{x})$$



# Visualising the hidden layer

- ➊ Linear transformation:  $\mathbf{Wx}$
- ➋ Translation:  $\mathbf{w}_0$
- ➌ Non-linear activation:  $\sigma$

$$y = \sigma(w_0^{(2)} + \mathbf{w}^{(2)} \cdot \tilde{\mathbf{x}})$$



# Update weights via backpropagation

## Backpropagation

Adapt weights by going backwards in the network

- Initialize weights
- Evaluate  $y_{predict} = f(\mathbf{x})$  and calculate  $\Delta = (y_{true} - y_{predict})$
- Update weights: Weight adaption often denoted via loss  $L \propto \Delta$   
→ minimize loss → find best fit

# Update weights via backpropagation

## Backpropagation

Adapt weights by going backwards in the network

- Initialize weights
- Evaluate  $y_{predict} = f(\mathbf{x})$  and calculate  $\Delta = (y_{true} - y_{predict})$
- Update weights: Weight adaption often denoted via loss  $L \propto \Delta$   
→ minimize loss → find best fit

# Update weights via backpropagation

## Backpropagation

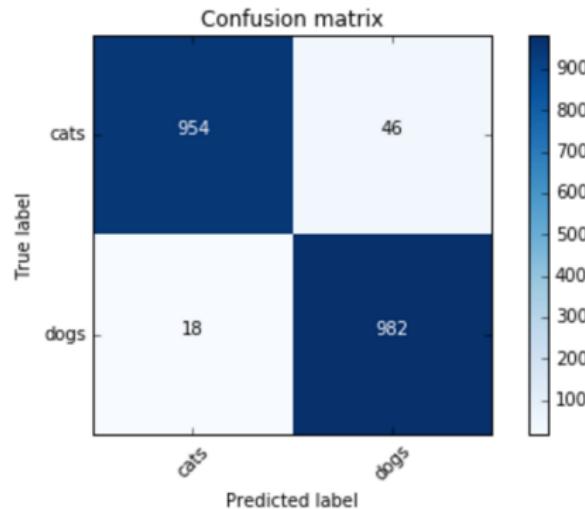
Adapt weights by going backwards in the network

- Initialize weights
- Evaluate  $y_{predict} = f(\mathbf{x})$  and calculate  $\Delta = (y_{true} - y_{predict})$
- Update weights: Weight adaption often denoted via loss  $L \propto \Delta$   
→ minimize loss → find best fit

# Quantify MVA output

## Evaluate classifier performance

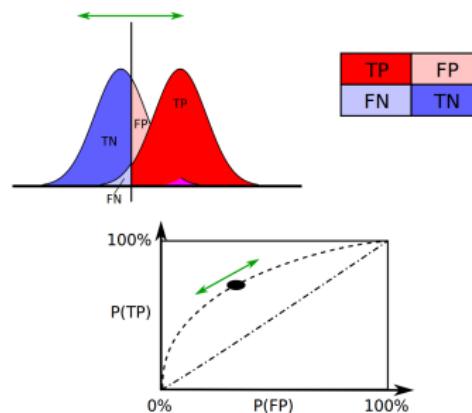
- Compare predicted class labels with the true ones → confusion matrix



# Quantify MVA output

## Evaluate classifier performance

- ROC curve: continuously scan y & plot background acceptance (FPR) vs. signal efficiency (TPR)



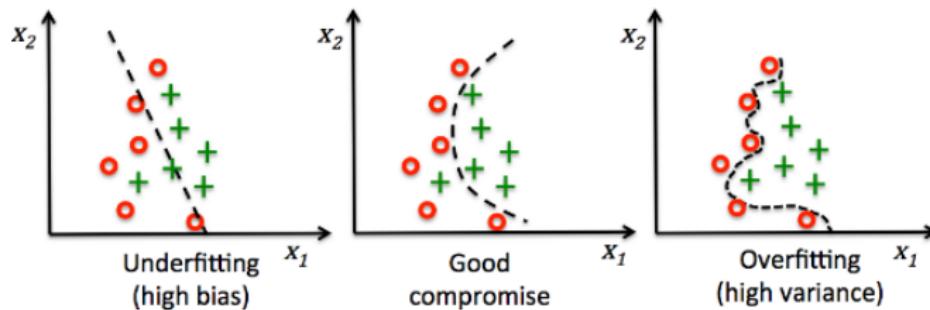
$$\text{FPR} = \text{FP}/(\text{FP} + \text{TN})$$

$$\text{TPR} = \text{TP}/(\text{TP} + \text{FN})$$

ROC-AUC: probability that classifier ranks randomly chosen positive sample higher than randomly chosen negative one

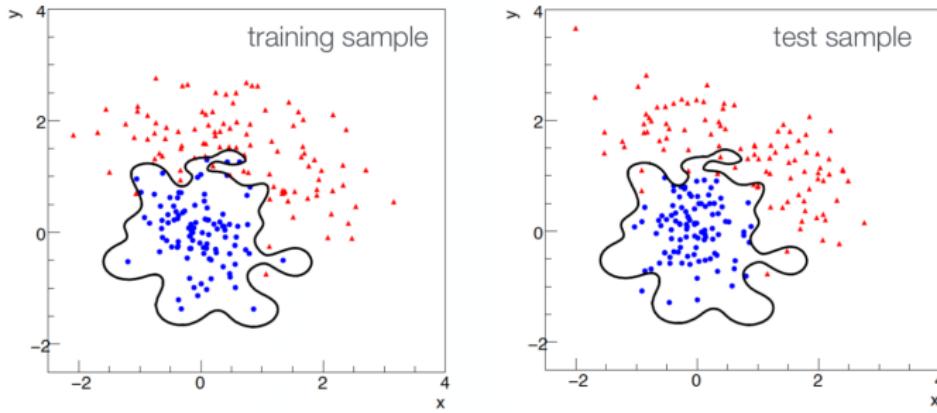
# Bias-variance tradeoff

- Small variance: Classifiers with low degrees of freedom are less prone to statistical fluctuations  
→ different training samples result in similar classification boundaries
- However: if data contain features that a model with few degrees of freedom cannot describe, a bias is introduced

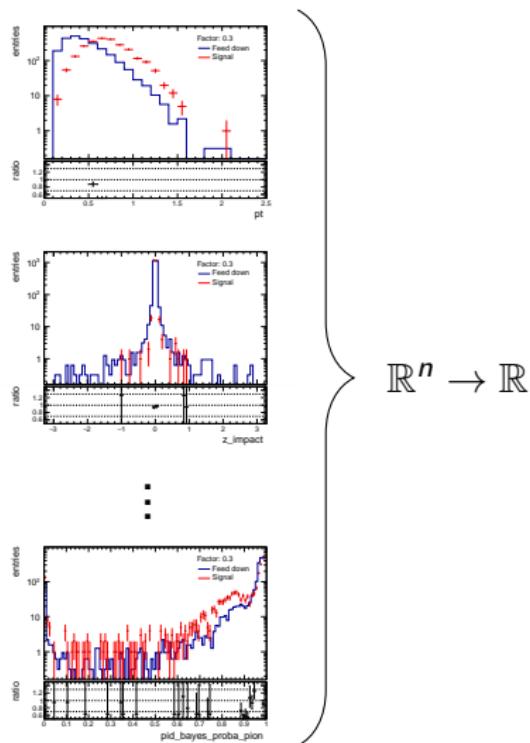


# Bias-variance tradeoff

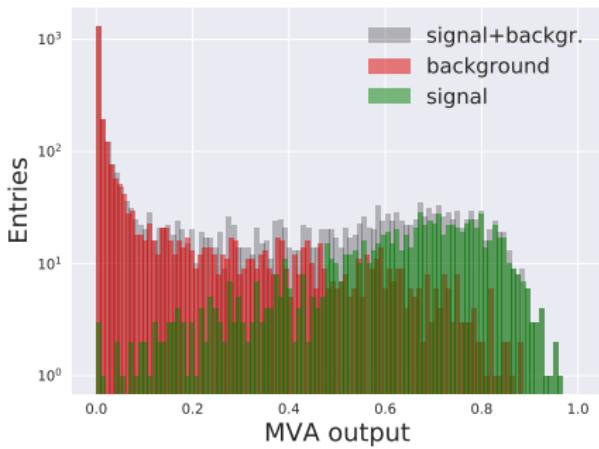
- Small variance: Classifiers with low degrees of freedom are less prone to statistical fluctuations  
→ different training samples result in similar classification boundaries
- However: if data contain features that a model with few degrees of freedom cannot describe, a bias is introduced



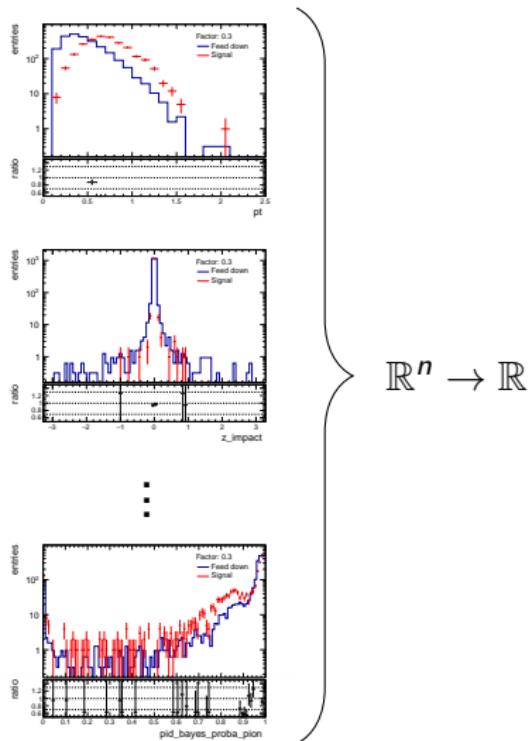
# Results



$$\mathbb{R}^n \rightarrow \mathbb{R}$$

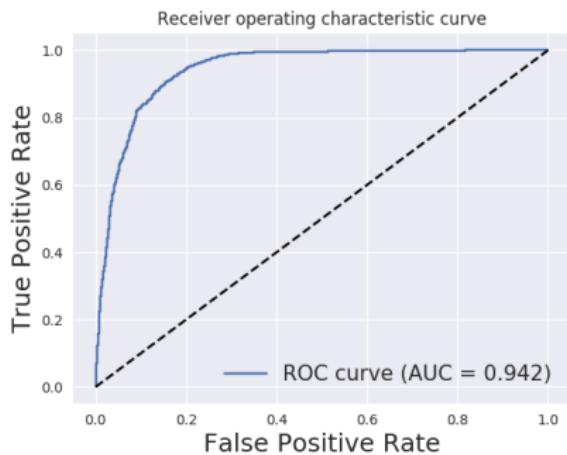


# Results

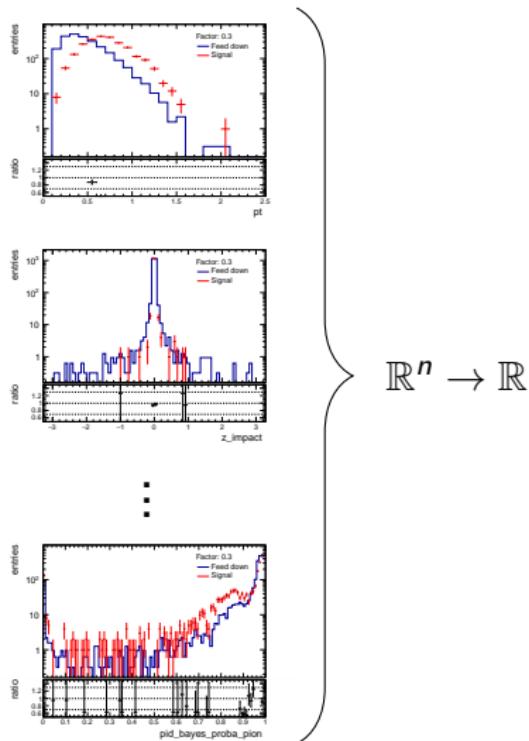


$$\mathbb{R}^n \rightarrow \mathbb{R}$$

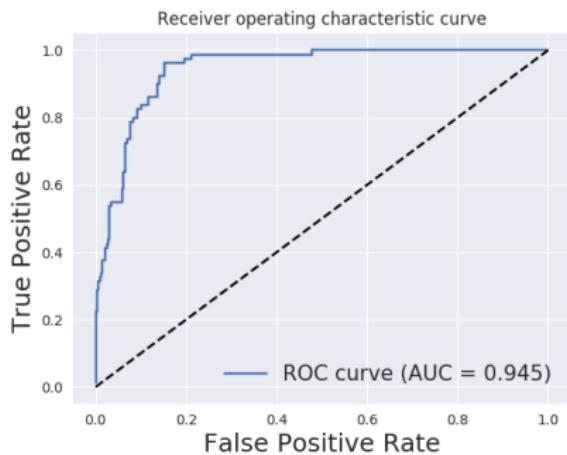
## ROC train



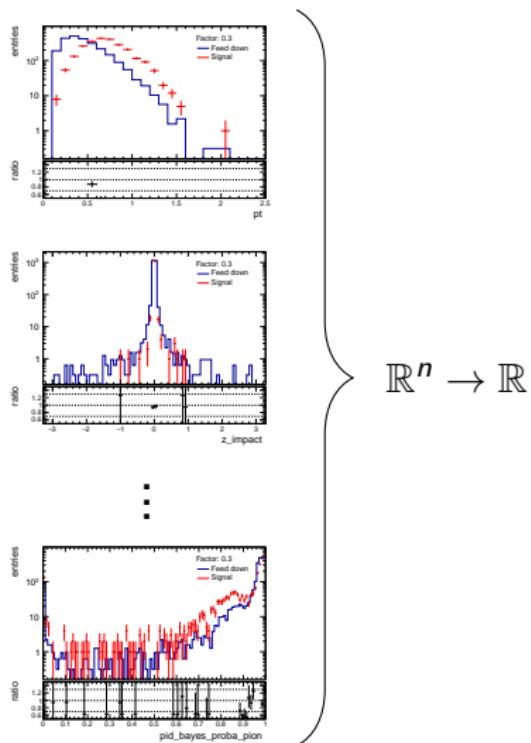
# Results



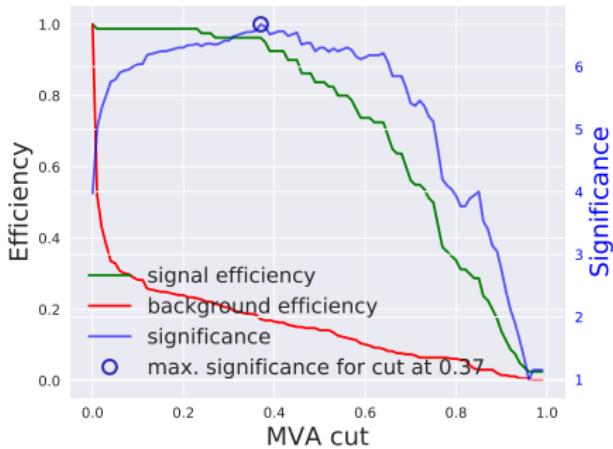
## ROC test



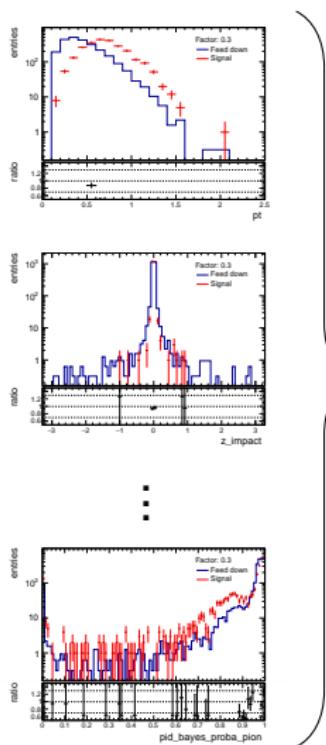
# Results



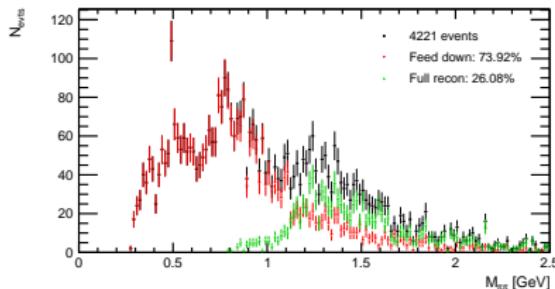
$$\text{Significance } S = N_{sig} / \sqrt{N_{sig} + N_{BG}}$$



# Results



no MVA cut



with MVA cut

