Thomas Fuchs (JPL, Caltech)

# Random Forests

# Random Forests: Citations

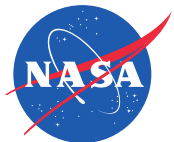**Total citations**    Cited by 13154

**Citations per year**



**Scholar articles**    Random forests
L Breiman - Machine learning, 2001
Cited by 13154 - Related articles - All 83 versions

Google Scholar Citations: 2014-08-21

# History

1983  **CART**                              *Breiman*



**Leo Breiman**
1928 - 2005

1996  **Bagging**                           *Breiman*

1996  **AdaBoost**                          *Freund & Schapire*

2001  **Random Forests**        *Breiman*

# History

1983  **CART**                              *Breiman*

1996  **Bagging**                       *Breiman*

**Leo Breiman**
1928 - 2005

1994  **Randomized Trees (WS)** *Amint & Geman*
1996  **AdaBoost**                *Freund & Schapire*
1997  **Randomized Trees**      *Amint & Geman*
1998  **Decision Forests**       *Ho*
1998  **Random split selection** *Dietterich*
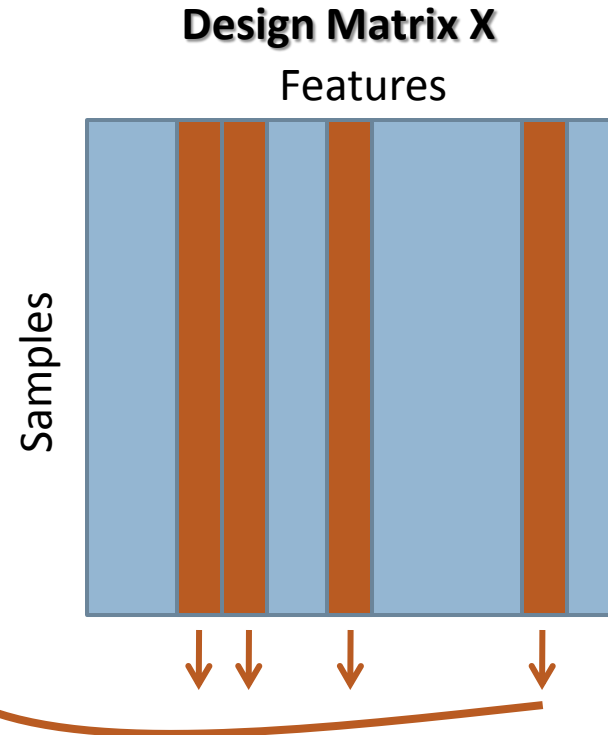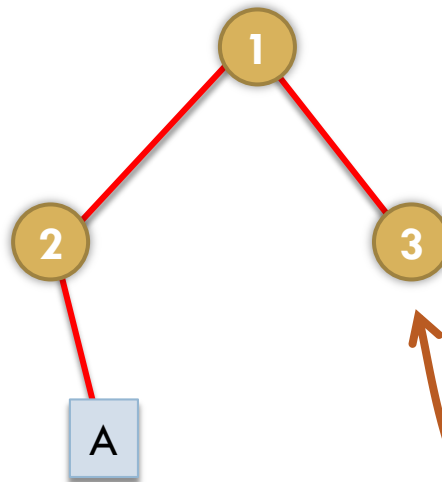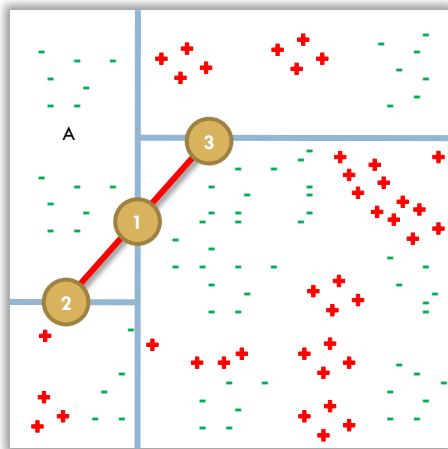2001  **Random Forests**       *Breiman*

# Random Forests    (Breiman 2001)

**Definition 1.1**   *A random forest is a classifier consisting of a collection of tree-structured classifiers   $\{h(\mathbf{x},\Theta_k), k=1, ...\}$   where the   $\{\Theta_k\}$   are independent identically distributed random vectors and each tree casts a unit vote for the most popular class at input*  $\mathbf{x}$ .

The common element in all of these procedures is that for the $k$th tree, a random vector $\Theta_k$ is generated, independent of the past random vectors $\Theta_{1,...}, \Theta_{k-1}$ but with the same distribution; and a tree is grown using the training set and $\Theta_k$, resulting in a classifier $h(\mathbf{x}, \Theta_k)$ where $\mathbf{x}$ is an input vector.

# Randomized Tree Learning



**Design Matrix X**

Features

Samples

At each node only a random subset of features is considered to choose the best split. Common splitting criteria are Entropy, Gini Index and misclassification rate.

# Random Forest Learning

$$Z = \{ \; A \quad B \quad C \quad D \quad E \quad F \; \}$$

Bootstrap Samples

$$Z^{*1} = \left\{ \begin{array}{ccc} B & E & A \\ B & C & A \end{array} \right\}$$

$$Z^{*2} = \left\{ \begin{array}{ccc} E & F & E \\ F & C & D \end{array} \right\}$$

$$Z^{*T} = \left\{ \begin{array}{ccc} E & F & A \\ F & C & D \end{array} \right\}$$

# Random Forest Learning

$Z = \{$ A B C D E F $\}$

Bootstrap Samples

$Z^{*1} = \{$ B E A / B C A $\}$ → Tree 1

$Z^{*2} = \{$ E F E / F C D $\}$ → Tree 2

$Z^{*T} = \{$ E F A / F C D $\}$ → Tree T

Random Forest Model

# Random Forest Classification

# Out Of Bag (OOB) Error

$Z = \{$ A B C D E F $\}$

Bootstrap Samples     OOB Samples   classify

$Z^{*1} = \{$ B E A / B C A $\}$ → $\{$ D / F $\}$ → Tree 1

$Z^{*2} = \{$ E F E / F C D $\}$ → $\{$ A / B $\}$ → Tree 1-2

$Z^{*T} = \{$ E F A / F C D $\}$ → $\{$ B $\}$ → Tree 1-T

# Out Of Bag Error



Convergence after 100 Trees

(a)

(b)

(c)

| | Example 1 - two-class spiral | Example 2 - four-class spiral | Example 3 - noisier four-class spiral |
|---|---|---|---|

Training points

Testing posteriors

Entropy images

Plot from [Criminisi et al. 2012]