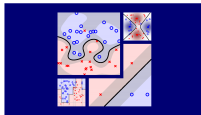


# Machine Learning Techniques (機器學習技法)



## Lecture 10: Random Forest

Hsuan-Tien Lin (林軒田)

[htlin@csie.ntu.edu.tw](mailto:htlin@csie.ntu.edu.tw)

Department of Computer Science  
& Information Engineering

National Taiwan University  
(國立台灣大學資訊工程系)



# Roadmap

- 1 Embedding Numerous Features: Kernel Models
- 2 Combining Predictive Features: Aggregation Models

## Lecture 9: Decision Tree

**recursive branching (purification)** for **conditional aggregation** of **constant hypotheses**

## Lecture 10: Random Forest

- Random Forest Algorithm
- Out-Of-Bag Estimate      bagging+ decision tree
- Feature Selection
- Random Forest in Action

- 3 Distilling Implicit Features: Extraction Models

## Recall: Bagging and Decision Tree

## Bagging

function **Bag**( $\mathcal{D}, \mathcal{A}$ )

For  $t = 1, 2, \dots, T$

- ① request size- $N'$  data  $\tilde{\mathcal{D}}_t$  by **bootstrapping** with  $\mathcal{D}$
- ② obtain base  $g_t$  by  $\mathcal{A}(\tilde{\mathcal{D}}_t)$

return  $G = \text{Uniform}(\{g_t\})$

bootstrap sample 得到不同  $g_t$ , uniform voting  
取 uniform 以后, 不同  $g_t$  的差异被消除, 输出  
结果稳定, 受输入数据的改变影响小。不同  
输入的结果  $G(x)$  方差较小

—**reduces variance**

by voting/averaging

## Decision Tree

function **DTree**( $\mathcal{D}$ )

if **termination** return base  $g_t$   
else

- ① learn  $b(x)$  and split  $\mathcal{D}$  to  $\mathcal{D}_c$  by  $b(x)$
- ② build  $G_c \leftarrow \text{DTree}(\mathcal{D}_c)$
- ③ return  $G(x) =$

$$\sum_{c=1}^C \mathbb{I}[b(x) = c] G_c(x)$$

—**large variance** 切分点不同, 结果很不  
同。输入敏感, 方差大  
especially if fully-grown

putting them together?

(i.e. **aggregation of aggregation :-)**)

## 基本的随机森林

Random Forest (RF) 不同 $g_t$ 只由不同 $D_t$ 产生

**random forest (RF) = bagging + fully-grown C&RT decision tree**

function **RandomForest**( $\mathcal{D}$ )

For  $t = 1, 2, \dots, T$

- ① request size- $N'$  data  $\tilde{\mathcal{D}}_t$  by **bootstrapping** with  $\mathcal{D}$
- ② obtain tree  $g_t$  by **DTree**( $\tilde{\mathcal{D}}_t$ )

return  $G = \text{Uniform}(\{g_t\})$

决策树 $g_t$ 的输入是bootstrap sample  
 $g_t$ 对输入特别敏感, 所以能产生多样性

而利用bagging, 可以有效降低方差, 使结果稳定

function **DTree**( $\mathcal{D}$ )

if **termination** return base  $g_t$

else

- ① learn  $b(\mathbf{x})$  and split  $\mathcal{D}$  to  $\mathcal{D}_c$  by  $b(\mathbf{x})$
- ② build  $G_c \leftarrow \text{DTree}(\mathcal{D}_c)$
- ③ return  $G(\mathbf{x}) =$

$$\sum_{c=1}^C \mathbb{I}[b(\mathbf{x}) = c] G_c(\mathbf{x})$$

决策树作为 $g_t(\mathbf{x})$

- highly **parallel/efficient** to learn
- **inherit pros** of C&RT 可以多个树在多个不同抽样上并行计算
- **eliminate cons** of fully-grown tree

# Diversifying by Feature Projection

recall: **data randomness** for **diversity** in **bagging**

之前gt多样性只靠不同抽样

randomly **sample  $N'$  examples** from  $\mathcal{D}$

another possibility for **diversity**:

也可以让每个gt用不同特征空间，增加多样性

randomly **sample  $d'$  features** from  $\mathbf{x}$

每颗树从所有 $n$ 个特征中随机选 $d'$ 个特征做决策树。增加diversity.

一般 $d' \ll n$

- when sampling index  $i_1, i_2, \dots, i_{d'}$ :  $\Phi(\mathbf{x}) = (x_{i_1}, x_{i_2}, \dots, x_{i_{d'}})$

- $\mathcal{Z} \in \mathbb{R}^{d'}$ : a **random subspace** of  $\mathcal{X} \in \mathbb{R}^d$

相当于把 $\mathbf{x}$ 变换到子空间上

- often  $d' \ll d$ , efficient for large  $d$   
—can be generally applied on other models

原文作者建议每棵树每个分割点都只随机抽 $d'$ 个特征，从这些特征中做选择

- original RF **re-sample new subspace for each  $b(\mathbf{x})$  in C&RT**

RF = **bagging** + **random-subspace C&RT**

gt除了不同的抽样数据，特征子空间的选择也不同，增加多样性。再一起bagging，降低方差

# Diversifying by Feature Expansion

randomly **sample  $d'$  features** from  $\mathbf{x}$ :  $\Phi(\mathbf{x}) = \mathbf{P} \cdot \mathbf{x}$   
 with **row  $i$  of  $\mathbf{P}$**  sampled randomly  $\in$  **natural basis**

more **powerful** features for **diversity**: row  $i$  other than natural basis

- **projection** (combination) with random row  $\mathbf{p}_i$  of  $\mathbf{P}$ :  $\phi_i(\mathbf{x}) = \mathbf{p}_i^T \mathbf{x}$
  - often consider **low-dimensional** projection:  
only  $d''$  **non-zero** components in  $\mathbf{p}_i$
  - includes **random subspace** as **special case**:  
 $d'' = 1$  and  $\mathbf{p}_i \in$  **natural basis**
  - original RF consider  $d'$  random **low-dimensional projections for each  $b(\mathbf{x})$**  in C&RT
- 每一个分割点，都把原始的 $n$ 维特征做一个随机的低维投影，看作新特征在新的 $d'$ 维投影数据上做特征切割。可以进一步增加 $gt$ 的随机性

```
random.seed(1)
random.random()
生成同一个随机数
random.seed(1)
random.random()
同一个树，同一个随机种子？
```

随机矩阵特征抽取进一步增加随机性，用随机矩阵对原始数据投影。sign(WX-theta)

RF = **bagging** + random-**combination** C&RT  
 —**randomness** everywhere!

# Fun Time

Within RF that contains random-combination C&RT trees, which of the following hypothesis is equivalent to each branching function  $b(\mathbf{x})$  within the tree?

- ① a constant
- ② a decision stump
- ③ a perceptron
- ④ none of the other choices

$\text{sign}(\mathbf{W}\mathbf{X}-\text{theta})$  包含了斜线

# Fun Time

Within RF that contains random-combination C&RT trees, which of the following hypothesis is equivalent to each branching function  $b(\mathbf{x})$  within the tree?

- ① a constant
- ② a decision stump
- ③ a perceptron
- ④ none of the other choices

Reference Answer: ③

In each  $b(\mathbf{x})$ , the input vector  $\mathbf{x}$  is first projected by a random vector  $\mathbf{v}$  and then thresholded to make a binary decision, which is exactly what a perceptron does.



# Bagging Revisited

## Bagging

function **Bag**( $\mathcal{D}, \mathcal{A}$ )

For  $t = 1, 2, \dots, T$

- ① request size- $N'$  data  $\tilde{\mathcal{D}}_t$   
by **bootstrapping** with  $\mathcal{D}$
- ② obtain base  $g_t$  by  $\mathcal{A}(\tilde{\mathcal{D}}_t)$

return  $G = \text{Uniform}(\{g_t\})$

$g_1$ 选到  
的数据

	$g_1$	$g_2$	$g_3$	$\dots$	$g_T$
$(\mathbf{x}_1, y_1)$	$\tilde{\mathcal{D}}_1$	*	$\tilde{\mathcal{D}}_3$		$\tilde{\mathcal{D}}_T$
$(\mathbf{x}_2, y_2)$	*	*	$\tilde{\mathcal{D}}_3$		$\tilde{\mathcal{D}}_T$
$(\mathbf{x}_3, y_3)$	*	$\tilde{\mathcal{D}}_2$	*		$\tilde{\mathcal{D}}_T$
$\dots$					
$(\mathbf{x}_N, y_N)$	$\tilde{\mathcal{D}}_1$	$\tilde{\mathcal{D}}_2$	*		*

\* in  $t$ -th column: not used for obtaining  $g_t$   
—called **out-of-bag (OOB) examples** of  $g_t$

$g_t$ 没有用到的样本: out-of-bag examples

# Number of OOB Examples

OOB (in  $\star$ )  $\iff$  not sampled after  $N'$  drawings 每批  $D_t$  ~ 抽  $N'$  次样本

if  $N' = N$

- probability for  $(\mathbf{x}_n, y_n)$  to be OOB for  $g_t$ :  $(1 - \frac{1}{N})^N$
- if  $N$  large: 某样本  $N$  次都没有被抽到

$$\left(1 - \frac{1}{N}\right)^N = \frac{1}{\left(\frac{N}{N-1}\right)^N} = \frac{1}{\left(1 + \frac{1}{N-1}\right)^N} \approx \frac{1}{e}$$

OOB size per  $g_t \approx \frac{1}{e}N$

每一轮，大概有1/3左右的数据抽不到，成为OOB

## OOB versus Validation

## OOB

	$g_1$	$g_2$	$g_3$	$\dots$	$g_T$
$(\mathbf{x}_1, y_1)$	$\tilde{D}_1$	★	$\tilde{D}_3$		$\tilde{D}_T$
$(\mathbf{x}_2, y_2)$	★	★	$\tilde{D}_3$		$\tilde{D}_T$
$(\mathbf{x}_3, y_3)$	★	$\tilde{D}_2$	★		$\tilde{D}_T$
$\dots$					
$(\mathbf{x}_N, y_N)$	$\tilde{D}_1$	★	★		★

## Validation

	$g_1$	$g_2$	$\dots$	$g_M$
$D_{\text{train}}$	$D_{\text{train}}$	$D_{\text{train}}$		$D_{\text{train}}$
$D_{\text{val}}$	$D_{\text{val}}$	$D_{\text{val}}$		$D_{\text{val}}$
$D_{\text{val}}$	$D_{\text{val}}$	$D_{\text{val}}$		$D_{\text{val}}$
$D_{\text{train}}$	$D_{\text{train}}$	$D_{\text{train}}$		$D_{\text{train}}$

- ★ like  $D_{\text{val}}$ : 'enough' random examples unused during training
- use ★ to validate  $g_t$ ? easy, but **rarely needed** 用OOV验证 $g_t$ , 但没必要。  
 $g_t$ 差不多就行
- use ★ to validate  $G$ ?  $E_{\text{oob}}(G) = \frac{1}{N} \sum_{n=1}^N \text{err}(y_n, G_n^-(\mathbf{x}_n))$ ,  
with  $G_n^-$  contains only trees that  $\mathbf{x}_n$  is OOB of,

每个样本 $(\mathbf{x}_n, y_n)$

可以用来验证所有(该样本是OOV样本)的模型组成的模型 $G_n^-$

可以作为一个valid sample, 得到一个valid预测结果。

所有的样本用这样的方式, 对对应的 $G$ -做测试

取平均, 可以在一定程度上验证

整体 $G$ 的表现( $G$ -上表现的平均)

such as  $G_N^-(\mathbf{x}) = \text{average}(g_2, g_3, g_T)$

$E_{\text{oob}}$ : self-validation of bagging/RF

可以自我valid

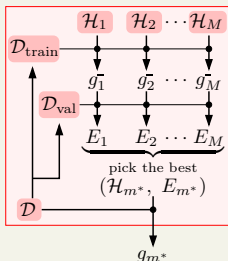
# Model Selection by OOB Error

## Previously: by Best $E_{\text{val}}$

$$g_{m^*} = \mathcal{A}_{m^*}(\mathcal{D})$$

$$m^* = \underset{1 \leq m \leq M}{\operatorname{argmin}} E_m$$

$$E_m = E_{\text{val}}(\mathcal{A}_m(\mathcal{D}_{\text{train}}))$$



## RF: by Best $E_{\text{OOB}}$

$$G_{m^*} = \text{RF}_{m^*}(\mathcal{D})$$

$$m^* = \underset{1 \leq m \leq M}{\operatorname{argmin}} E_m$$

$$E_m = E_{\text{OOB}}(\text{RF}_m(\mathcal{D}))$$

- use  $E_{\text{OOB}}$  for **self-validation** —of RF **parameters** such as  $d''$   
不需要切成2部分
- **no re-training** needed

只需要利用D中OOV数据，做self-valid,得到的  $E_{\text{OOB}}(G)$  就能较准确反映随机森林的性能。  
可以用来选择RF参数:  $g_t$  树深,  $g_t$  特征抽取维度

$E_{\text{OOB}}$  often **accurate** in practice

## Fun Time

For a data set with  $N = 1126$ , what is the probability that  $(\mathbf{x}_{1126}, y_{1126})$  is not sampled after bootstrapping  $N' = N$  samples from the data set?

- 1 0.113
- 2 0.368
- 3 0.632
- 4 0.887

## Fun Time

For a data set with  $N = 1126$ , what is the probability that  $(\mathbf{x}_{1126}, y_{1126})$  is not sampled after bootstrapping  $N' = N$  samples from the data set?

- ① 0.113
- ② 0.368
- ③ 0.632
- ④ 0.887

Reference Answer: ②

The value of  $(1 - \frac{1}{N})^N$  with  $N = 1126$  is about 0.367716, which is close to  $\frac{1}{e} = 0.367879$ .

# Feature Selection

for  $\mathbf{x} = (x_1, x_2, \dots, x_d)$ , want to remove      应用：删除冗余或不相关特征

- **redundant** features: like keeping one of 'age' and 'full birthday'
- **irrelevant** features: like insurance type for cancer prediction

and only 'learn' **subset-transform**  $\Phi(\mathbf{x}) = (x_{i_1}, x_{i_2}, x_{i_{d'}})$   
with  $d' < d$  for  $g(\Phi(\mathbf{x}))$

选了以后高效。但不好选，有可能选错或者overfitting

advantages:

- **efficiency**: simpler hypothesis and shorter prediction time
- **generalization**: 'feature noise' removed
- **interpretability**

disadvantages:

- **computation**: 组合爆炸  
难遍历
- **overfit**: 'combinatorial' optimization in training
- **mis-interpretability**

决策树本身  
是拿特征来  
分割数据。

decision tree: a rare model  
with **built-in feature selection**

相当于选了  
不同特征

# Feature Selection by Importance

idea: if possible to calculate

**importance**( $i$ ) for  $i = 1, 2, \dots, d$  按重要性选择

then can select  $i_1, i_2, \dots, i_{d'}$  of top- $d'$  **importance**

## importance by linear model

$$\text{score} = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^d w_i x_i$$

像我之前做的  
系数 $w_i$ 大的特征，比较重要

- intuitive estimate: **importance**( $i$ ) =  $|w_i|$  with some 'good'  $\mathbf{w}$
- getting 'good'  $\mathbf{w}$ : learned from data
- non-linear models? often **much harder**

next: 'easy' feature selection in RF



# Feature Importance by Permutation Test

idea: random test

—if feature  $i$  needed, ‘random’ values of  $x_{n,i}$  degrades performance

重要的维度被输入随机值，模型的表现一定会严重下降。用该维度打乱后的数据和原数据比较，考察每一维的重要程度

- which random values?

- uniform, Gaussian, ...:  $P(x_i)$  changed
- bootstrap, **permutation** (of  $\{x_{n,i}\}_{n=1}^N$ ):  $P(x_i)$  approximately remained

为了保证要测量的维度 $i$ 数据分布不变

- **permutation** test:

只将所有该维的数据打乱后，重新插入该维，只不过对应不同的数据。形成新数据 $\mathcal{D}_p$ 。表现与原始数据 $\mathcal{D}$ 差距越大，表现越差，该特征越重要

$$\text{importance}(i) = \text{performance}(\mathcal{D}) - \text{performance}(\mathcal{D}^{(p)})$$

with  $\mathcal{D}^{(p)}$  is  $\mathcal{D}$  with  $\{x_{n,i}\}$  replaced by **permuted**  $\{x_{n,i}\}_{n=1}^N$

**permutation** test: a general statistical tool for arbitrary non-linear models like RF

# Feature Importance in Original Random Forest

**permutation test:** 不重新训练，只是验证时用打乱的数据。  
每个gt, 第*i*维被gt的所有OOV数据permutation. 仍然未被污染

$$\text{importance}(i) = \text{performance}(\mathcal{D}) - \text{performance}(\mathcal{D}^{(p)})$$

with  $\mathcal{D}^{(p)}$  is  $\mathcal{D}$  with  $\{x_{n,i}\}$  replaced by **permuted**  $\{x_{n,i}\}_{n=1}^N$   
 比如gt OOV:  $x_1, x_3, x_n$   
 $x_1, x_3, x_n$ 的*i*维permutation  
 $gt+1$ , 也是OOV数据的第*i*维交换

- **performance**( $\mathcal{D}^{(p)}$ ): needs re-training and **validation** in general
- ‘**escaping**’ **validation**? **OOB** in RF 本来需要用 $\mathcal{D}_p$ , 重新训练+测试, 完整过一遍

- original RF solution:  $\text{importance}(i) = E_{\text{ooB}}(G) - E_{\text{ooB}}^{(p)}(G)$ ,  
 where  $E_{\text{ooB}}^{(p)}$  comes from replacing each request of  $x_{n,i}$  by a

**permuted OOB** value 偷个懒, 不重新训练得到 $G(\mathcal{D}_p)$ 了, 还用 $G$   
 $E_{\text{ooB}}(G_p) \rightarrow E_{\text{ooB}}(p)(G)$ : 只是验证时用permutation后的新数据

原来算 $E_{\text{ooB}}$ 时, 对于 $x_n$ , 需要算很多gt( $x_n$ ) 得到 $G_n$ -  $x_n$ 对gt来说, 是OOV的数据  
 这里 $x_n$ 用新数据算。用*i*维数据切分时,  $x_{n,i}$ 是gt所有的OOV数据*i*维的取值。

用来做特征选择, 特别好用!

**RF feature selection via permutation + OOB:**  
 often efficient and promising in practice

这样可以保证替换后的 $x_{n,i}$ 对gt来说, 仍然是OOV数据, 未被污染过

## Fun Time

For RF, if the 1126-th feature within the data set is a constant 5566, what would  $\text{importance}(i)$  be?

- 1 0
- 2 1
- 3 1126
- 4 5566

# Fun Time

For RF, if the 1126-th feature within the data set is a constant 5566, what would importance( $i$ ) be?

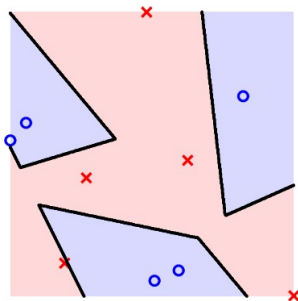
- 1 0
- 2 1
- 3 1126
- 4 5566

Reference Answer: 1

When a feature is a constant, permutation does not change its value. Then,  $E_{\text{oob}}(G)$  and  $E_{\text{oob}}^{(p)}(G)$  are the same, and thus importance( $i$ ) = 0.

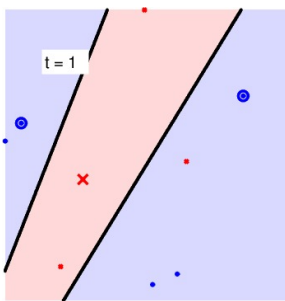
# A Simple Data Set

$g_{C\&RT}$   
with random combination



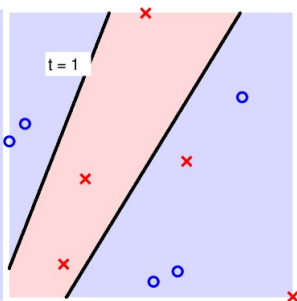
单一决策树  
只是切分时候用的是随机组合的特征  
所以是斜线

$g_t (N' = N/2)$



bootstrap 一半的数据量  
扔到同一个 $g_t$ 里。  
大圈是bootstrap sample到的样本

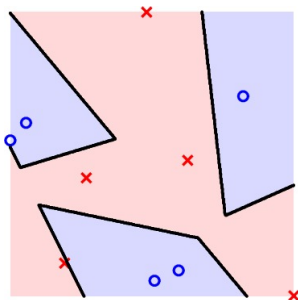
$G$  with first  $t$  trees



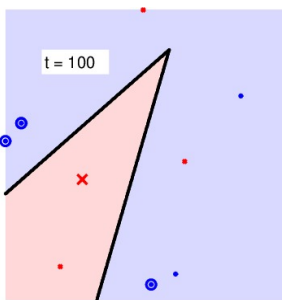
$G \quad t=1$  一棵树 同左图

# A Simple Data Set

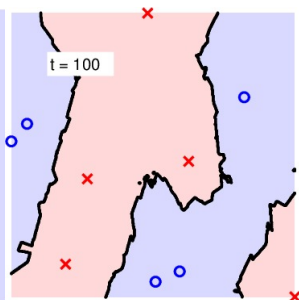
$g_{C\&RT}$   
with random combination



$g_t (N' = N/2)$



$G$  with first  $t$  trees

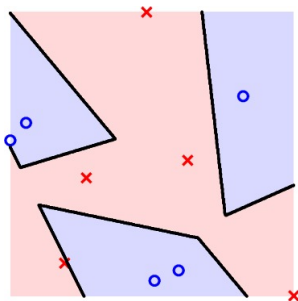


一棵树bagging  
t=100次取平均  
(只改变bootstrap样本  $N'=N/2$ )

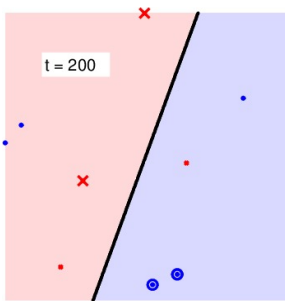
随机森林  
100颗树

# A Simple Data Set

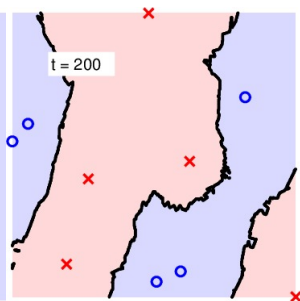
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$



$G$  with first  $t$  trees

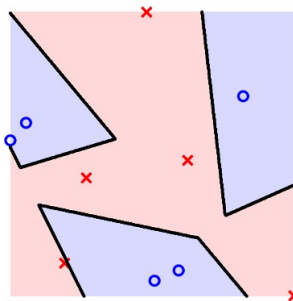


bagging

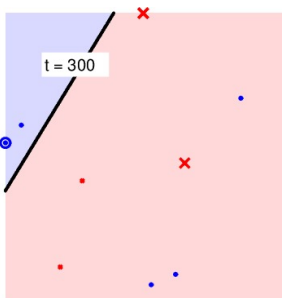
很稳定 有些错

# A Simple Data Set

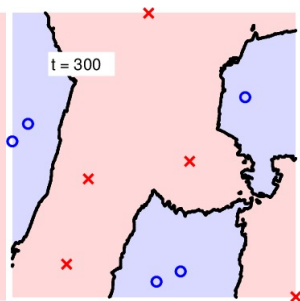
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$



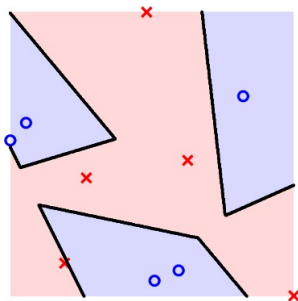
$G$  with first  $t$  trees



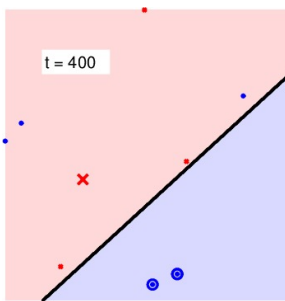


# A Simple Data Set

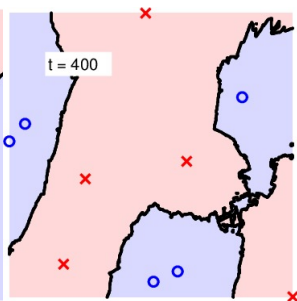
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

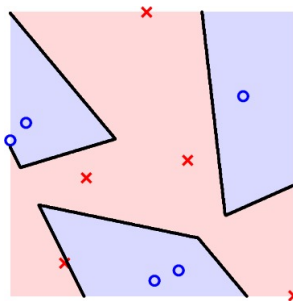


$G$  with first  $t$  trees

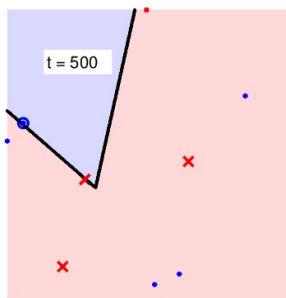


# A Simple Data Set

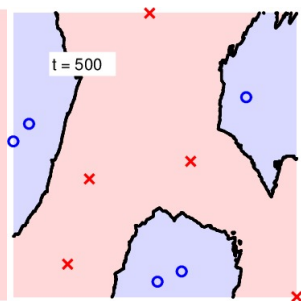
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

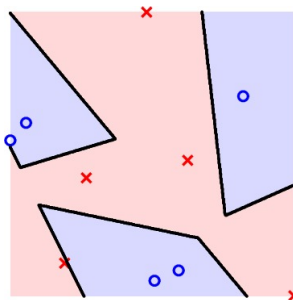


$G$  with first  $t$  trees

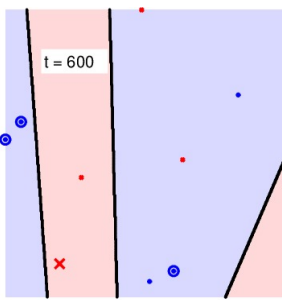


# A Simple Data Set

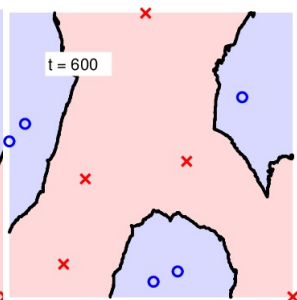
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

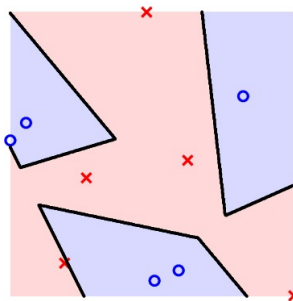


$G$  with first  $t$  trees

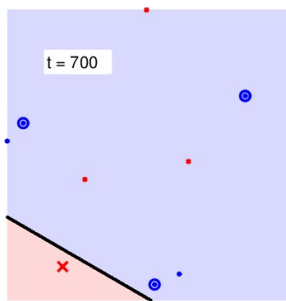


# A Simple Data Set

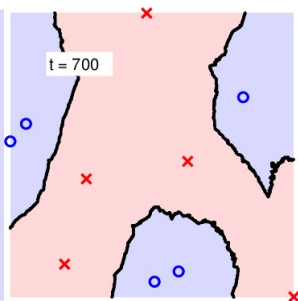
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

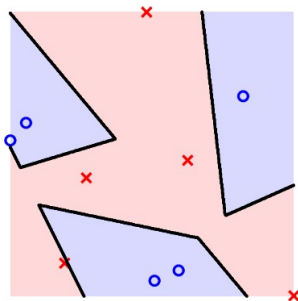


$G$  with first  $t$  trees

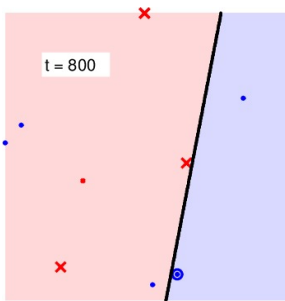


# A Simple Data Set

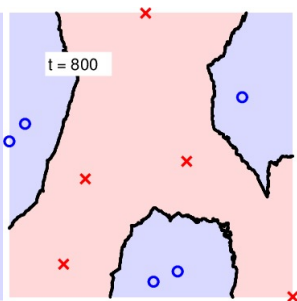
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

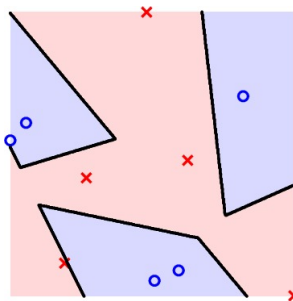


$G$  with first  $t$  trees

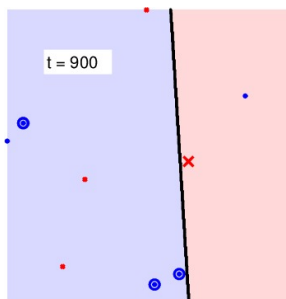


# A Simple Data Set

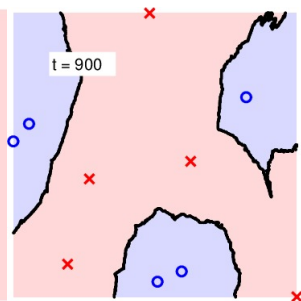
$g_{\text{C\&RT}}$   
with random combination



$g_t (N' = N/2)$

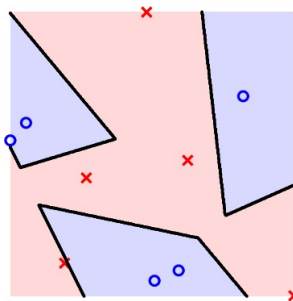


$G$  with first  $t$  trees



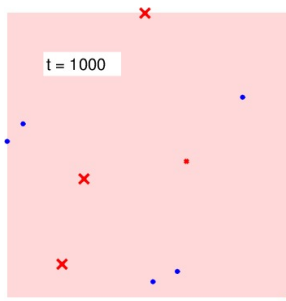
# A Simple Data Set

$g_{C\&RT}$   
with random combination



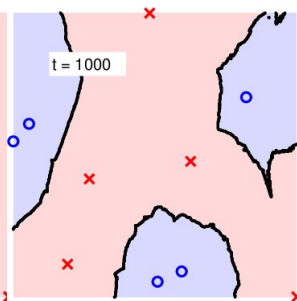
单颗树

$g_t (N' = N/2)$



一颗树 bagging 1000次

$G$  with first  $t$  trees

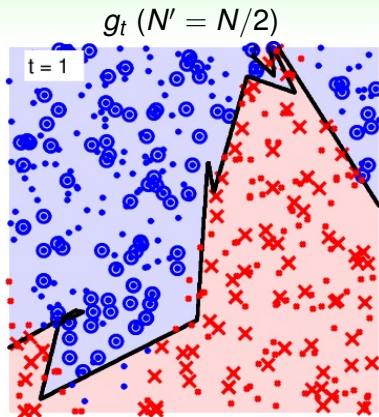


随机森林 1000颗

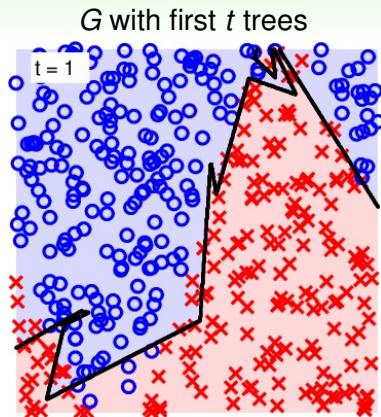
随机森林相比一棵树。边界更平滑，而且类似SVM, 分类间隔最大。远好于一棵树用bootstrap做bagging

**‘smooth’ and large-margin-like boundary  
with many trees**

# A Complicated Data Set



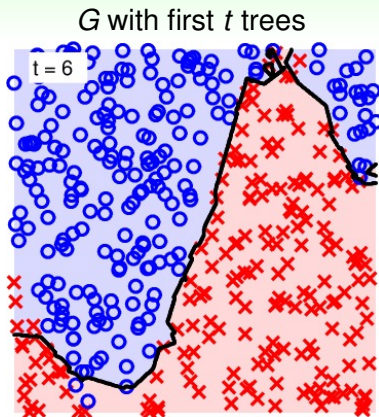
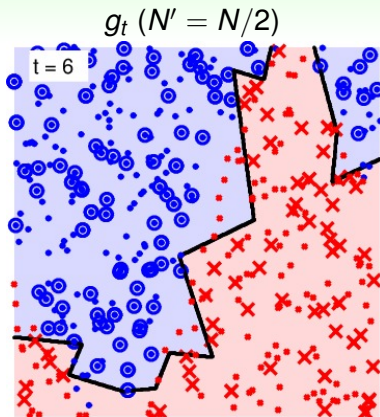
一棵树 bagging



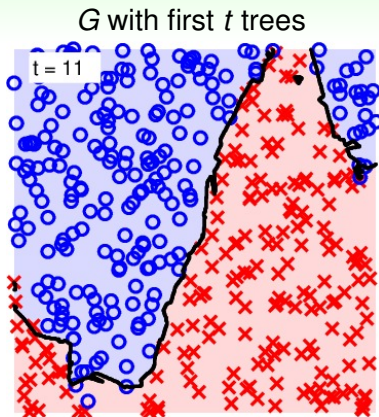
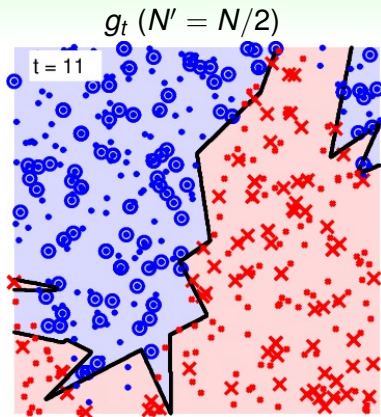
随机森林



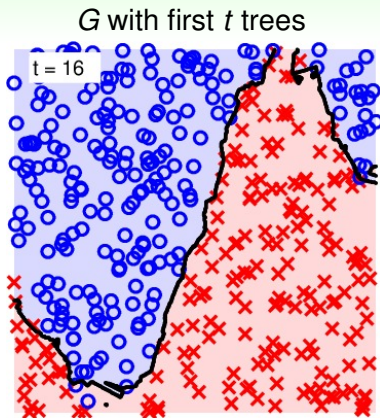
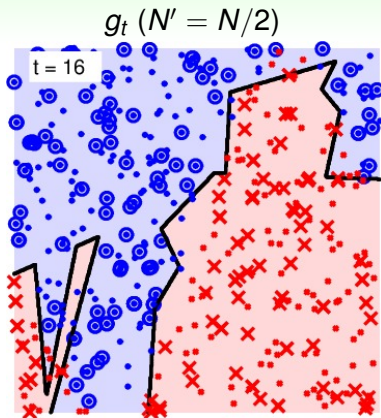
# A Complicated Data Set



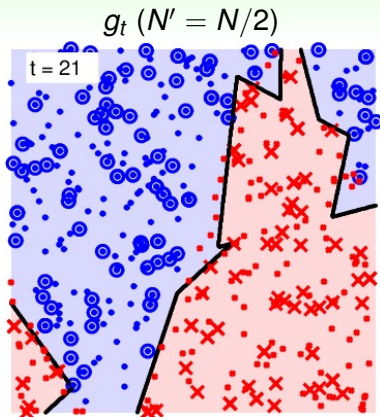
# A Complicated Data Set



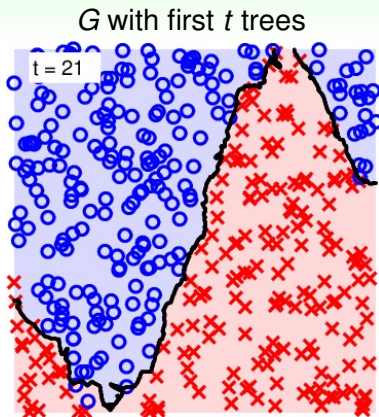
# A Complicated Data Set



# A Complicated Data Set



一棵树 bagging

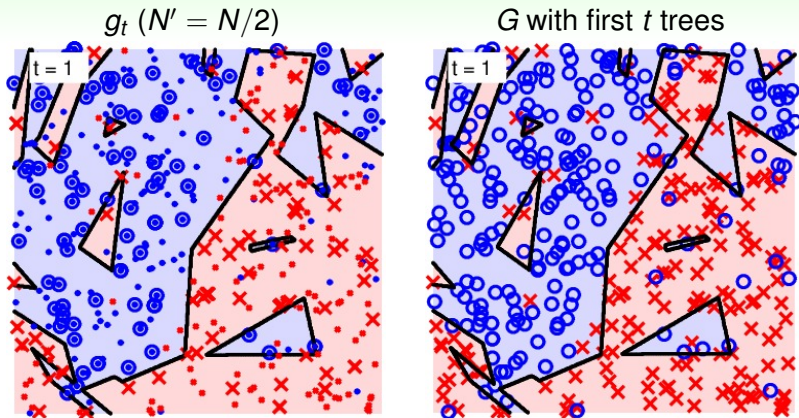


随机森林 多棵树bagging

平滑, 鲁棒, 非线性很多

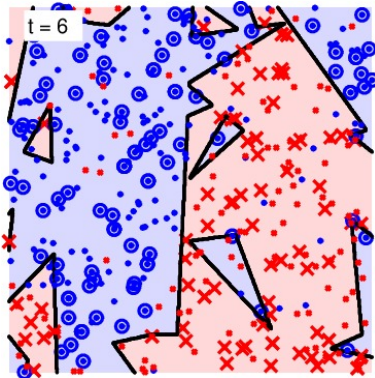
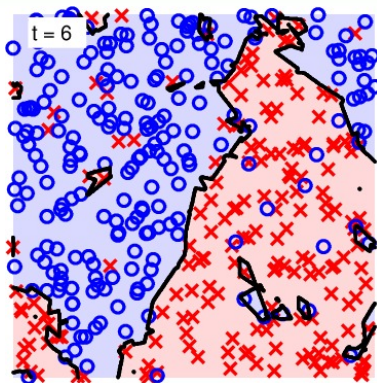
**'easy yet robust' nonlinear model**

# A Complicated and Noisy Data Set



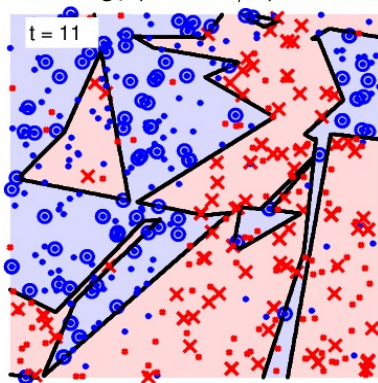
帶有噪音的数据。一棵树容易overfitting

# A Complicated and Noisy Data Set

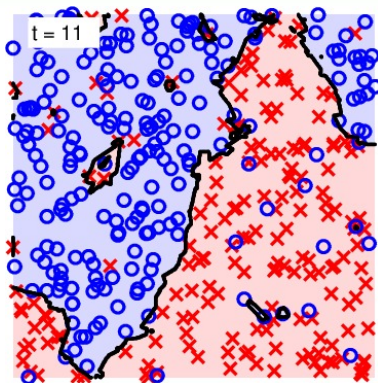
 $g_t (N' = N/2)$  $G$  with first  $t$  trees

# A Complicated and Noisy Data Set

$g_t (N' = N/2)$

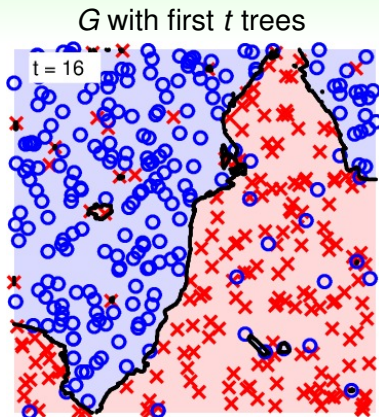
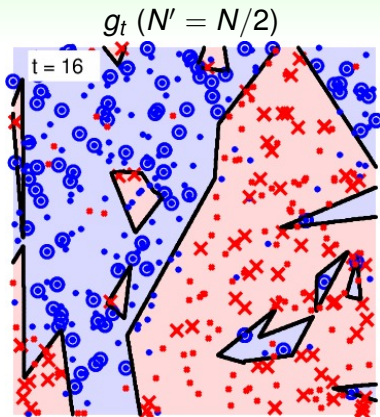


$G$  with first  $t$  trees



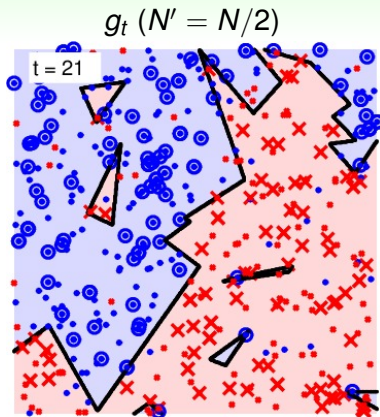


# A Complicated and Noisy Data Set

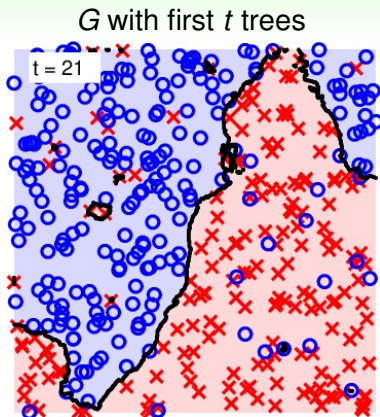




# A Complicated and Noisy Data Set



一棵树bagging (bootstrap抽样多次)



随机森林

可以在一定程度上通过多棵树投票，消除噪音，得到稳定的分类边界。不容易过拟合

**noise corrected** by voting

# How Many Trees Needed?

almost every theory: the more, **the 'better'**  
assuming **good**  $\bar{g} = \lim_{T \rightarrow \infty} G$

## Our NTU Experience

- KDDCup 2013 Track 1 (**yes, NTU is world champion again! :-)**): predicting author-paper relation
- $E_{\text{val}}$  of **thousands** of trees: [0.015, 0.019] depending **on seed**;  $E_{\text{out}}$  of top 20 teams: [0.014, 0.019]
- decision: take **12000 trees** with **seed 1**

结果不是特别稳定，受随机性的影响大。可能需要增加树的数目，来让结果保持稳定。  
树够多，结果稳定。 同时选择好的种子

cons of RF: may need lots of trees **if the whole random process too unstable**  
—should double-check **stability of  $G$**   
to ensure **enough trees**

## Fun Time

Which of the following is **not** the best use of Random Forest?

- ① train each tree with bootstrapped data
- ② use  $E_{\text{oob}}$  to validate the performance
- ③ conduct feature selection with permutation test
- ④ fix the number of trees,  $T$ , to the lucky number 1126

## Fun Time

Which of the following is **not** the best use of Random Forest?

- ① train each tree with bootstrapped data
- ② use  $E_{\text{oob}}$  to validate the performance
- ③ conduct feature selection with permutation test
- ④ fix the number of trees,  $T$ , to the lucky number 1126

Reference Answer: ④

A good value of  $T$  can depend on the nature of the data and the stability of the whole random process.

# Summary

- 1 Embedding Numerous Features: Kernel Models
- 2 Combining Predictive Features: Aggregation Models

## Lecture 10: Random Forest

- Random Forest Algorithm  
**bag of trees on randomly projected subspaces**
- Out-Of-Bag Estimate  
**self-validation with OOB examples**
- Feature Selection  
**permutation test for feature importance**
- Random Forest in Action  
**'smooth' boundary with many trees**

- **next: boosted decision trees beyond classification**

- 3 Distilling Implicit Features: Extraction Models