

Canonical Correlation Forests

Tom Rainforth, Frank Wood

{twgr,fwood}@robots.ox.ac.uk, www.robots.ox.ac.uk/~twgr



UNIVERSITY OF
OXFORD

Replacement for Random Forest

- Canonical correlation forests (CCFs) [1] are a new decision tree ensemble learning method for classification
- CCFs are composed of canonical correlation trees (CCTs) which use hyperplane splits based on canonical correlation components
- CCFs require no parameter tuning
- CCFs outperform both the state-of-art tree ensemble methods random forest (RF) [2] and rotation forest [3] significantly

Trees and Forests Review

- Decision trees hierarchically divide the input space and assign local models to the leafs
- Classical decision tree training algorithms greedily search the possible space of axis-aligned unique splits
- Combining individual trees to form a forest improves performance

Canonical Correlation Analysis

- CCA [4] is used to give pairs of projections that maximise the correlation between the features X and the class labels Y
 $\{\Phi, \Omega\} = \text{CCA}(X, Y) = \text{argmax}_{a,b} (\text{corr}(Xa, Yb))$
- $\min(\text{rank}(X), \text{rank}(Y))$ pairs are produced by adding the constraint that new components are uncorrelated with previous components
- e.g.

$$\left\{ \begin{bmatrix} -2.11 & -2.49 \\ 0.52 & 0.93 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ -1.34 & -2.49 \\ -2.28 & -0.38 \end{bmatrix} \right\} = \text{CCA} \left(\begin{bmatrix} 1 & 0.5 \\ 2 & 2 \\ 3 & 4.5 \\ 4 & 8 \\ 5 & 12.5 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \right)$$

- Unaffected by affine transformation or rotation

CCF Training

Algorithm 1 Canonical correlation forest training algorithm

```
1: Inputs:  $X^r \in \mathbb{R}^{N \times D^r}$ ,  $X^c \in \mathbb{S}^{N \times D^c}$ ,  $\mathcal{Y} \in \mathbb{I}^{N \times K}$ ,  $L \in \mathbb{Z}^+$ ,  $\lambda \in \mathbb{Z}^+$ ,  $\delta$  represents non-ordinal space
2: Convert  $X^c$  to 1-of-K encoding  $X^b \in \mathbb{I}^{N \times D^b}$ ,  $X = \{X^r, X^b\}$ 
3:  $\mu_{(d)} = \sum_{n=1}^N X_{(n,d)} / N$ ,  $\sigma_{(d)} = \sqrt{\sum_{n=1}^N X_{(n,d)}^2 - \mu_{(d)}^2 / (N-1)}$   $\forall d$ 
4:  $X_{(i,d)} \leftarrow (X_{(i,d)} - \mu_{(d)}) / \sigma_{(d)}$   $\forall d$ , set missing values in  $X$  to 0
5: if  $\lambda < (D^r + D^b)$  then  $b = \text{true}$  else  $b = \text{false}$  end if
6: for  $i = 1; L$  do
7:   if  $b$  then  $\{X', \mathcal{Y}'\} \leftarrow \{X, \mathcal{Y}\}$  else  $\{X', \mathcal{Y}'\} \leftarrow$  sample with replacement  $N$  rows from  $\{X, \mathcal{Y}\}$  end if
8:    $[\cdot, \Psi, \Theta] = \text{GROWTREE}(X', \mathcal{Y}', \{1, \dots, D^r + D^b\}, \lambda, b)$ 
9:    $t_i = \{\Psi, \Theta\}$ 
10: end for
11: return  $T = \{t_i\}_{i=1, \dots, L}$ 
```

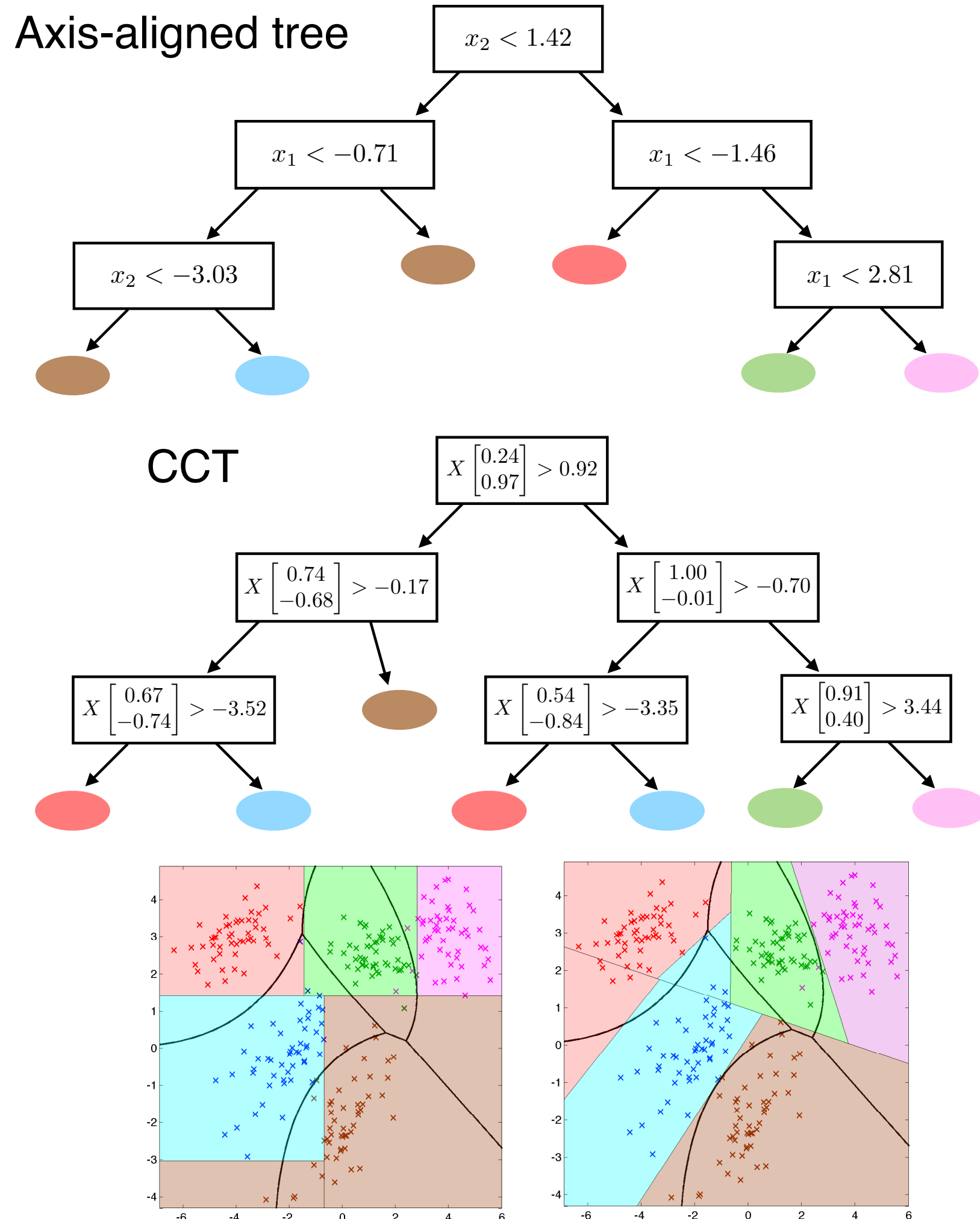
References

- Tom Rainforth, and Frank Wood. Canonical Correlation Forests, 2015. Preprint and code available at <http://robots.ox.ac.uk/~twgr>
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- Juan José Rodríguez, Ludmila I Kuncheva, and Carlos J Alonso. Rotation forest: A new classifier ensemble method. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(10):1619–1630, 2006.
- Magnus Borga. Canonical correlation: a tutorial. *On line tutorial* <http://people.imt.liu.se/magnus/cca>, 4, 2001.
- Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- Tin Kam Ho. The random subspace method for constructing decision forests. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(8):832–844, 1998.
- M. Lichman. UCI machine learning repository, 2013. <https://archive.ics.uci.edu/ml/datasets.html>

Coordinate Free Splitting

- CCTs allow non axis aligned partitions, $X\phi > s$ for features X , projection vector ϕ and split point s

Axis-aligned tree



Corresponding input space partitionings with Bayes optimal decision surface in bold

Projection Bootstrap - A New Idea

- Instead of using tree bagging [5], perform CCA on a local bootstrap sample of the data
- Use all data points to select the split point and best projection in the projected space

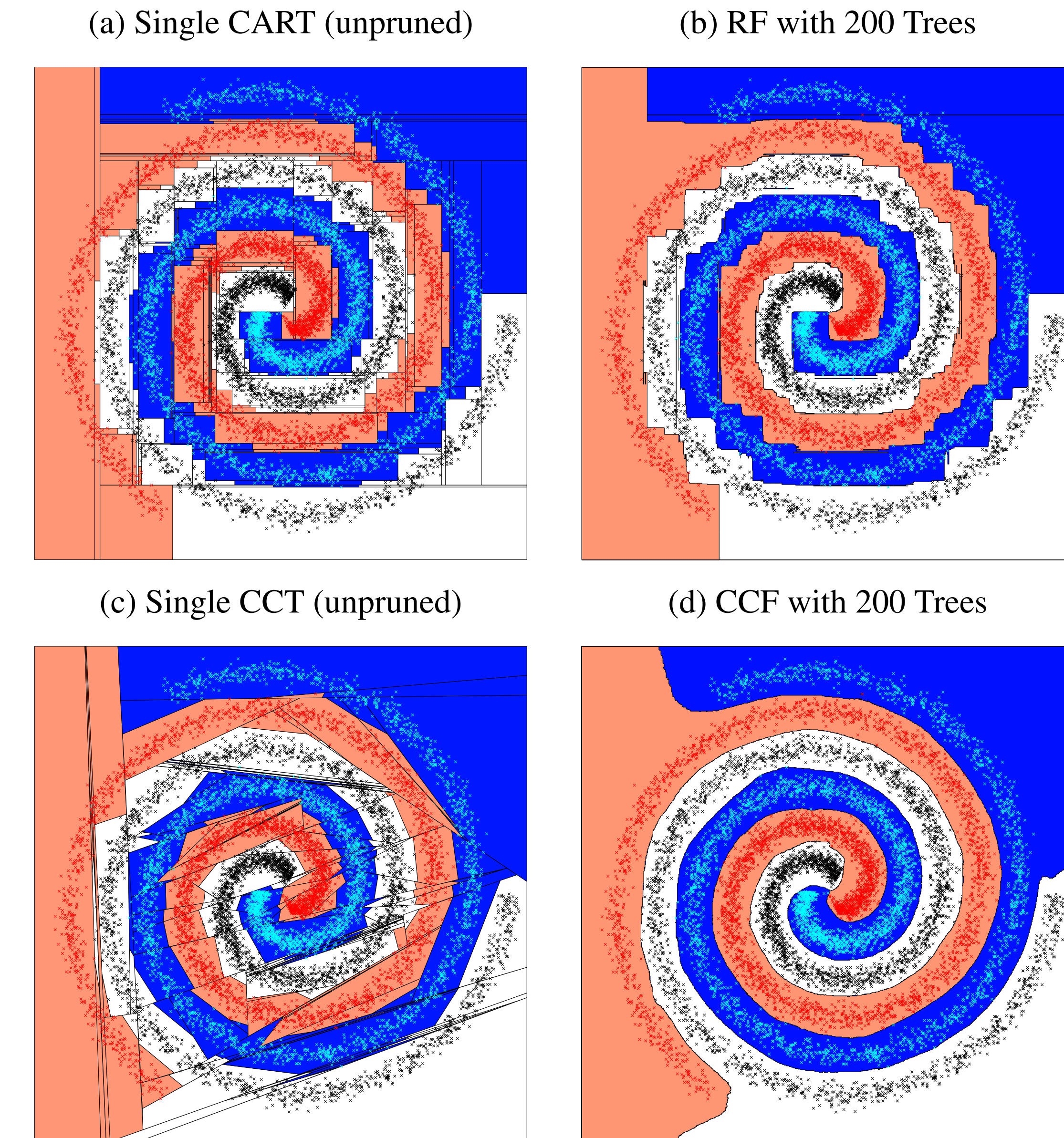
Growing CCTs

- Randomly subsample features [6]
- On local bootstrap sample, carry out CCA between features and class labels
- Calculate best split in projected space using original data (projection bootstrapping)
- If no improvement is possible assign as leaf

Algorithm 2 GROWTREE

```
1: Inputs  $X^j \in \mathbb{R}^{N^j \times D^j}$ ,  $\mathcal{Y}^j \in \mathbb{I}^{N^j \times K}$ ,  $\mathcal{D}^j \subseteq \{1, \dots, D^r + D^c\}$ ,  $\lambda \in \mathbb{Z}^+$ ,  $b \in \{\text{true}, \text{false}\}$ 
2: Set current node index  $j$  to an unique node identifier (0 for root node)
3: Sample  $\delta \subseteq \mathcal{D}^j$  by taking  $\min(\lambda, |\mathcal{D}^j|)$  samples without replacement from  $\mathcal{D}^j$ 
4: While  $\delta$  contains features without variation, eliminate these from  $\mathcal{D}^j$  and  $\delta$  and resample
5:  $\gamma = \delta$  mapped to the column indices of  $X^j$  in accordance with the 1-of-K encoding of  $X^c$ 
6: if  $b$  then  $\{X', \mathcal{Y}'\} \leftarrow$  sample with replacement  $N^j$  rows from  $\{X_{(\gamma, \gamma)}, \mathcal{Y}^j\}$ 
7: else  $\{X', \mathcal{Y}'\} \leftarrow \{X, \mathcal{Y}\}$  end if
8: if all rows in  $X'$  or  $\mathcal{Y}'$  are identical then
9:   if all rows in  $X_{(\gamma, \gamma)}$  or  $\mathcal{Y}^j$  are identical then return  $[j, \emptyset, \text{LABEL}(\mathcal{Y}^j)]$  end if
10:   $\{X', \mathcal{Y}'\} \leftarrow \{X'_{(\gamma, \gamma)}, \mathcal{Y}^j\}$ 
11: end if
12: if  $X'$  contains only two unique rows then
13:   $\mathcal{X}' = \text{UNIQUEROWS}(X')$ 
14:   $\phi_{(\gamma, \gamma)} \leftarrow \mathcal{X}'_{(2, \gamma)} - \mathcal{X}'_{(1, \gamma)}$ ,  $\phi_{\gamma(\gamma, \gamma)} \leftarrow 0$ 
15: else
16:   $[A, \cdot] = \text{CCA}(X', \mathcal{Y}')$ 
17:   $R_{(\gamma, \gamma)} \leftarrow A$ ,  $R_{(\gamma, \gamma)} \leftarrow 0$ 
18:   $U = X^j R$ 
19:   $[\cdot, s_j, \text{gain}] = \text{FINDBESTSPLIT}(U)$ 
20:  if  $\text{gain} \leq 0$  then return  $[j, \emptyset, \text{LABEL}(\mathcal{Y}^j)]$  end if
21:   $\phi_j = R_{(\gamma, \delta)}$ 
22: end if
23:  $\tau_j = \{n \in \{1, \dots, N^j\} : X'_{(n, \gamma)} \phi_j \leq s_j\}$ ,  $\tau_r = \{1, \dots, N^j\} \setminus \tau_j$ 
24:  $[\chi_{(j, 1)}, \Psi_i, \Theta_i] = \text{GROWTREE}(X'_{(\tau_j, \gamma)}, \mathcal{Y}'_{(\tau_j, \gamma)}, \mathcal{D}^j, \lambda, b)$ 
25:  $[\chi_{(j, 2)}, \Psi_r, \Theta_r] = \text{GROWTREE}(X'_{(\tau_r, \gamma)}, \mathcal{Y}'_{(\tau_r, \gamma)}, \mathcal{D}^j, \lambda, b)$ 
26:  $\psi_j = \{\chi_{(j, 1)}, \chi_{(j, 2)}, \phi_j, s_j\}$ 
27: return  $[j, \{\psi_j \cup \Psi_i \cup \Psi_r\}, \{\Theta_i \cup \Theta_r\}]$ 
```

CCFs have Smoother Decision Surfaces



Experiments

- 37 datasets [7]
- Compare against random forest, rotation forest and CCF-Bag which uses bagging instead of the projection bootstrap

	CCF	CCF-Bag	RF	Rotation Forest
CCF	-	2	2	7
CCF-Bag	18	-	1	14
RF	26	27	-	25
Rotation Forest	13	10	4	-

Number of significant victories column vs row at 1% level of Wilcoxon signed ranked test

Conclusions

- CCFs outperforms all other methods significantly
- Computationally less expensive than rotation forest and similar to RF
 - RF poor at dealing with such correlations
 - Rotation forests can only adapt to global correlations
- New benchmark for out-of-box tree ensemble classification
- Concepts introduced apply to forest regression models

Results

Data set	K	N	D^r	D^c	CCF	CCF-Bag	RF	Rotation Forest
Balance scale	3	625	0	4	91.06 \pm 3.74	91.26 \pm 3.48	83.81 \pm 4.23	92.71 \pm 3.43
Banknote	2	1372	0	4	100.00 \pm 0.00	100.00 \pm 0.00	99.28 \pm 0.75	100.00 \pm 0.00
Breast tissue	6	106	0	9	71.58 \pm 11.79	71.09 \pm 12.38	69.03 \pm 12.99	71.58 \pm 12.70
Climate crashes	2	360	0	18	94.17 \pm 4.04	93.56 \pm 4.09	92.87 \pm 4.18	94.04 \pm 3.89
Fertility	2	100	0	9	86.73 \pm 9.38	86.47 \pm 9.70	86.40 \pm 9.78	87.67 \pm 8.78
Heart-SPECT	2	267	0	22	82.84 \pm 7.13	81.98 \pm 6.71	81.38 \pm 7.27	82.49 \pm 7.29
Heart-SPECTF	2	267	0	44	81.46 \pm 7.34	81.98 \pm 6.89	81.01 \pm 6.71	81.23 \pm 7.29
Hill valley	2	1212	0	100	100.00 \pm 0.00	100.00 \pm 0.00	61.02 \pm 4.34	93.74 \pm 2.66
Hill valley noisy	2	1212	0	100	94.99 \pm 1.85	94.38 \pm 2.02	57.98 \pm 4.41	88.75 \pm 2.84
ILPD	2	640	0	10	71.97 \pm 5.05	72.03 \pm 5.21	70.34 \pm 4.95	70.98 \pm 5.20
Ionosphere	2	351	0	33	95.12 \pm 3.63	94.38 \pm 3.78	93.56 \pm 3.89	94.30 \pm 3.51
Iris	3	150	0	4	97.56 \pm 3.89	97.69 \pm 3.78	94.93 \pm 5.39	95.82 \pm 5.10
Landsat satellite	2	6435	0	36	91.76 \pm 1.08	91.30 \pm 1.05	91.84 \pm 1.01	92.19 \pm 0.99
Letter	26	20000	0	16	97.75 \pm 0.33	97.42 \pm 0.36	96.64 \pm 0.38	97.52 \pm 0.32
Libras	15	360	0	90	89.70 \pm 4.75	88.65 \pm 5.23	81.30 \pm 5.97	90.28 \pm 4.79
MAGIC	2	19020	0	10	88.42 \pm 0.72	88.29 \pm 0.72	88.15 \pm 0.74	87.36 \pm 0.72
Nursery	5	12960	0	8	99.96 \pm 0.07	99.91 \pm 0.11	99.67 \pm 0.19	99.96 \pm 0.06
ORL	40	400	0	10304	97.82 \pm 2.25	97.38 \pm 2.75	97.55 \pm 2.52	NaN \pm NaN
Optical digits	10	5620	0	64	98.71 \pm 0.45	98.56 \pm 0.46	98.35 \pm 0.49	98.69 \pm 0.40
Parkinsons	2	195	0	22	93.90 \pm 5.43	92.27 \pm 6.01	90.97 \pm 6.04	92.39 \pm 5.44
Pen digits	10	10992	0	16	99.60 \pm 0.19	99.54 \pm 0.21	99.17 \pm 0.29	99.51 \pm 0.23
Polya	2	9255	0	169	78.82 \pm 1.31	78.72 \pm 1.28	78.72 \pm 1.36	79.79 \pm 1.32
Seeds	3	210	0	7	95.21 \pm 4.73	94.57 \pm 5.18	93.62 \pm 5.23	95.14 \pm 4.53
Skin seg	2	245057	0	3	99.97 \pm 0.01	99.97 \pm 0.01	99.96 \pm 0.01	99.96 \pm 0.01
Soybean	19	683	13	22	94.58 \pm 2.94	94.18 \pm 3.14	94.44 \pm 3.08	94.40 \pm 2.92
Spirals	3	10000	0	2	99.73 \pm 0.16	99.73 \pm 0.16	98.78 \pm 0.33	98.98 \pm 0.34
Splice	3	3190	60	0	96.90 \pm 0.97	96.71 \pm 1.12	96.88 \pm 0.93	95.74 \pm 1.18
Vehicle	4	846	0	18	82.68 \pm 3.93	82.68 \pm 4.08	74.74 \pm 4.64	79.09 \pm 4.39
Vowel-c	11	990	2	10	99.06 \pm 0.95	98.72 \pm 1.08	97.35 \pm 1.72	99.00 \pm 0.95
Vowel-n	11	990	0	10	98.01 \pm 1.32	97.18 \pm 1.64	96.75 \pm 1.77	98.52 \pm 1.23
Waveform (1)	3	5000	0	21	86.42 \pm 1.58	86.52 \pm 1.49	85.04 \pm 1.63	86.44 \pm 1.56
Waveform (2)	3	5000	0	40	86.64 \pm 1.59	86.69 \pm 1.66	85.31 \pm 1.66	86.64 \pm 1.64
Wholesale-c	2	440	1	7	91.48 \pm 3.80	91.61 \pm 3.99	91.61 \pm 3.99	91.44 \pm 4.12
Wholesale-r	3	440	0	7	69.44 \pm 6.18	71.03 \pm 6.27	71.05 \pm 6.19	71.82 \pm 6.19
Wisconsin cancer	2	699	0	9	96.71 \pm 2.10	96.79 \pm 1.96	96.87 \pm 1.98	97.19 \pm 1.82
Yeast	2	1484	0	8	61.85 \pm 4.04	62.72 \pm 3.80	62.28 \pm 4.07	62.75 \pm 4.10
Zoo	7	101	0	16	96.73 \pm 5.73	96.27 \pm 6.19	95.20 \pm 6.42	94.33 \pm 6.70

Above: % test accuracy 15, 10-fold cross validations, best method in bold, \bullet / \circ indicate CCF significantly better / worse at 1% level

Below: corresponding box plots

