

**Comparison and identification for rhizomes and leaves of *Paris yunnanensis* based on Fourier transform mid infrared spectroscopy combined with chemometrics**

Yi-Fei Pei <sup>a, b</sup>, Qing-Zhi Zhang <sup>b</sup>, Zhi-Tian Zuo <sup>a\*</sup> and Yuan-Zhong Wang <sup>a\*</sup>

<sup>a</sup> Institute of Medicinal Plants, Yunnan Academy of Agricultural Sciences, Kunming 650200, PR China.

<sup>b</sup> College of Traditional Chinese Medicine, Yunnan University of Traditional Chinese Medicine, Kunming 650500, PR China.

\*Corresponding author:

**Mr Yuan-Zhong Wang** and **Mrs Zhi-Tian Zuo**, Institute of Medicine Plants, Yunnan Academy of Agricultural Science, 2238, Beijing Road, Panlong District, Kunming 650200, China. Tel: +86 871-65033575, Fax: +86 871-65033441, E-mail: boletus@126.com (**Mr Yuan-Zhong Wang**), yaaszztian@126.com (**Mrs Zhi-Tian Zuo**)

## Figure captions

Fig. S1 VIP scores of FT-MIR data of leaves for regional difference: (a) raw dataset, (b) SNV-SD dataset.

Fig. S2 The  $n_{\text{tree}}$  and  $m_{\text{try}}$  screening of RF models of *P. yunnanensis* samples before variables ranked by permutation accuracy importance: (a)  $n_{\text{tree}}$  of raw leaves dataset, (b)  $n_{\text{tree}}$  of SNV-SD leaves dataset, (c)  $m_{\text{try}}$  of raw leaves dataset, (d)  $m_{\text{try}}$  of SNV-SD leaves dataset.

Fig. S3 The 10-fold cross validation error rates of RF model (sequentially reduce each five variables) based on *P. yunnanensis* samples: (a) raw leaves dataset, (b) SNV-SD leaves dataset.

Fig. S4 The  $n_{\text{tree}}$  and  $m_{\text{try}}$  screening of RF models of *P. yunnanensis* samples after variables ranked by permutation accuracy importance: (a)  $n_{\text{tree}}$  of raw leaves dataset, (b)  $n_{\text{tree}}$  of SNV-SD leaves dataset, (c)  $m_{\text{try}}$  of raw leaves dataset, (d)  $m_{\text{try}}$  of SNV-SD leaves dataset.

Fig. S5 The  $n_{\text{tree}}$  and  $m_{\text{try}}$  screening of RF models of *P. yunnanensis* samples before variables ranked by permutation accuracy importance: (a)  $n_{\text{tree}}$  of raw data fusion dataset, (b)  $n_{\text{tree}}$  of SNV-SD data fusion dataset, (c)  $m_{\text{try}}$  of raw data fusion dataset, (d)  $m_{\text{try}}$  of SNV-SD data fusion dataset.

Fig. S6 The 10-fold cross validation error rates of RF model (sequentially reduce each five variables) based on *P. yunnanensis* samples: (a) raw data fusion dataset, (b) SNV-SD data fusion dataset.

Fig. S7 The  $n_{\text{tree}}$  and  $m_{\text{try}}$  screening of RF models of *P. yunnanensis* samples after variables ranked by permutation accuracy importance: (a)  $n_{\text{tree}}$  of raw data fusion dataset, (b)  $n_{\text{tree}}$  of SNV-SD data fusion dataset, (c)  $m_{\text{try}}$  of raw data fusion dataset, (d)  $m_{\text{try}}$  of SNV-SD data fusion dataset.

## Table captions

Table S1 The major parameters of raw and preprocessing **calibration** models based on *P. yunnanensis* samples combined with rhizomes FT-MIR spectra.

Table S2 The major parameters of PLS-DA and RF models of each class based on raw and SNV-SD rhizomes FT-MIR spectra datasets of *P. yunnanensis* samples.

Table S3 The major parameters of raw and preprocessing **calibration** models based on *P. yunnanensis* samples combined with leaves FT-MIR spectra.

Table S4 The major parameters of PLS-DA and RF models of each class based on raw and SNV-SD leave FT-MIR spectra datasets of *P. yunnanensis* samples.

**Table S5 The geographical location of *P. yunnanensis* samples.**

**Table S6 The sample size of calibration set and validation set for each class.**

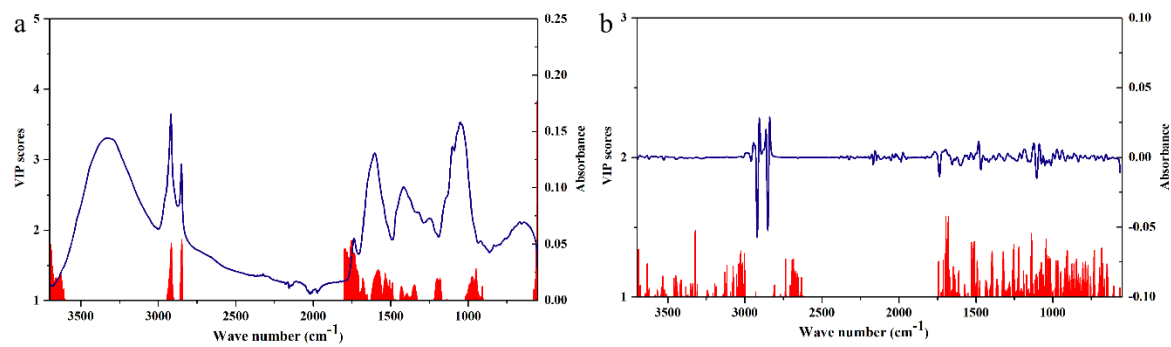


Fig. S1

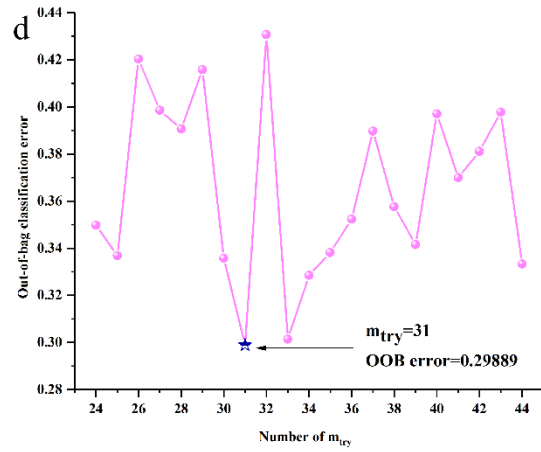
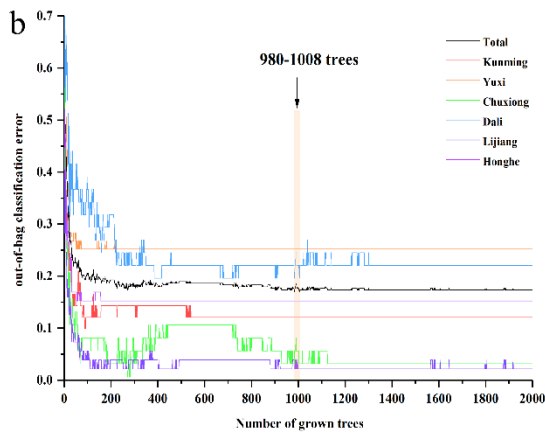
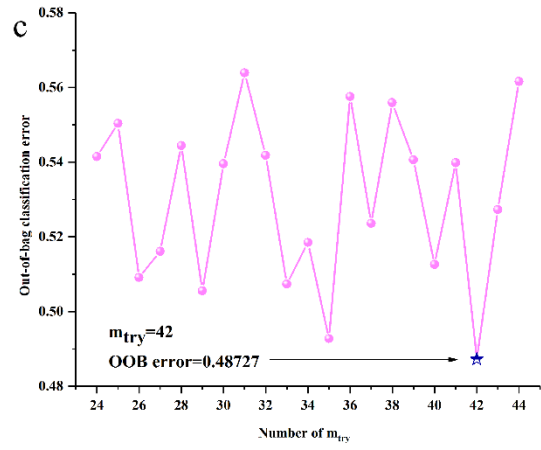
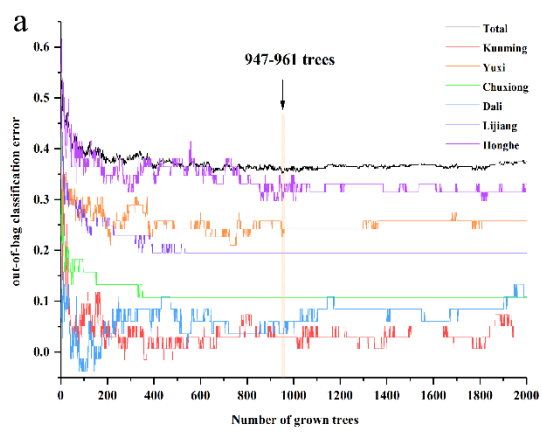


Fig. S2

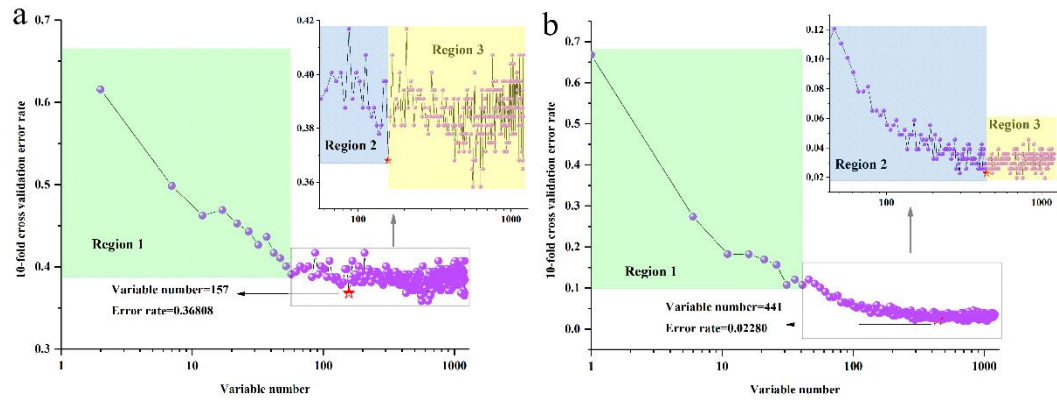


Fig. S3

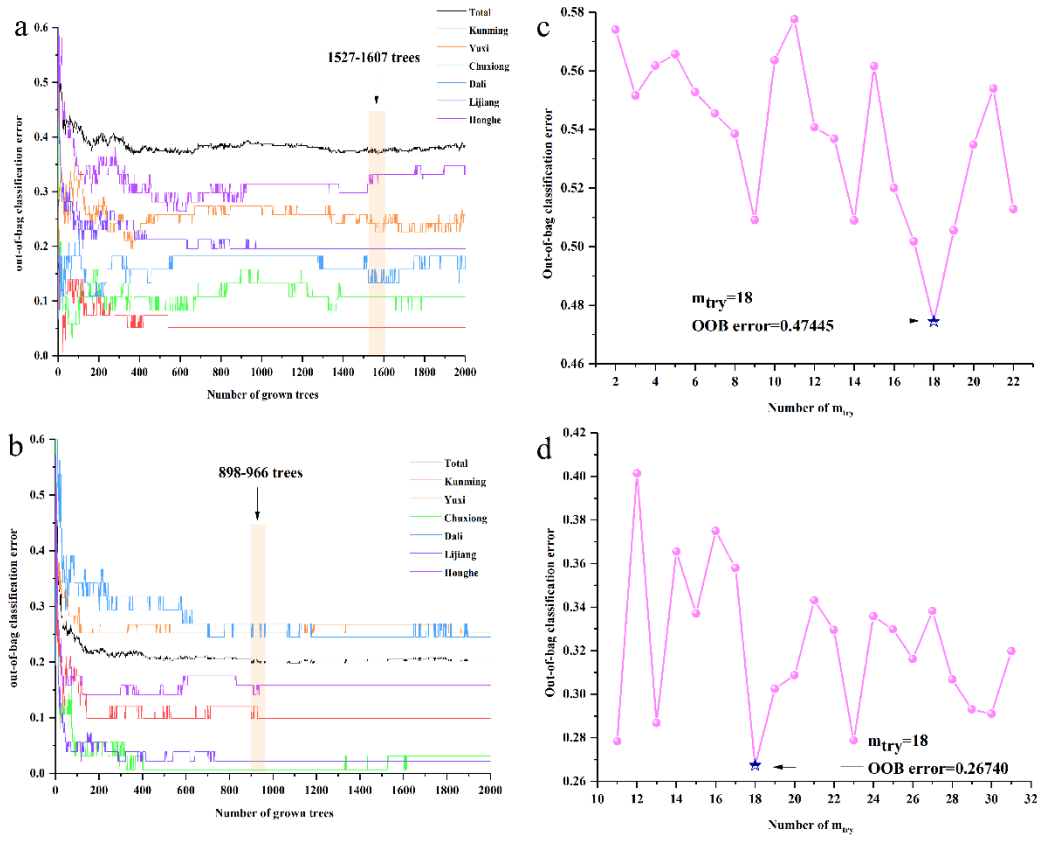


Fig. S4

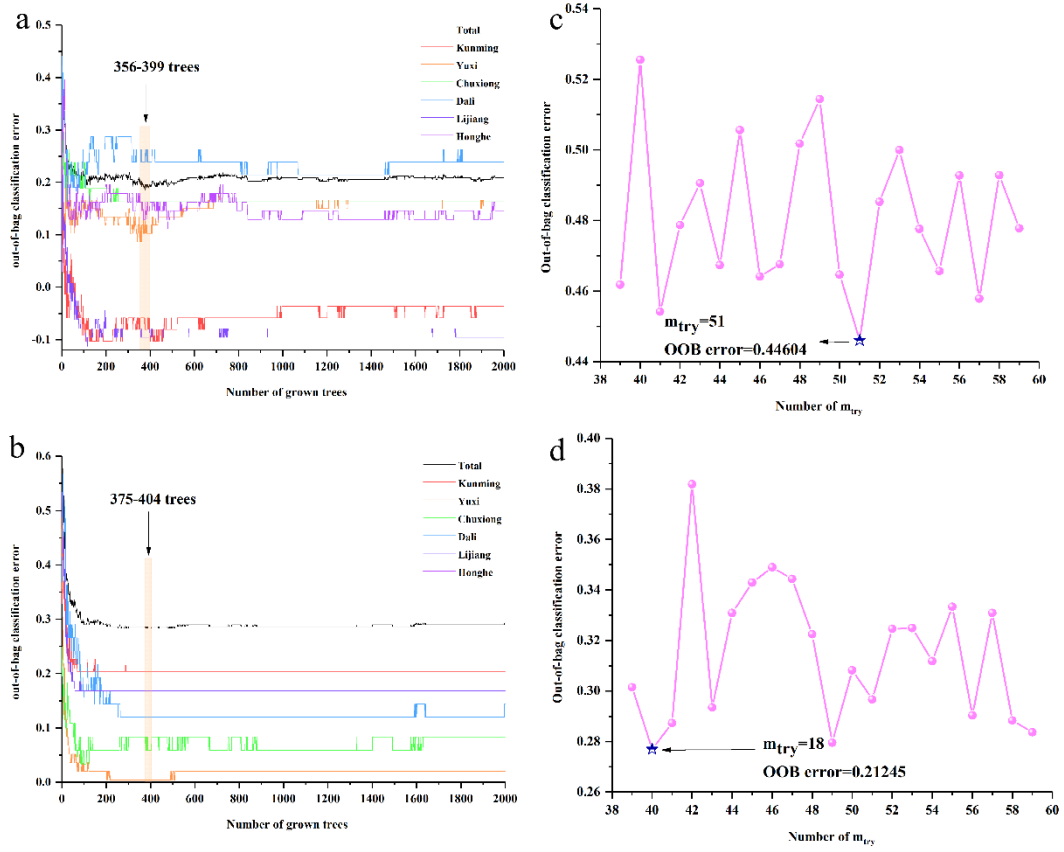


Fig. S5



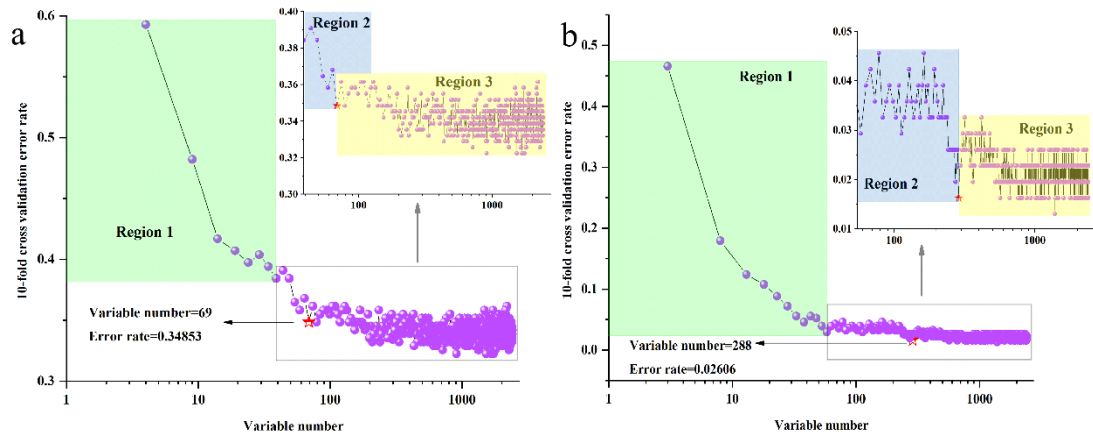


Fig. S6

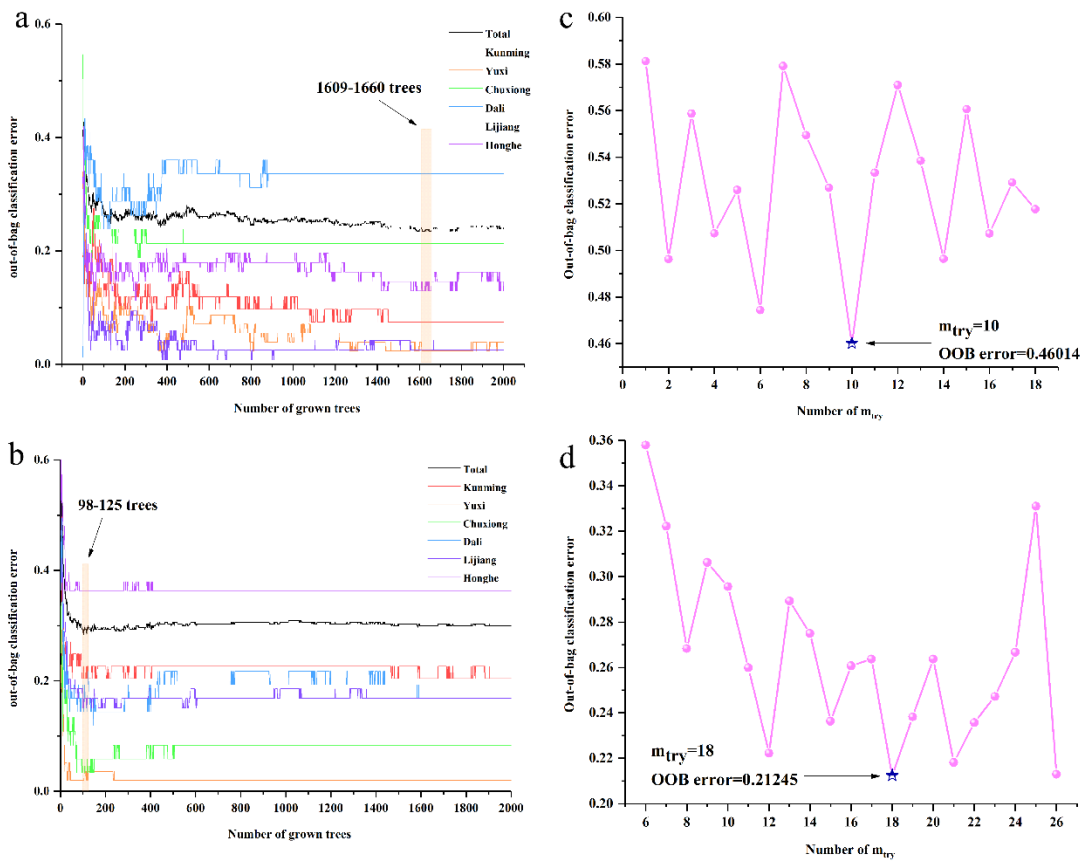


Fig. S7

Table S1

Preprocessing	LVs	R <sup>2</sup>	Q <sup>2</sup>	RMSEE	RMSECV	Accuracy
Raw	18	0.714	0.584	0.201239	0.250368	95.44%
SNV	18	0.696	0.56	0.208316	0.254995	95.77%
SNV-FD	19	0.775	0.614	0.179515	0.239488	99.02%
SNV-SD	14	0.816	0.674	0.162205	0.219739	99.67%
SD	16	0.824	0.63	0.174812	0.237413	100%

Table S2

Preprocessing	Set	Classes <sup>a</sup>	PLS-DA				RF			
			SENS	SPEC	ACC	MCC	SENS	SPEC	ACC	MCC
Raw	Calibration	1	1	1	1	1	0.6	0.962	0.909	0.61
		2	0.968	0.992	0.987	0.96	0.698	0.955	0.902	0.688
		3	0.775	1	0.971	0.866	0.625	0.974	0.928	0.66
		4	0.927	0.992	0.984	0.929	0.561	0.951	0.899	0.541
		5	1	0.972	0.977	0.931	0.69	0.9	0.86	0.565
		6	1	0.988	0.99	0.97	0.764	0.861	0.844	0.553
	Validation	1	1	0.992	0.994	0.975	0.87	0.977	0.961	0.847
		2	0.75	0.992	0.942	0.816	0.469	1	0.89	0.642
		3	0.619	1	0.948	0.764	0.238	1	0.897	0.461
		4	0.95	1	0.994	0.971	0.8	1	0.974	0.881
		5	1	0.937	0.948	0.857	0.931	0.849	0.865	0.666
		6	1	0.944	0.955	0.875	0.9333	0.824	0.845	0.64
SNV-SD	Calibration	1	1	1	1	1	0.933	0.992	0.984	0.934
		2	1	1	1	1	0.921	0.996	0.98	0.939
		3	0.975	1	0.997	0.986	0.8	0.989	0.964	0.835
		4	1	1	1	1	0.854	0.981	0.964	0.844
		5	1	0.996	0.997	0.989	0.897	0.948	0.938	0.809
		6	1	1	1	1	0.967	0.976	0.974	0.92
	Validation	1	1	1	1	1	1	1	1	1
		2	1	1	1	1	0.969	1	0.994	0.98
		3	0.905	1	0.987	0.944	0.905	1	0.987	0.944
		4	1	0.993	0.994	0.972	0.85	0.993	0.974	0.882
		5	1	1	1	1	0.966	0.96	0.961	0.882
		6	1	0.992	0.994	0.98	1	0.992	0.994	0.98

Table S3

Preprocessing	LVs	R <sup>2</sup>	Q <sup>2</sup>	RMSEE	RMSECV	Accuracy
Raw	18	0.698	0.559	0.203855	0.203855	95.77%
SNV	17	0.719	0.549	0.198488	0.256876	97.72%
SNV-FD	17	0.807	0.622	0.165141	0.240036	100%
SNV-SD	15	0.876	0.754	0.133634	0.195234	100%
SD	14	0.837	0.704	0.151242	0.210209	99.67%



Table S5

Region	Location	Sample size	Latitude (°N)	Longitude (°E)
1	Wuhua, Kunming	68	25.042165	102.704412
2	Hongta, Yuxi	95	24.43105	102.44098
3	Yaoan, Chuxiong	61	25.522293	101.375931
4	Weishan, Dali	61	25.307049	100.316085
5	Gucheng, Lijiang	87	26.874046	100.190409
6	Yuanyang, Honghe	90	23.007286	103.025416

Table S6

	Class 1	Class 2	Class 3	Class 4	Class 5	Class 6
Calibration set	45	63	40	41	58	60
Validation set	23	32	21	20	29	30
Total	68	95	61	61	87	90