**Machine** **Learning**

**Sergei** **Gleyzer**

**Lecture** **I**

**UPRM Lectures**

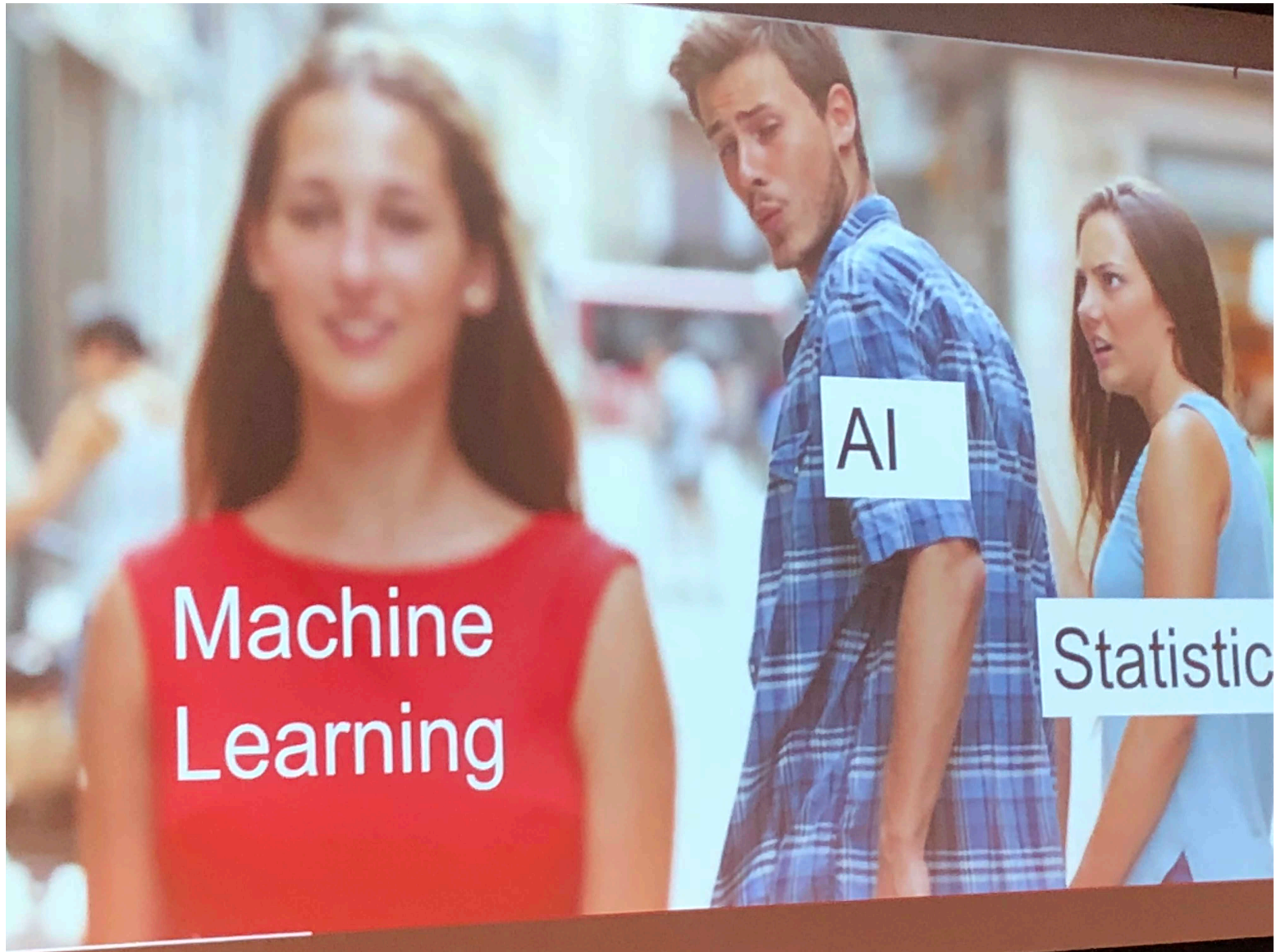**April 24, 2019**

# Today's Outline

- **What is Machine Learning**
- **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**

# 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
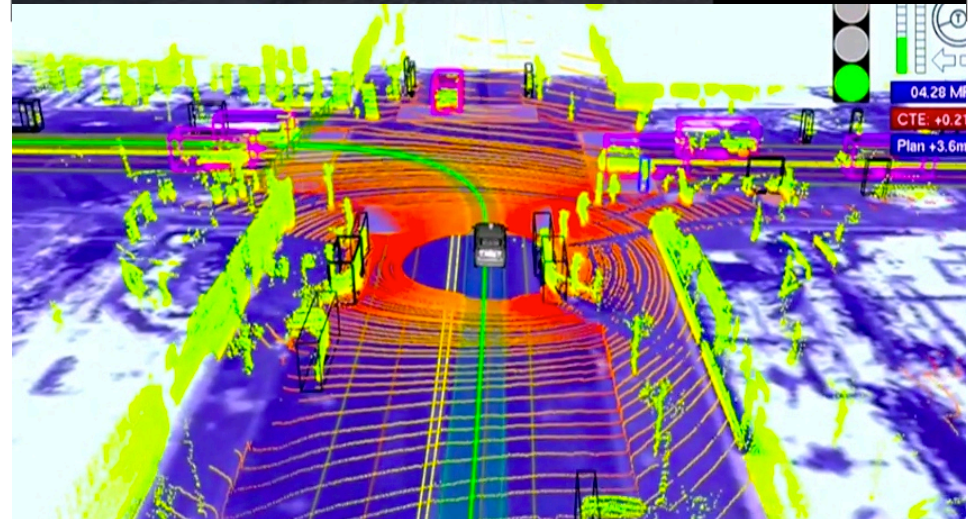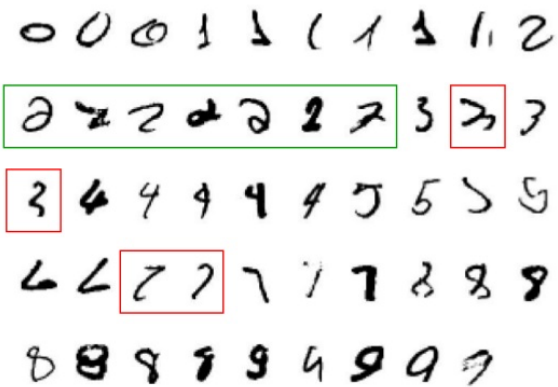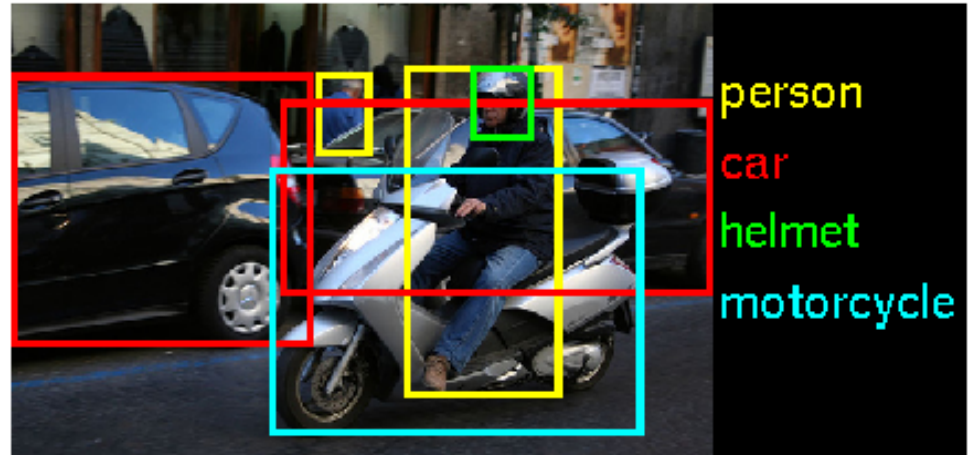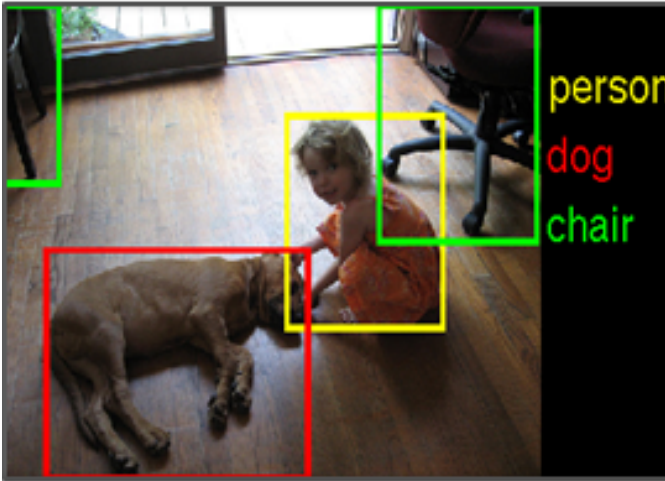
# Machine Learning

**General Approach:**

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

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

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

# Examples

# Machine Learning

**Choose**

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

Constraint        $C$

Loss function*        $L$

$f(x, w^*)$        $C(w)$        $F$

**Method**

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

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

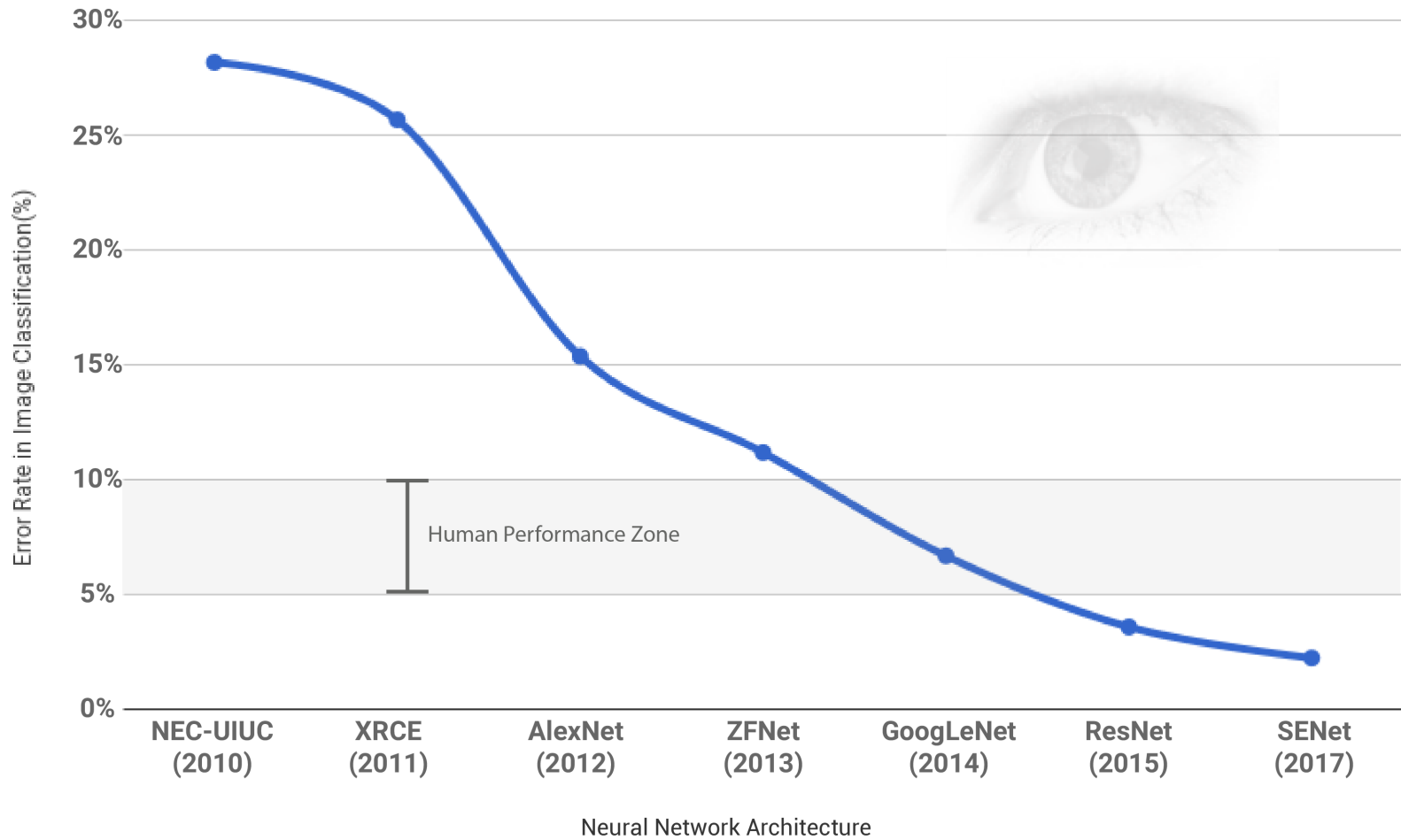1950s:        First methods invented

1960-80s: Focus on knowledge

1990s:        Computing power, new learning
                    methods, data-centrism

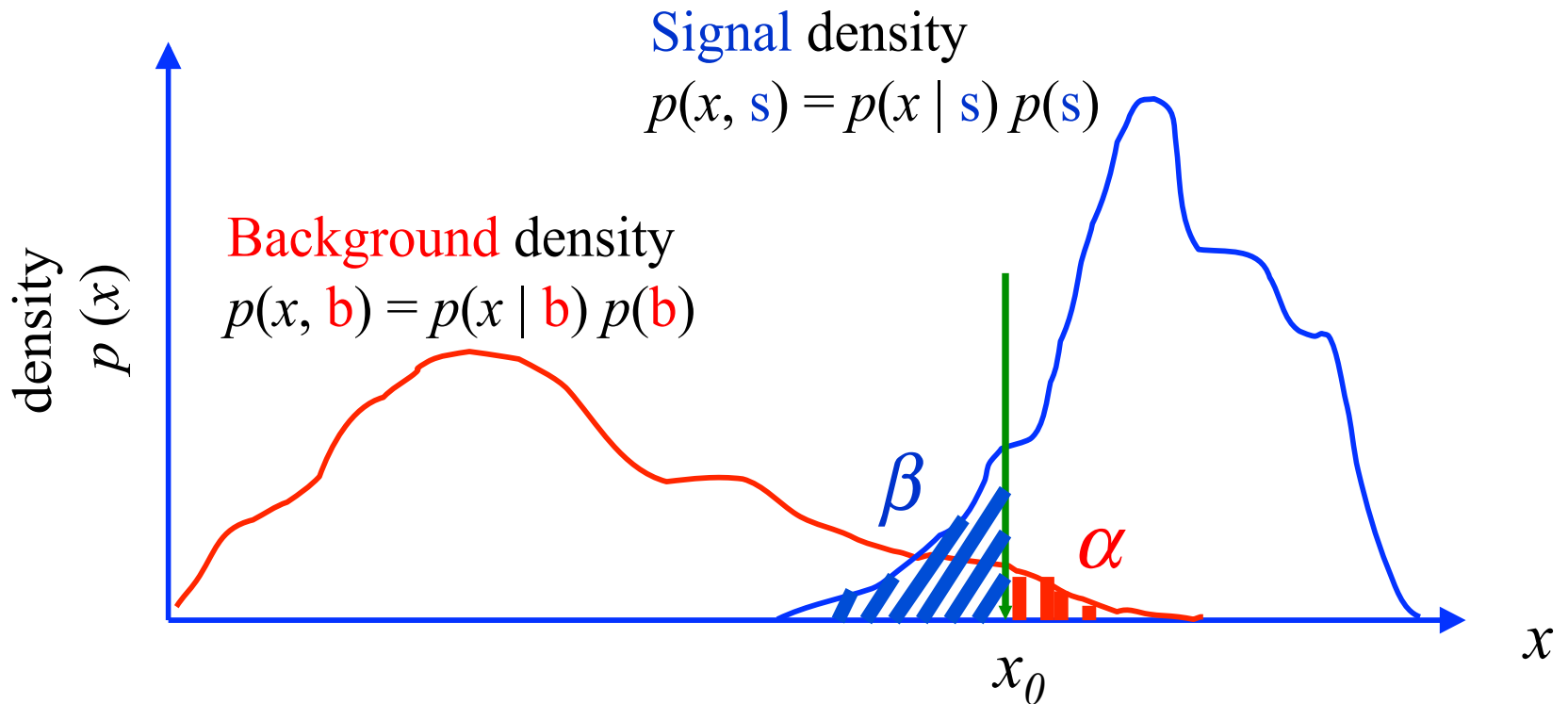2000-10s: Wider use research and industry

2010s:        Learning improvement, dedicated
                    hardware, deep learning

# Diving Deeper

# Classification Theory

# Classification Theory



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

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

$\beta$

$\alpha$

$x_0$

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

density $p(x)$

# Classification Theory
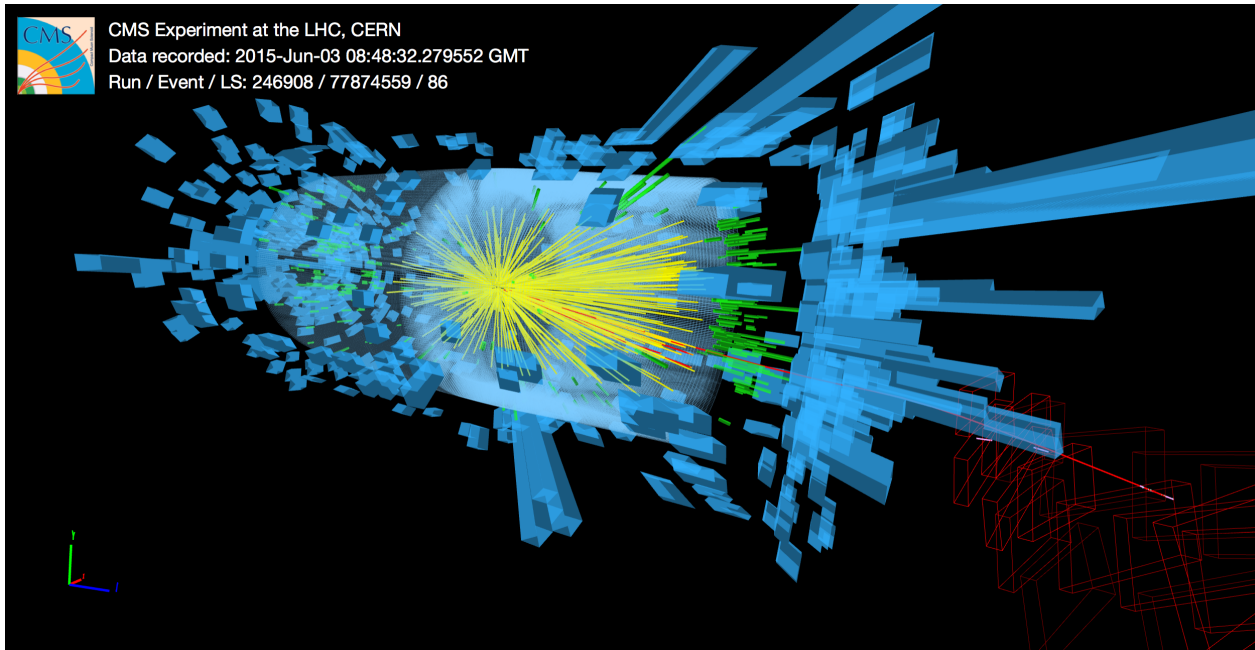
The total loss $L$ arising from classification errors is given by

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

<span style="color:red">Cost of background misclassification</span>

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

<span style="color:blue">Cost of signal misclassification</span>

where $f(x) = 0$ defines a decision boundary
such that $f(x) > 0$ defines the acceptance region

$H(f)$ is the Heaviside step function:
$$H(f) = 1 \text{ if } f > 0, \, 0 \text{ otherwise}$$
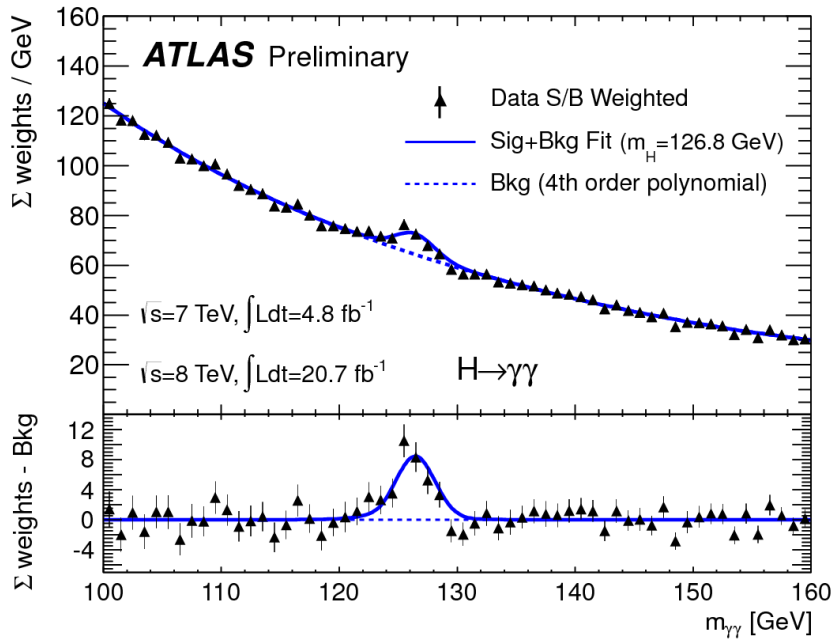
# **Classification in Practice**
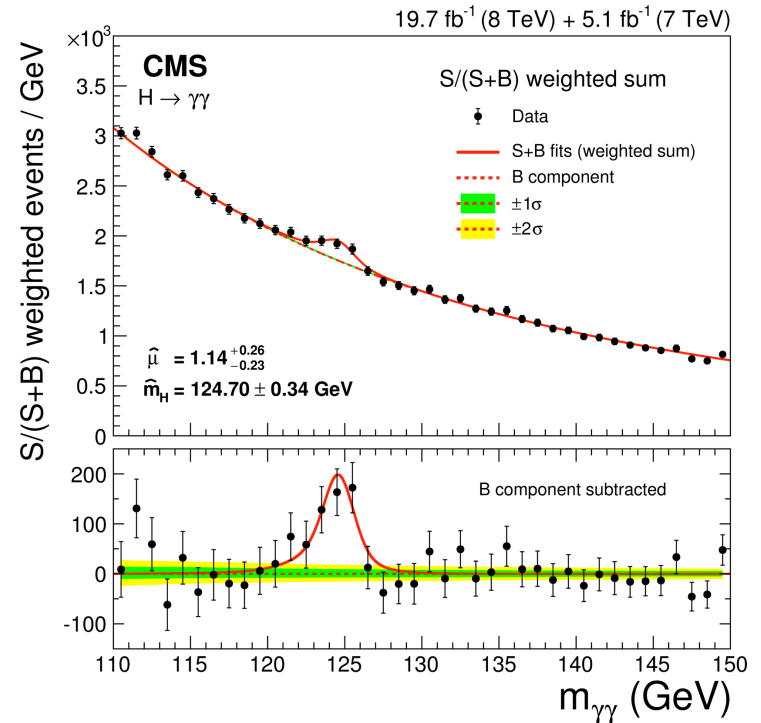
# In Particle Physics

# Higgs Boson Discovery
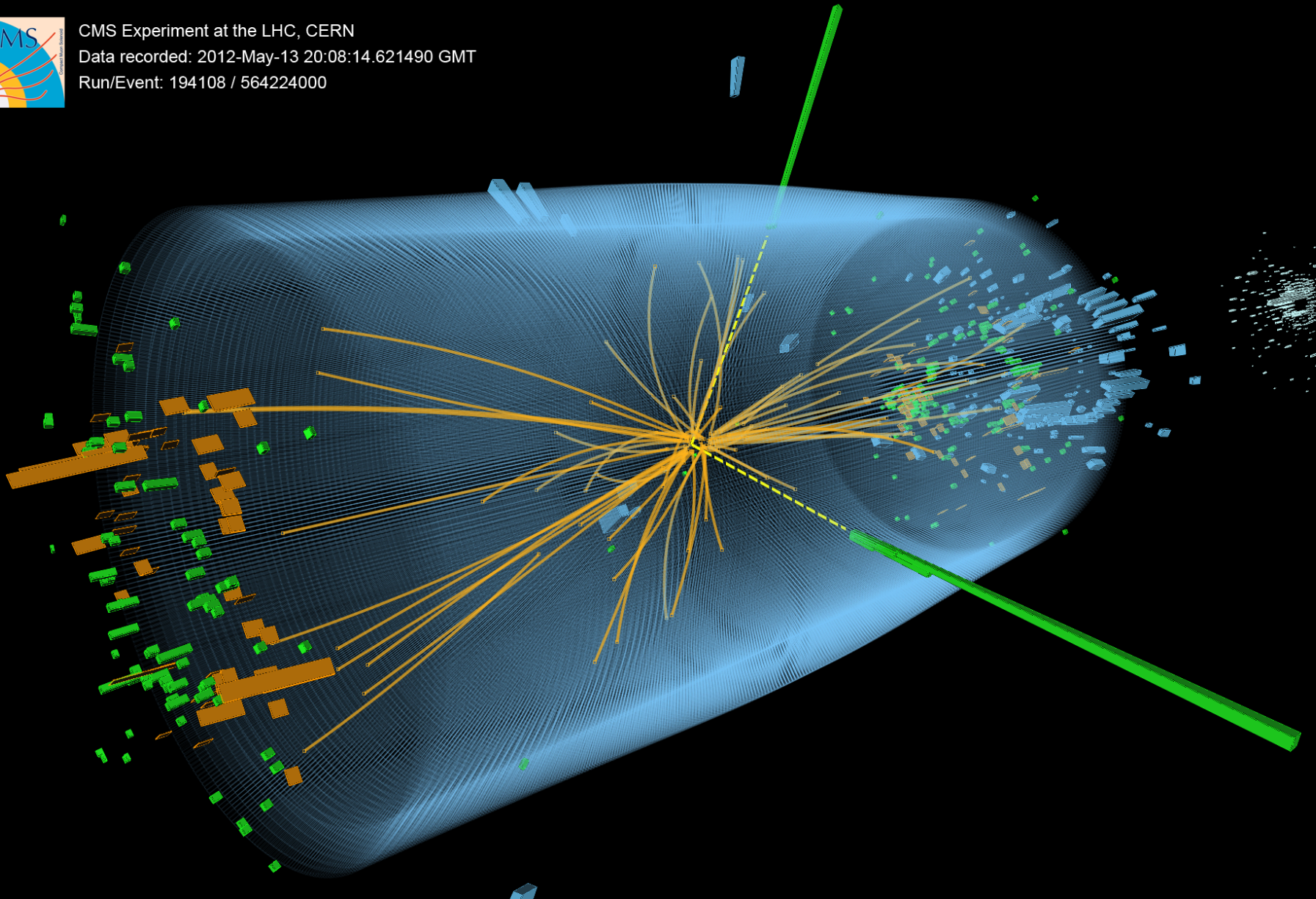


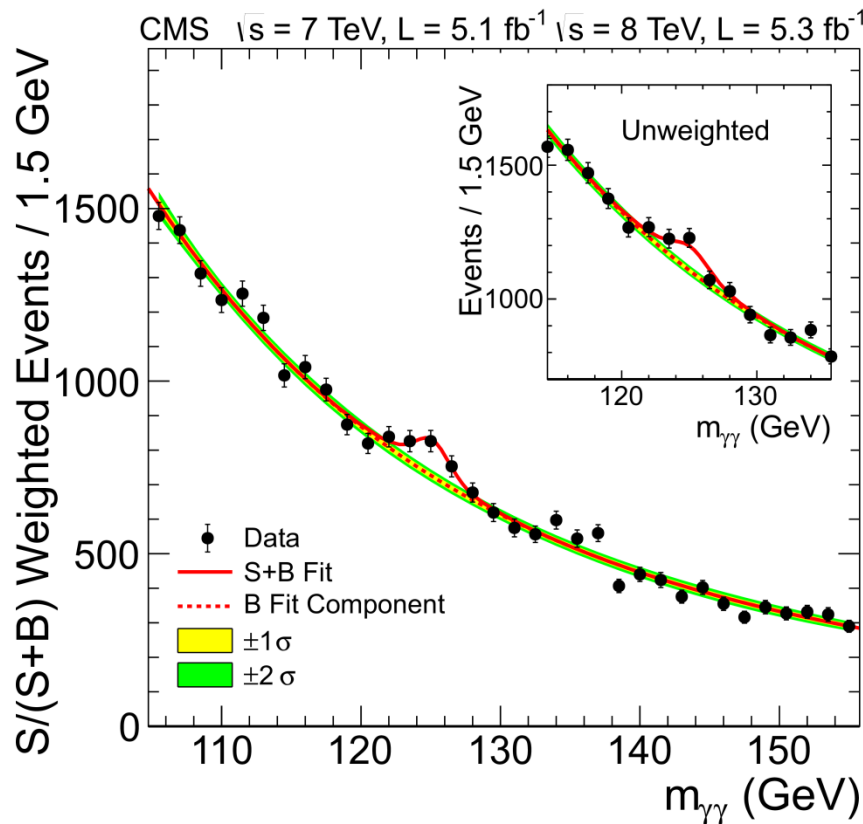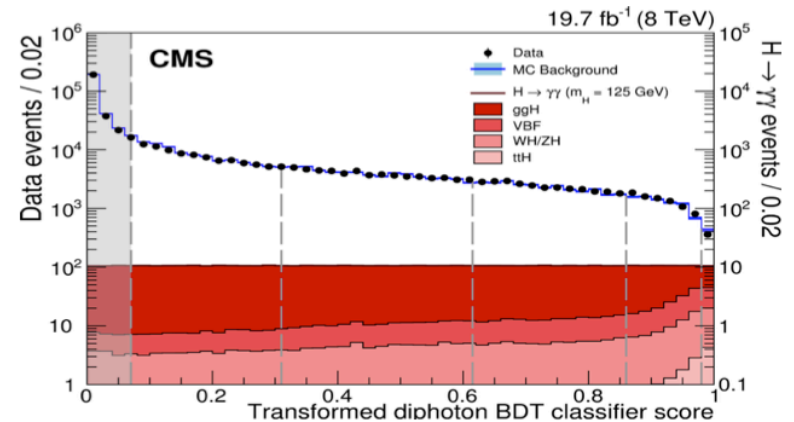July 4, 2012

# Higgs to di-photons



**ATLAS**

**CMS**

CMS Experiment at the LHC, CERN
Data recorded: 2012-May-13 20:08:14.621490 GMT
Run/Event: 194108 / 564224000

# in Higgs Discovery



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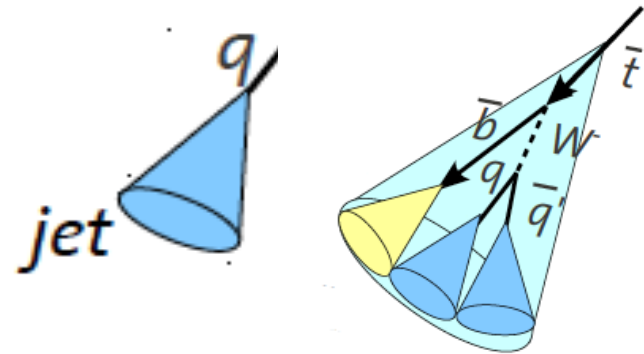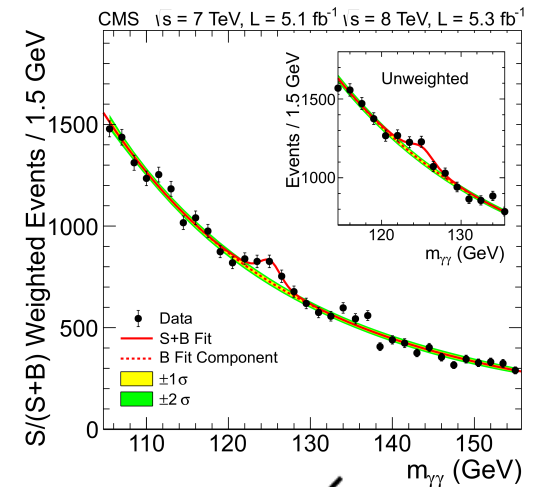- Event selection

**Improvement in analysis from all four areas**

# ML in HEP Today

**Machine learning already at forefront of what we do:**

- Physics object **identification**
- Event type **classification**
- Object properties **regression**

**Exp**  **quickly**

# 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



**Inputs**

Pt_Jet1Jet2

< 80.46         > 80.46

Ht_AllJets              QTimesEta

< 140.1   > 140.1              < 1.12

BGND        Shat          BGND       DeltaRJet1Jet2

< 349.3   > 349.3   > 3.05      < 3.05

SIGNAL   BGND   SIGNAL   BGND

**Outputs**

# Learning Types

**Supervised Learning**

Labeled examples with known classes

Examples: cats/dogs

**Unsupervised Learning**

Un-labeled data

Examples: clustering, anomaly detection

# **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
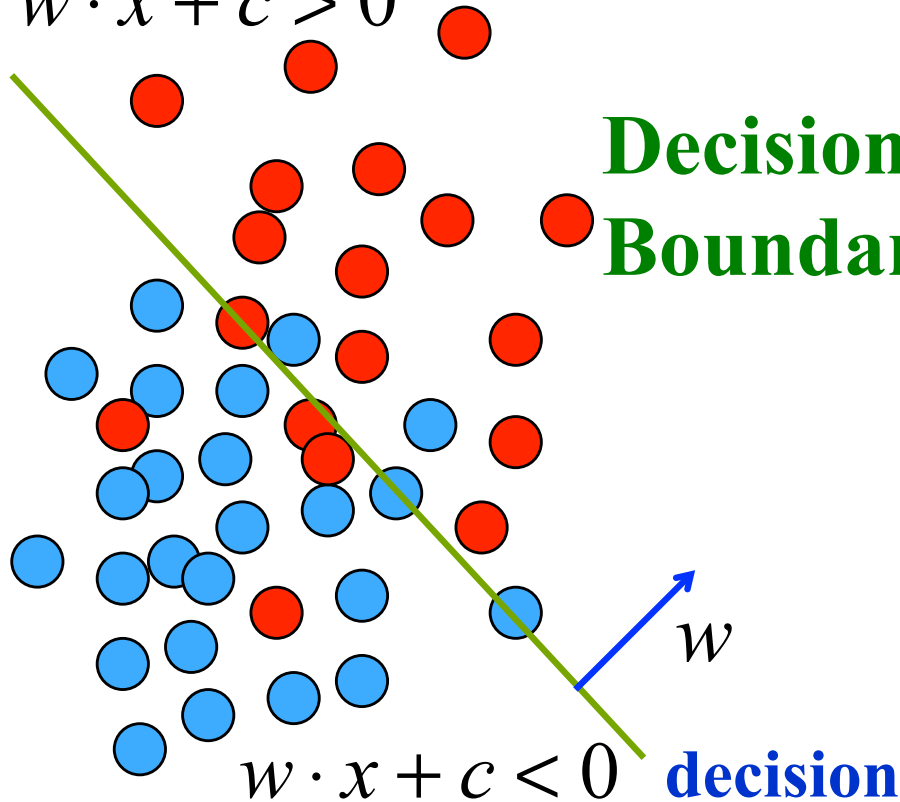
# ML Algorithms

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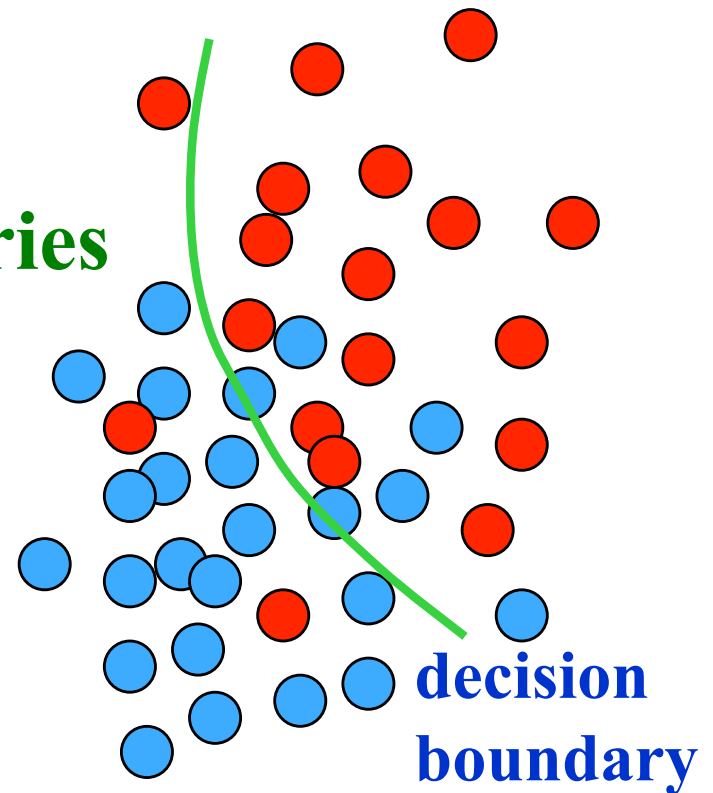- **Genetic algorithms**
- **Deep learning neural networks**

# Linear and Quadratic

**Linear (Fisher)**

$w \cdot x + c > 0$

**Quadratic**

**Decision Boundaries**

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

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

$w$

$w \cdot x + c < 0$ **decision**
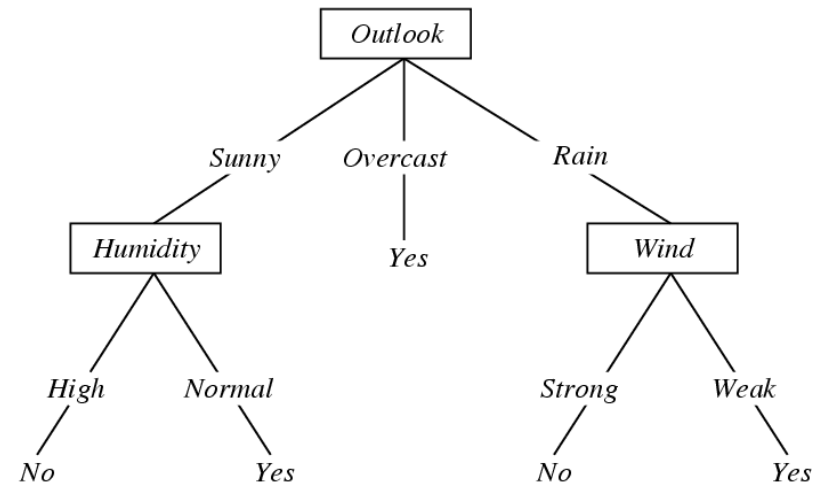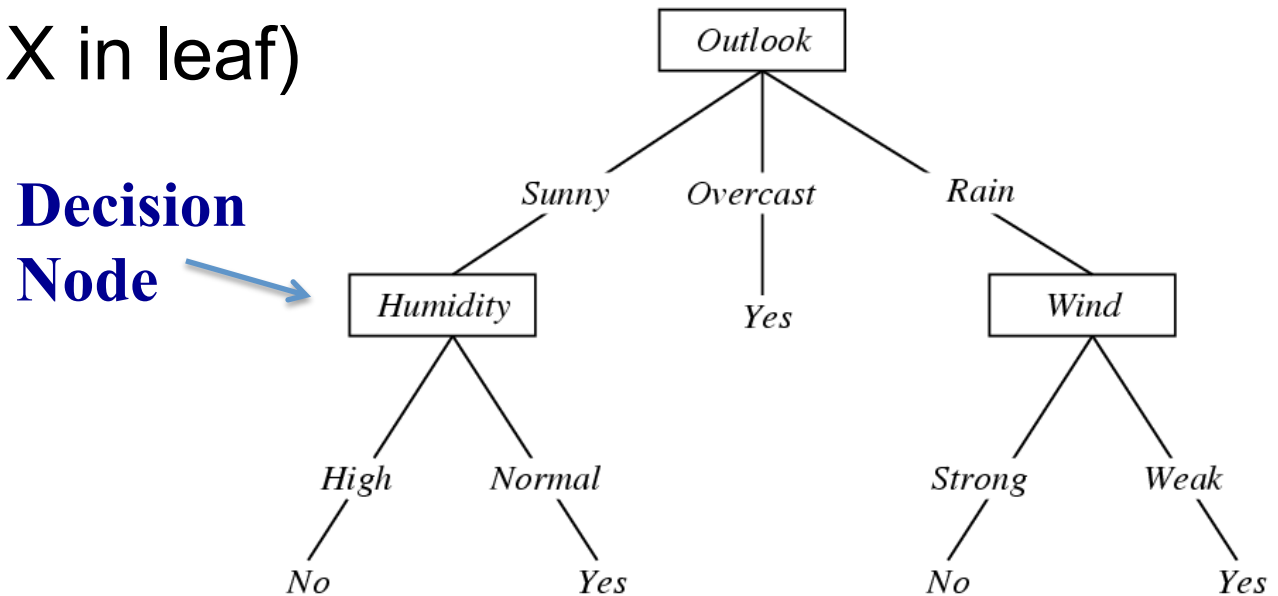
**decision boundary**

# **Binary Decision Trees**

# Decision Trees

- **Decision trees** are recursively constructed **multidimensional histograms**
  - Each leaf associated to the value (**class**) of f(x) to be approximated

  - Golf-Playing Tree: f(outlook, humidity, wind, temperature)

# Decision Trees

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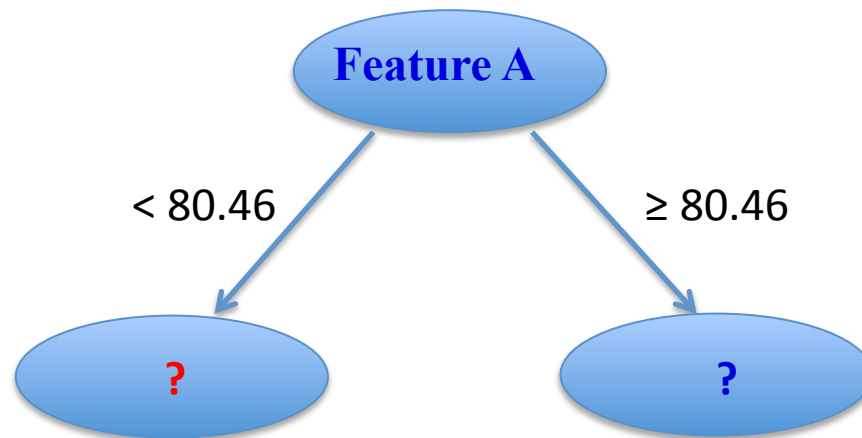  - Or P(Y|X in leaf)

**Decision Node**

# Decision Tree Learning

- Set of possible instances: **X**
  - **instance** is a feature vector

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

- Unknown target function f: **X** → **Y**
  - **Y** is discrete valued (class)

# Decision Trees

## Building a tree:

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

  - A cut on value that maximizes class separation (binary branching)
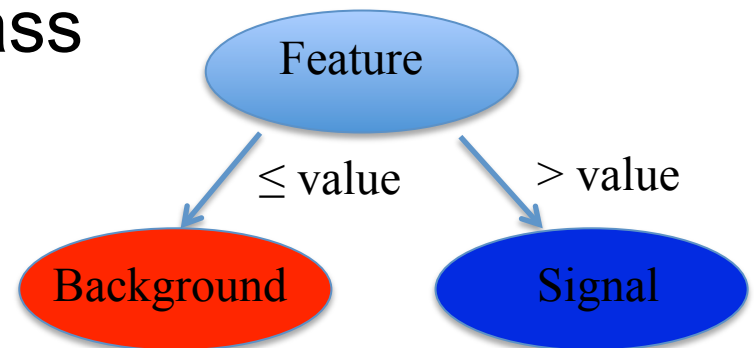
# Decision Trees

Choose **split** that leads to greatest separation among <span style="color:red">classes</span>

- Based on the information gained
  - Build regions of increasing purity
  - Stop when no improvement from branching
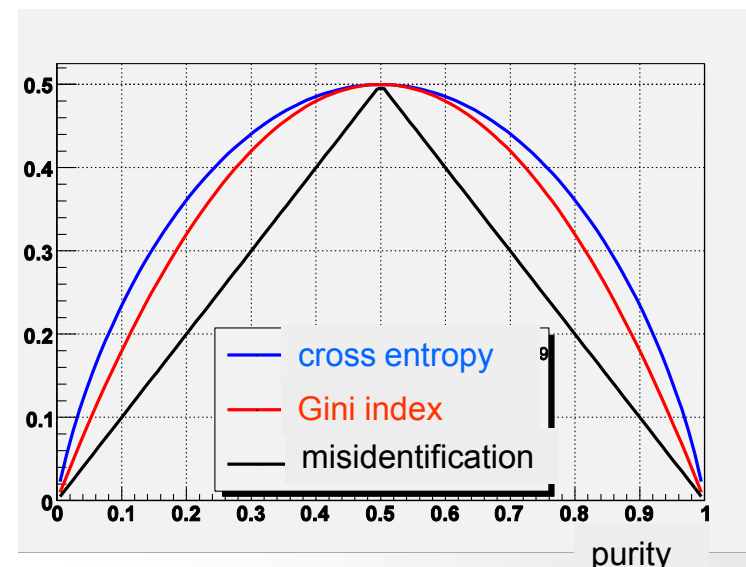  - Reach terminal node (leaf) and assign purity-based class

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



Feature

≤ value        > value

Background        Signal

# Separation Gain

## Popular Separation Gain Measures

- Cross-Entropy

  -p ln p + (1-p) ln(1-p)

- Gini Index

  p ( 1 – p )



Want to lower entropy from split

# Pruning

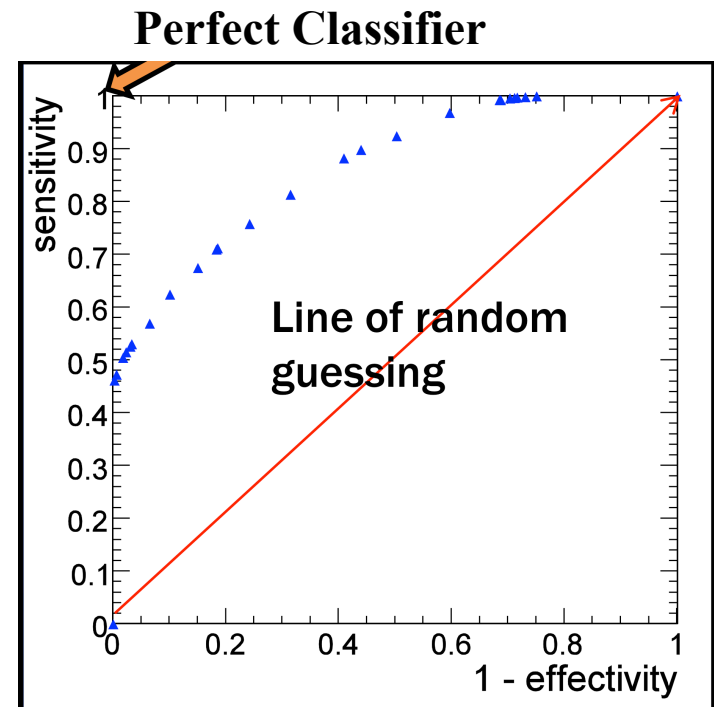Decision trees can become large and complex and risk over-fitting the data

- **Pruning:** remove parts of tree that are less powerful or noisy
  - start from the leaves and work back up
  - pruned trees smaller in size, easier to interpret

# Classifier Performance

## Receiver Operating Characteristic (ROC)

**Commonly used metric**

Shows the *relationship* between correctly classified positive cases (sensitivity) and incorrectly classified negative cases (1-effectivity)

**Perfect Classifier**

# Summary

- **Machine learning focuses on algorithms capable of learning**
  - Basic methods: linear, quadratic
  - Tree-based methods:
    - Decision Trees, Rules
  - Ensembles:
    - Boosted Decision Trees, Random Forests
  - Next lecture: Deep learning algorithms