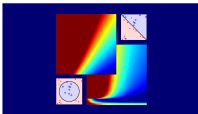


Machine Learning Techniques

(機器學習技巧)



Lecture 10: Random Forest

Hsuan-Tien Lin (林軒田)

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)



Agenda

Lecture 10: Random Forest

- Basic Random Forest Algorithm
- Out-Of-Bag Estimate
- Feature Selection

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

—reduces variance
by voting/averaging

Decision Tree

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 $G(\mathbf{x}) =$

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

—large variance
especially if fully-grown

putting them together?

(i.e. aggregate two aggregation models :-)

Random Forest

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 base G_t by **DTree**($\tilde{\mathcal{D}}_t$)

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

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

- highly parallel/efficient to learn
- inherit pros of C&RT
- eliminate cons of fully-grown tree

Diversifying by Feature Projection

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

randomly sample N examples from \mathcal{D}

other possibility:

randomly sample d' features from \mathbf{x}

- chosen 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 $\mathcal{X} \in \mathbb{R}^d$
- often $d' \ll d$, efficient when d large
— can be generally used for other learning models
- original RF re-sample new subspace for each $b(\mathbf{x})$ in C&RT

RF = bagging + random-subspace C&RT

Diversifying by Feature Expansion

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

more powerful features: row i of P other than natural basis

- low-dimensional random projection (combination) with \mathbf{u} :

$$\phi_i(\mathbf{x}) = \sum_{j=1}^{d''} u_j x_j$$

- includes random subspace as a special case: $d'' = 1$ and $u_1 = 1$
- original RF consider d' random projections for each $b(\mathbf{x})$ in C&RT

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

Fun Time

Bagging Revisited

Bagging

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

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	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)	*	$\tilde{\mathcal{D}}_1$	*		$\tilde{\mathcal{D}}_T$
\dots					
(\mathbf{x}_N, y_N)	$\tilde{\mathcal{D}}_1$	$\tilde{\mathcal{D}}_2$	*		*

*: not used for obtaining g_t
 —called **out-of-bag (OOB) examples**

Number of OOB Examples

OOB (in \star) \iff not sampled after N drawings

- probability for (\mathbf{x}_n, y_n) to be OOB for g_t : $(1 - \frac{1}{N})^N$
- if N large:

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

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

OOB versus Validation

OOB

	g_1	g_2	g_3	...	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}_1	*		\tilde{D}_T
...					
(\mathbf{x}_N, y_N)	\tilde{D}_1	\tilde{D}_2	*		*

Validation

	g_1^-	g_2^-	...	g_M^-
	D_{train}	D_{train}		D_{train}
	D_{val}	D_{val}		D_{val}
	D_{val}	D_{val}		D_{val}
	D_{train}	D_{train}		D_{train}

- * like D_{val} : 'enough' random examples unused during training
- use * to validate g_t ? easy, but **rarely needed (why?)**
- 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

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

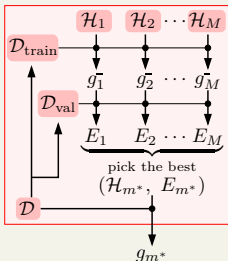
Model Selection by OOB Error

Previously: by Best E_{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_{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_{OOB} for **self-validation**
- **no re-training** needed

E_{OOB} often **accurate** in practice

Fun Time

Feature Selection

for $\mathbf{x} = (x_1, x_2, \dots, x_d)$, want to remove

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

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

advantages:

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

disadvantages:

- computation: 'combinatorial' optimization in training
- overfit: 'combinatorial' selection
- mis-interpretability

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

Feature Selection by Importance

idea: if possible to estimate

importance(i) for $i = 1, 2, \dots, d$

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

Linear Model

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

- intuitive estimate: importance(i) = $|w_i|$ with some 'good' \mathbf{w}
- 'good' \mathbf{w} : learned with full data
- non-linear models? often not easy

next: feature selection in RF

Feature Importance by Permutation Test

idea: random test

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

- which random values?
 - uniform, Gaussian, . . . : $P(x_i)$ changed
 - bootstrap, permutation (of $\{x_{n,i}\}_{n=1}^N$): $P(x_i) \approx$ remained
- permutation test:

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

with \mathcal{D}_p containing permuted $\{x_{n,i}\}_{n=1}^N$

permutation test: a general statistical tool that
can be used for arbitrary non-linear models like
RF

Feature Importance in Original Random Forest

permutation test:

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

with \mathcal{D}_p containing permuted $\{x_{n,i}\}_{n=1}^N$

- calculating performance needs re-training and validating on each \mathcal{D}_p in general
- how to 'escape' validation? OOB in RF
- original RF solution:

$$\text{importance}(i) = E_{\text{ooB}}(G, \mathcal{D}) - E_{\text{ooB}}(G, \mathcal{D}_p)$$

with \mathcal{D}_p 'dynamically' containing permuted $\{x_{n,i} : n \text{ OOB}\}$ for g_t

original RF solution often efficient and promising in practice

Fun Time

Summary

Lecture 10: Random Forest

- Basic Random Forest Algorithm
- Out-Of-Bag Estimate
- Feature Selection