# Machine Learning

**Sergei** **Gleyzer**

**PART** **I**

**TAE 2017 Lectures**

**Sep. 4, 2017**

# **Outline**

- **What is Machine Learning**
- **in Particle Physics**
- **in Theory**
- **in Practice**

# **Machine Learning Basics**

# Machine Learning

## What is Machine Learning?

- Study of algorithms that
  improve their <u>performance</u> **P**
  for a given <u>task</u> **T**
  with more <u>experience</u> **E**

**Sample tasks: identifying faces, Higgs bosons**

# **Machine Learning**

## **General Approach:**

Given **training** data $T_D$ = {y, **x**} = (y,x)$_1$...(y,x)$_N$,

**function space** {f} and a
**constraint** on these functions

Teach a machine to learn the **mapping** y = f(x)

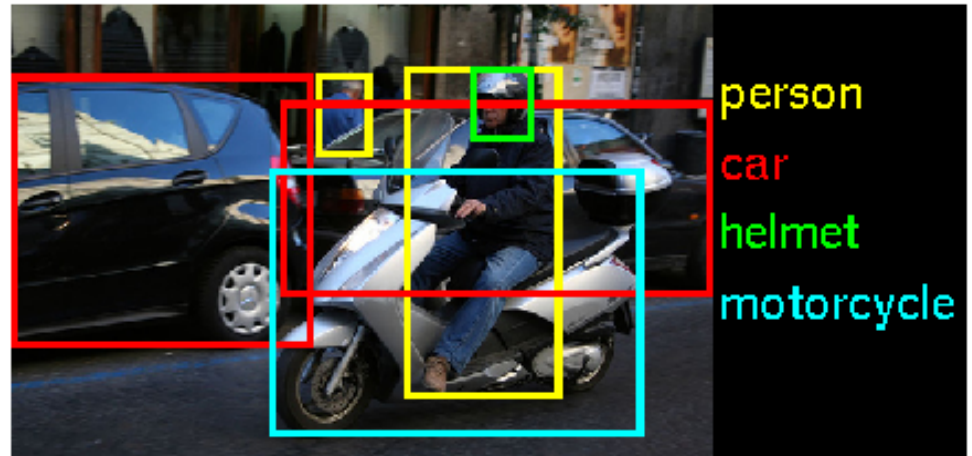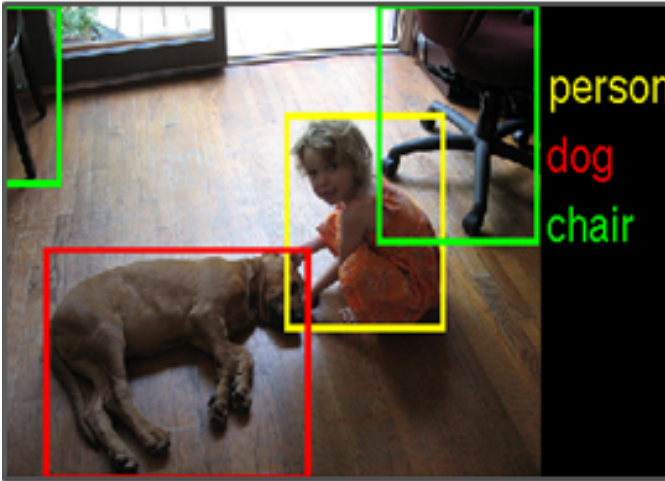# In Computer Science

## Already the preferred approach to:

- Speech recognition, natural language processing
- Computer vision, Robot control
- Medical outcomes analysis

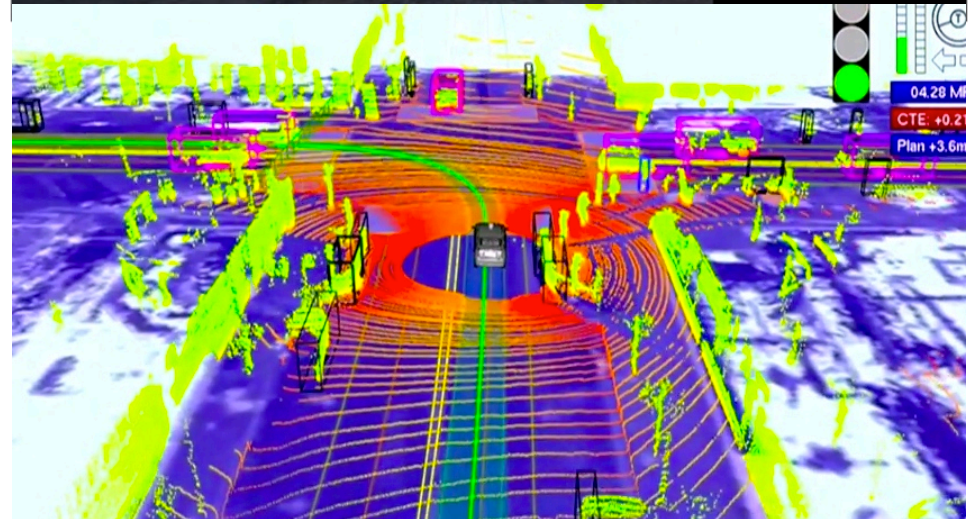## Growing fast

- Improved algorithms
- Increased data capture
- Software too complex to write by hand

# **Examples**

# Machine Learning

**Choose**

Function space $\quad\quad F = \{\, f(x, w)\,\}$

Constraint $\quad\quad\quad\quad C$

Loss function* $\quad\quad\quad L$

$$f(x, w^*) \quad\quad C(w)$$

$$F$$

**Method**

Find $f(x)$ by minimizing the empirical risk $R(w)$

$$R[f_w] = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i, w)) \quad\quad \text{subject to the constraint } C(w)$$

*The loss function measures the cost of choosing badly

# **Machine Learning**

Many methods (e.g., neural networks, boosted decision trees, rule-based systems, random forests,…) use the

quadratic loss

$$L(y, f(x, w)) = [y - f(x, w)]^2$$

and choose $f(x, w^*)$ by minimizing the

***constrained*** mean square empirical risk

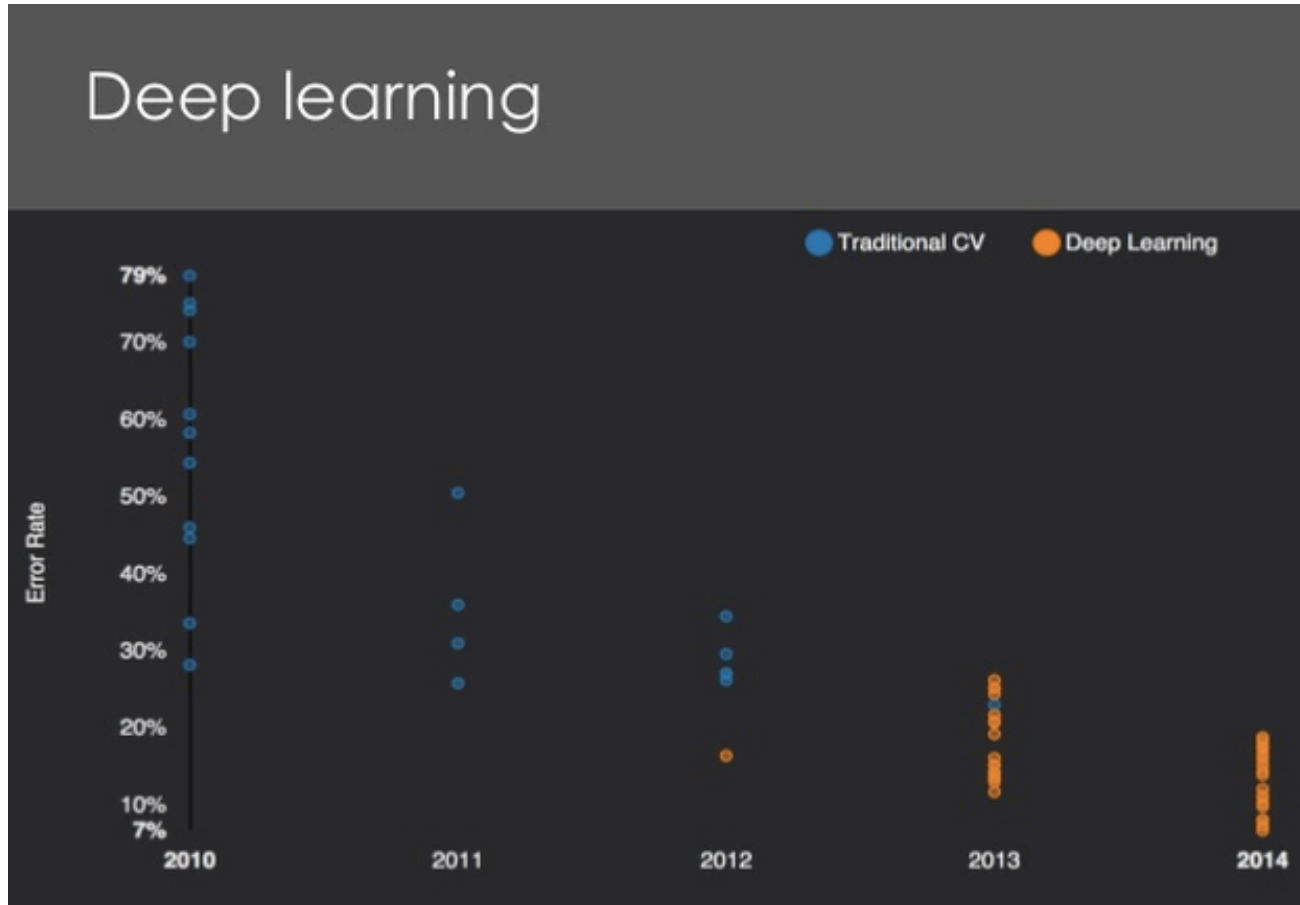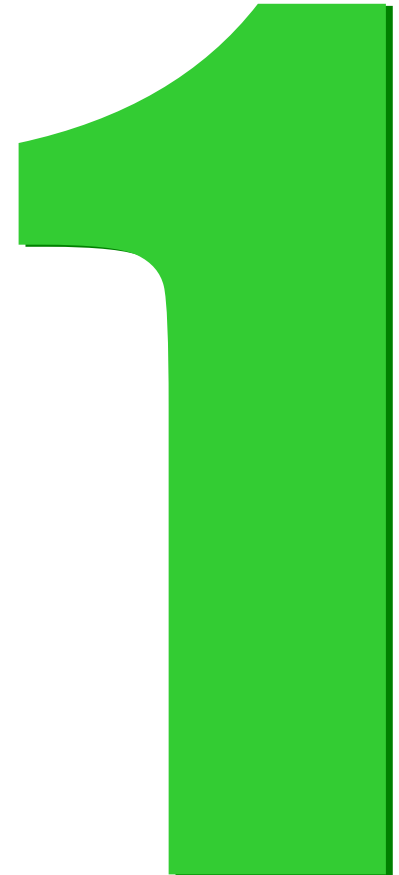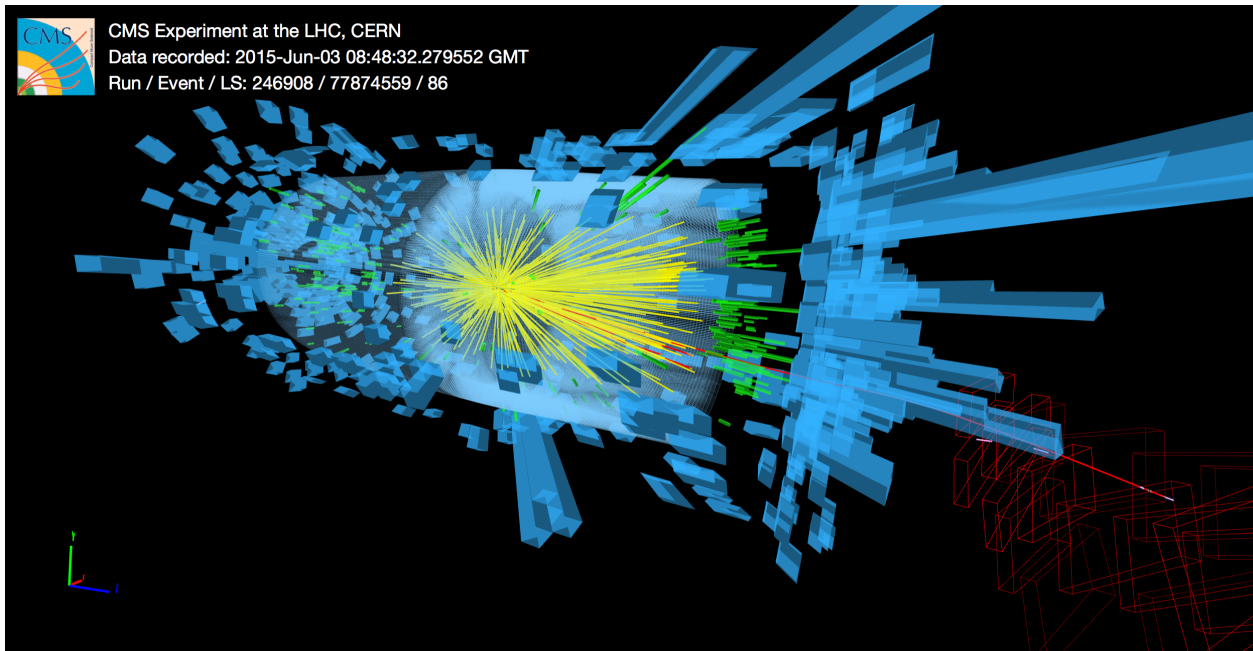$$R[f_w] = \frac{1}{N} \sum_{i=1}^{N} [y_i - f(x_i, w)]^2 + C(w)$$

# History

**1950s:** First methods invented

**1960-80s:** Slow growth, focus on knowledge

**1990s:** Growth of computing power, new learning methods, data-centric

**2000-10s:** Wider use in research and industry

**2010s:** Learning improvement, dedicated hardware, deeper learning

# Diving Deeper



**Huge Progress**

**In Particle Physics**

# Higgs Boson Discovery
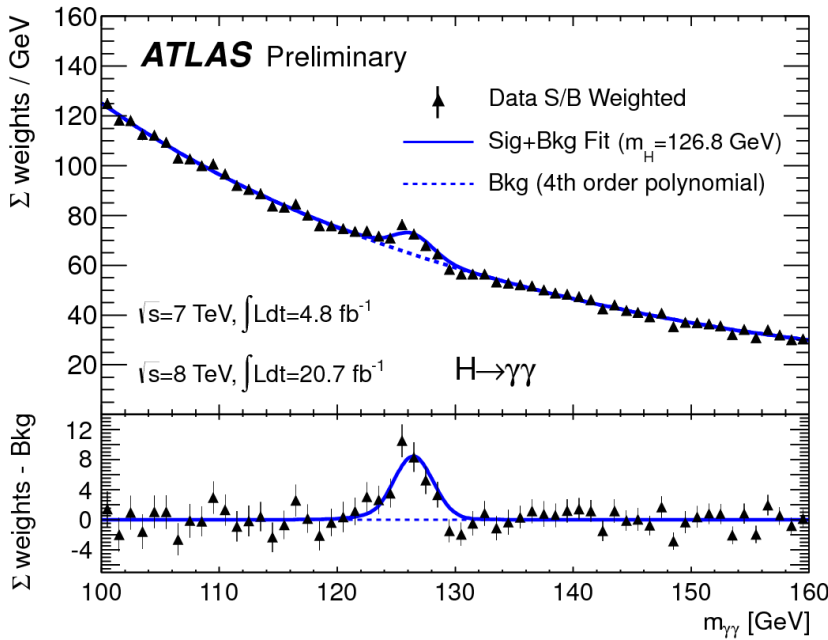


July 4, 2012

# Higgs to di-photons



**ATLAS**

**CMS**

# Higgs → 4 leptons



**ATLAS**

**CMS**

# Higgs → 4 leptons

**Signal**

$m_{Z_1}$

$m_{Z_2}$

$$pp \rightarrow H \rightarrow ZZ \rightarrow \ell^+ \ell^- \ell'^+ \ell'^-$$

**Background**

$$pp \rightarrow ZZ \rightarrow \ell^+ \ell^- \ell'^+ \ell'^-$$

$$x = (m_{Z1}, m_{Z2})$$

Run:      204769
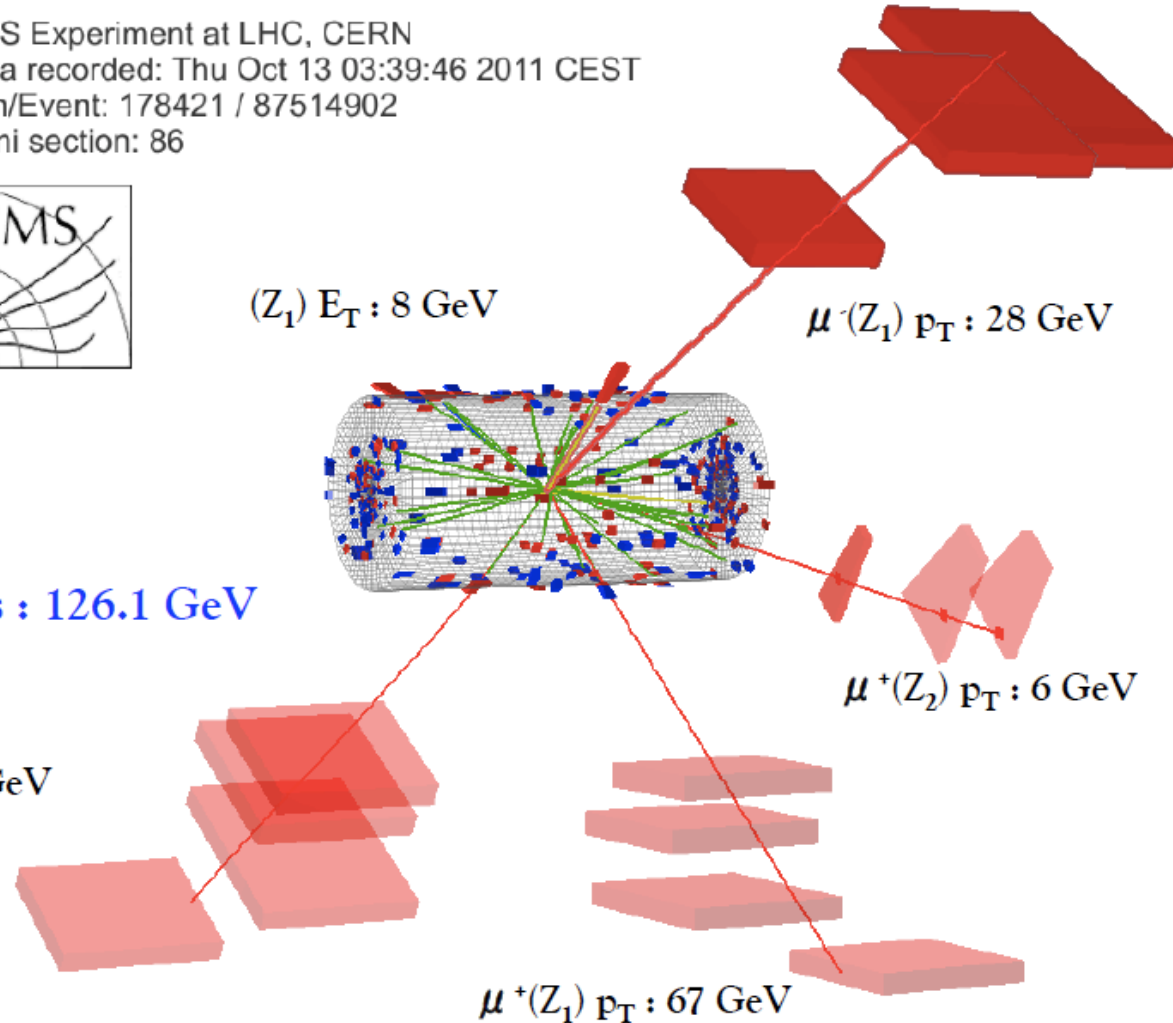Event:  71902630
Date: 2012-06-10
Time: 13:24:31 CEST

# 4-lepton event CMS

CMS Experiment at LHC, CERN
Data recorded: Thu Oct 13 03:39:46 2011 CEST
Run/Event: 178421 / 87514902
Lumi section: 86

$(Z_1)$ $E_T$ : 8 GeV
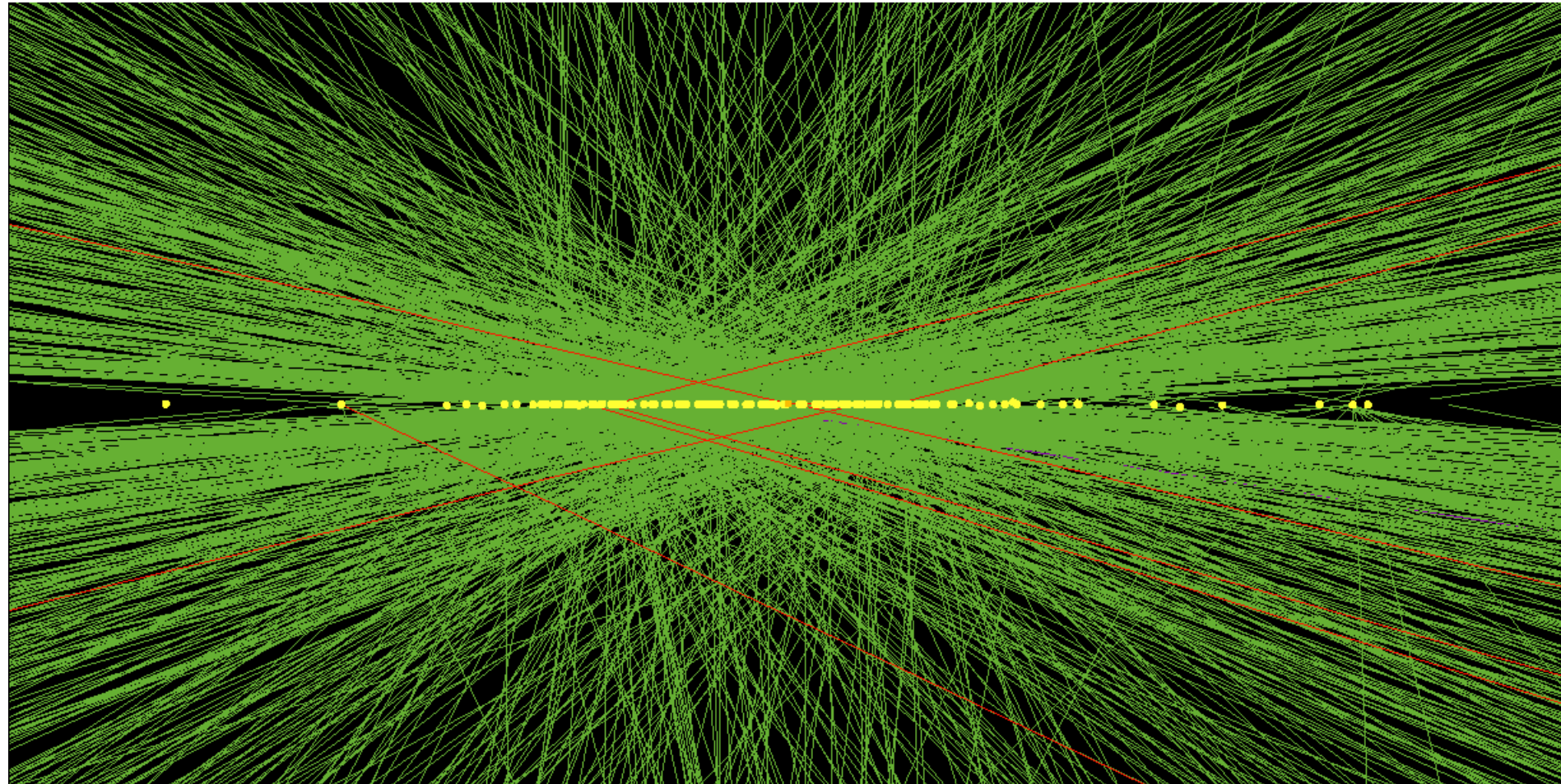
$\mu^-(Z_1)$ $p_T$ : 28 GeV

7 TeV DATA

$4\mu + \gamma$ Mass : 126.1 GeV

$\mu^+(Z_2)$ $p_T$ : 6 GeV

$\mu^-(Z_2)$ $p_T$ : 14 GeV

$\mu^+(Z_1)$ $p_T$ : 67 GeV

# Event Complexity

# Event Filtering

**$10^8$ sensors**

**Trigger Rate**

# **Applications**

## I. Classification

- **Particle Identification**
- **Pa**_____ **gnition (tracks)**
- **Se**_____ **New Physics**
- **Da**_____ **Monitoring**

# Applications

## II. Function estimation

- **Particle energy better estimated with ML methods**

- **ML Regression**



parametric

BDT

H→γγ MC
Illustration only

# Classification Theory

# **Classification Theory**

Signal density
$$p(x, \mathrm{s}) = p(x \mid \mathrm{s})\, p(\mathrm{s})$$

Background density
$$p(x, \mathrm{b}) = p(x \mid \mathrm{b})\, p(\mathrm{b})$$



density $p\,(x)$

$\beta$

$\alpha$

$x_0$

$x$

Optimality criterion: minimize the error rate, $\alpha + \beta$

# **Classification Theory**

The total loss *L* arising from classification errors is given by

$$L = L_b \int H(f)\, p(x, b)\, dx$$

Cost of background misclassification

$$+\, L_s \int [1 - H(f)]\, p(x, s)\, dx$$

Cost of signal misclassification

where $f(x) = 0$ defines a decision boundary
such that $f(x) > 0$ defines the acceptance region
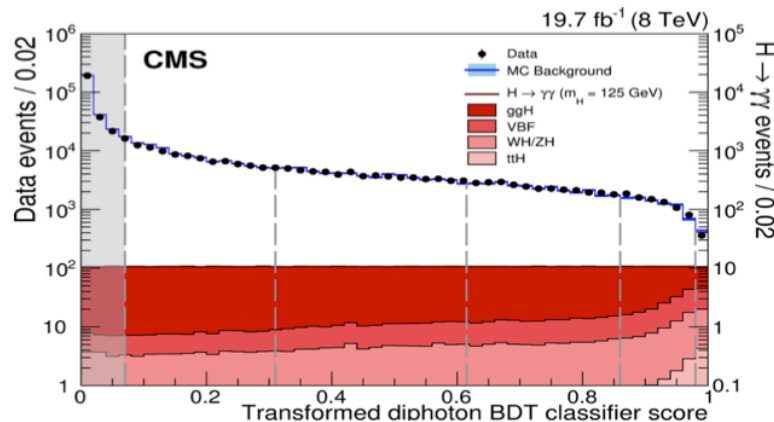
$H(f)$ is the Heaviside step function:
$$H(f) = 1 \text{ if } f > 0,\ 0 \text{ otherwise}$$

# Classification in Practice
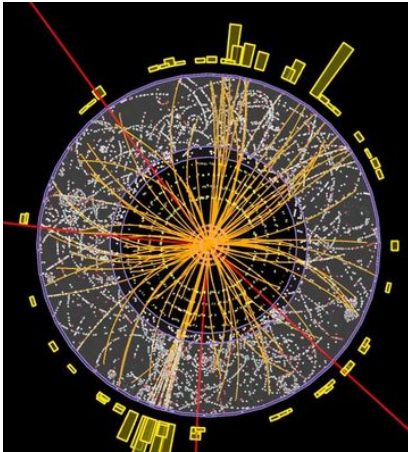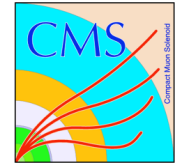
# in Higgs Discovery



- Identification of particles
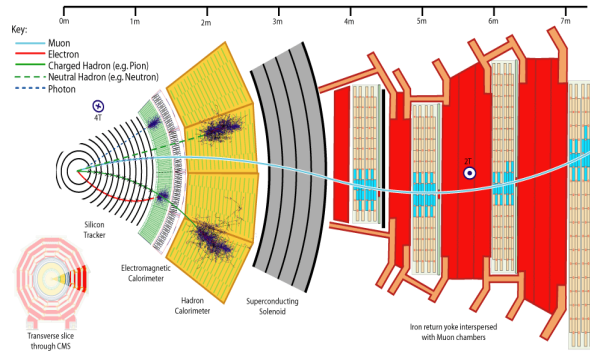- Identification of interactions
- Energy regression
- Event selection

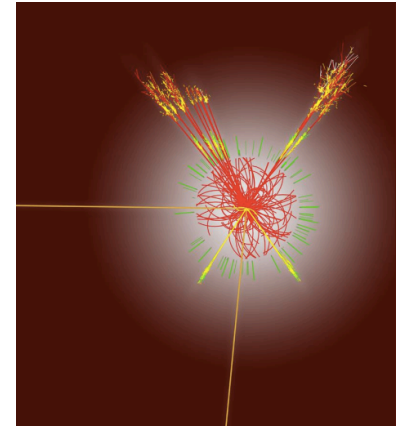**Improvement in analysis from all four areas**

# Interesting areas



**Tracking**

**Fast Simulation**

**Object Identification**

**Event Filtering**

**Imaging Techniques**

**Simulation**

Sergei V. Gleyzer
TAE 2017 Lecture

# CONSTRUCTING CLASSIFIERS

# **Classification**

**Distinguish f(x)**, **g(x)** using Training set of observations

{**inputs** , **outputs**}

Pass observations to a learning algorithm
    neural network, decision tree

that produces **outputs** in response to **inputs**

Use another set of observations to evaluate

# Classification

**Primary Goal:**

Achieve **lowest probability** of error

on unseen cases $\{<x^{(i)}, y^{(i)}>\}$

**Approach:**

Inductively learn from labeled examples

(where classes are known)

# ML Algorithms

- **Fisher, Quadratic**
- **Naïve Bayes (Likelihood)**
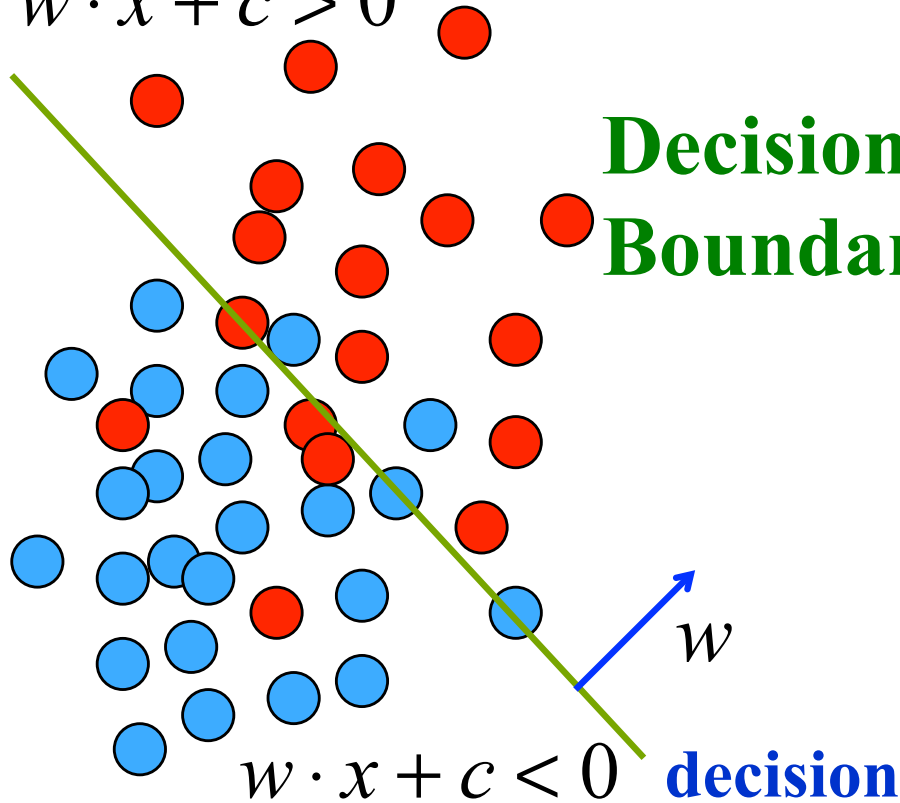- **Kernel Density Estimation**
- **Random Grid Search**
- **Rule ensembles**
- **Boosted decision trees**
- **Random forests**
- **Deep learning neural networks**
- **Support vector machines**
- **Genetic algorithms**

# Linear and Quadratic

## Linear (Fisher)

$$w \cdot x + c > 0$$

$$w \cdot x + c < 0 \quad \text{decision}$$

$w$

**Decision Boundaries**

## Quadratic

$$\lambda(x) = \ln \frac{G(x \mid \mu_s, \Sigma)}{G(x \mid \mu_b, \Sigma)} \rightarrow$$

$$w \propto \Sigma^{-1}(\mu_s - \mu_b)$$

**decision boundary**

# Binary Decision Trees

# Decision Trees

- **Decision trees** are **multidimensional histograms**
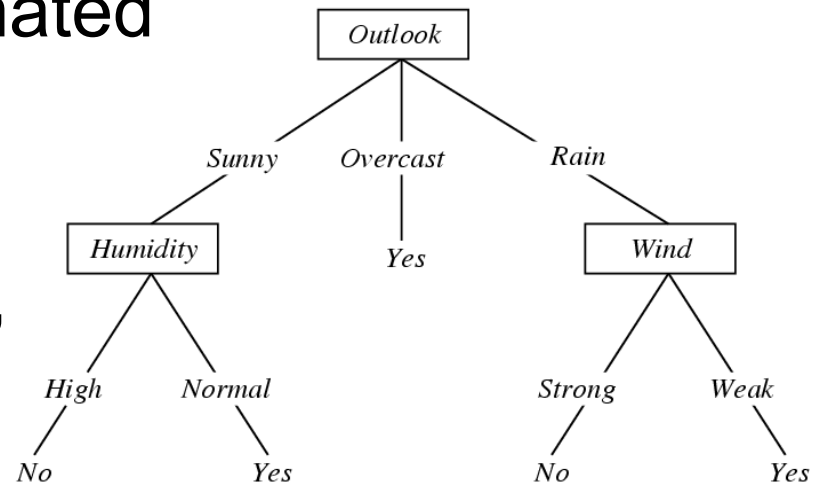  - Recursively constructed bins
  - Each associated to the value (or **class**) of f(x) to be approximated
  - **Golf-Playing** Decision Tree: f(outlook, humidity, wind, T)
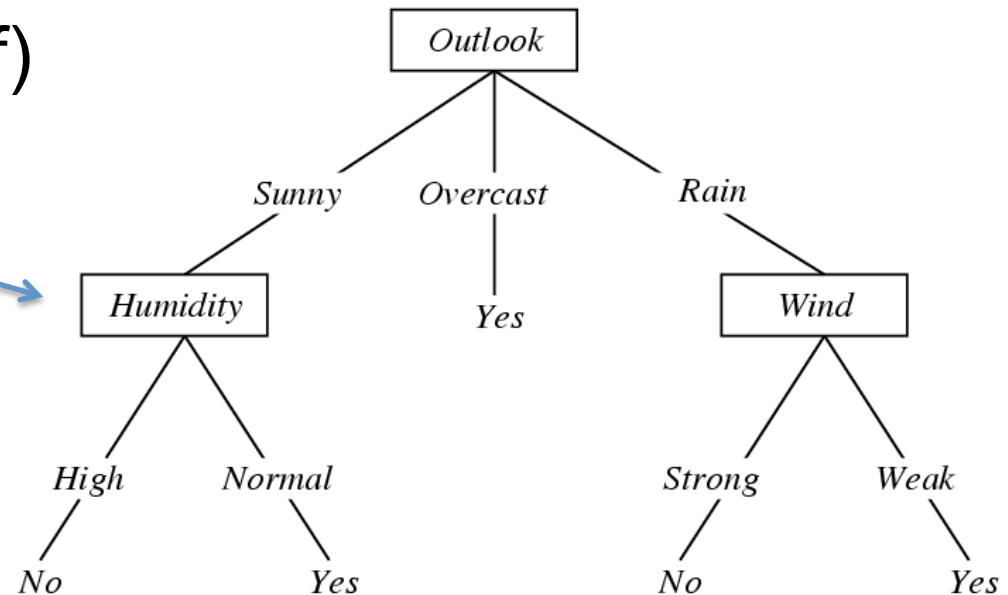
# **Decision Trees**

- Each **internal** node: test one attribute $X_i$
- Each **branch**: selects one value for $X_i$
- Each **leaf** node: predict Y
  - Or P(Y|X in leaf)

**Decision Node** →
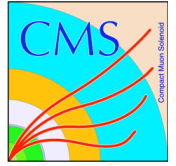
# Decision Tree Learning

- Unknown **target function** f: **X→Y**
  - **Y** is discrete valued (class)
- Set of possible instances **X**
  - each **instance** is a feature vector

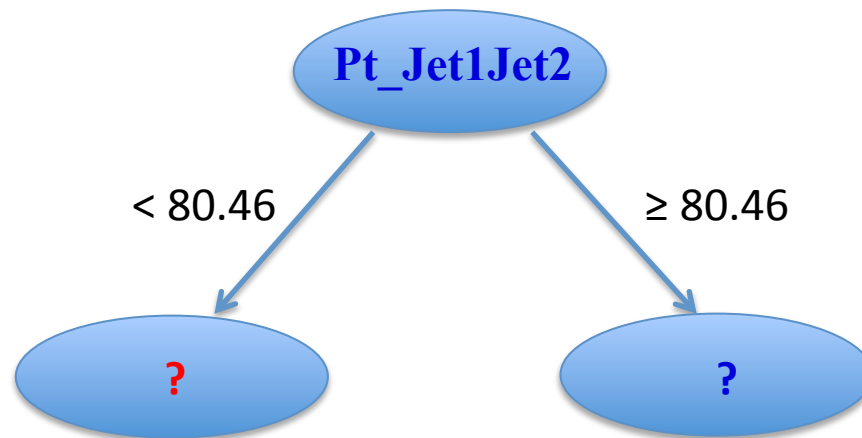  e.g. <Humidity = High, Wind = weak, Outlook = rain, Temp = hot>

# **Decision Tree Learning**

## Input:

– Training examples $\{<x^i, y^i>\}$

## Output

– Hypothesis $h \in H$ that

best approximates target function f

– Tree sorts x to leaf, which assigns y

# **Decision Trees**

## **Building a tree:**

- Scan along each variable and propose a **DECISION**

  – A cut on value that maximizes class separation (binary branching)



Pt_Jet1Jet2

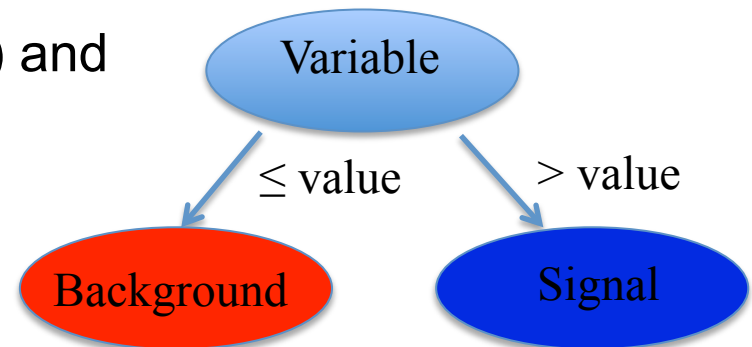< 80.46            ≥ 80.46

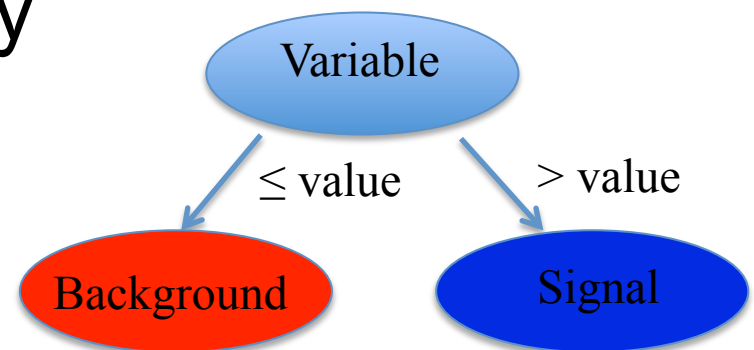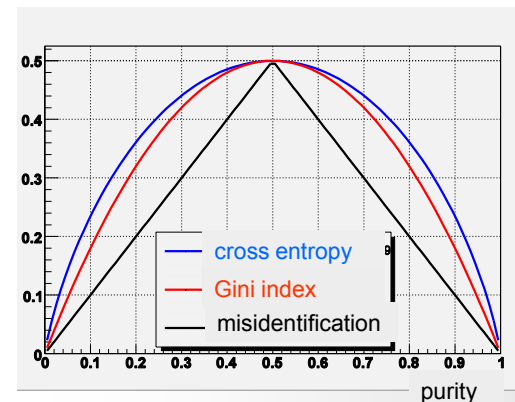?            ?

# Decision Trees

- Choose **decision** that leads to greatest separation among classes **signal**/ **background**

  - Based on the information gained from split

    - Build regions of increasing purity
    - Stop when no further improvement from additional branching
    - Reach terminal node (leaf) and assign purity-based class

$$\frac{N_{signal}}{N_{signal} + N_{background}}$$

# Separation Gain

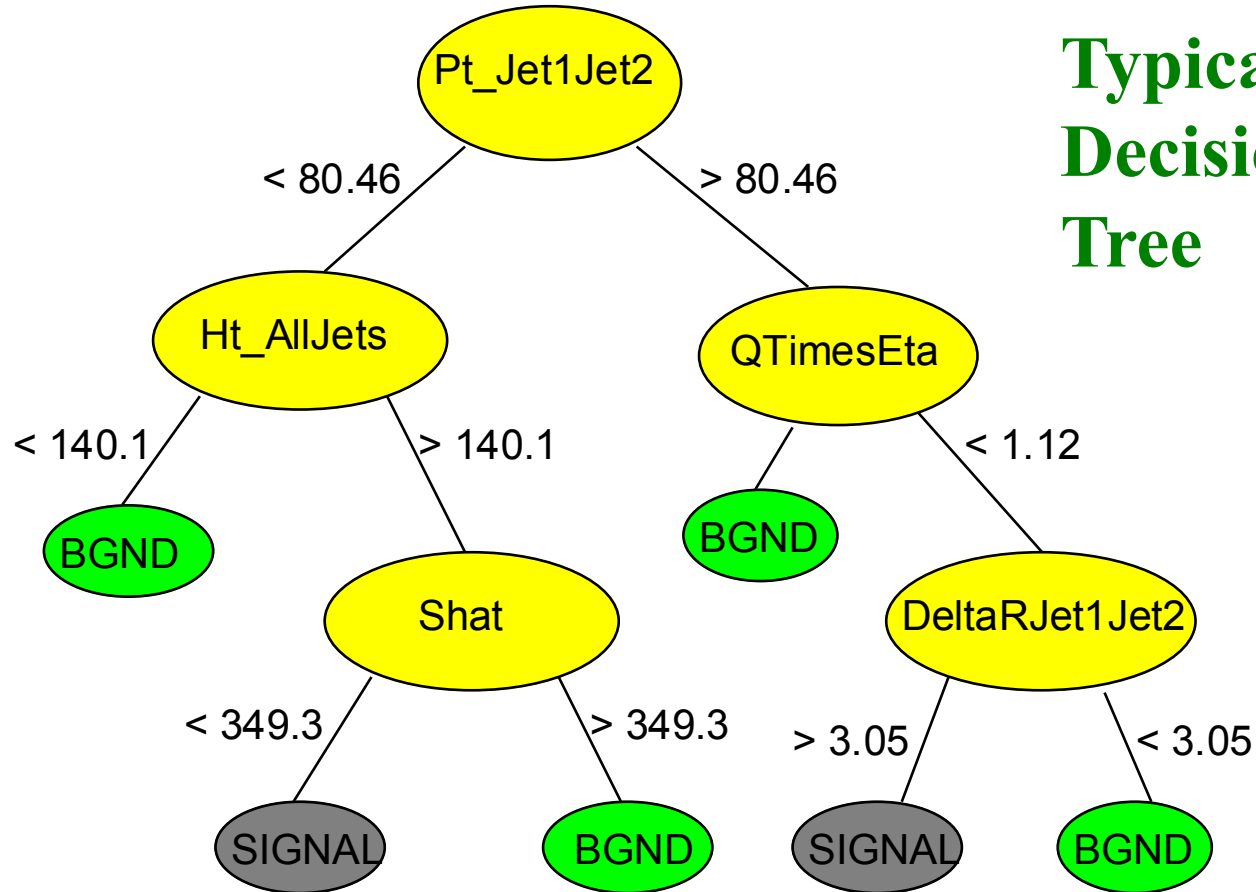## Measures of Separation Gain

- Cross-Entropy
  - $-(p\ln p + (1-p)\ln(1-p))$
- Gini Index
  - $p(1-p)$
- Want to lower entropy due to split

# **Representation**



**Typical Decision Tree**

# Pruning

Decision trees can become large and complex and risk over-fitting the data

**Pruning:** remove parts of the tree that are less powerful or possibly noisy

– start from the leaves and work back up

Pruned trees smaller in size, easier to interpret

# Summary

- **Machine Learning is a very powerful field with an expanding number of applications**
  - Basic Methods: Linear, Quadratic, Decision Trees, Decision Rules
  - More advanced methods next time
  - Many methods available, good to experiment