

Leaf Classification from Local Boundary Analysis

Anne Jorstad, Applied Math and Scientific Computation, UMD
Dr. David Jacobs, Computer Science, UMD



The Electronic Field Guide for Plants¹ in action.

Background

A previous leaf classification system uses the Inner-Distance Shape Context (IDSC) to compare the global shapes of leaves. The algorithm works very well, except when the global shape of two leaves are very similar. I have developed a wavelet-based algorithm which uses local boundary features to improve classification for these cases.



(a) *Cephalanthus occidentalis*
(smooth boundary)



(b) *Carpinus caroliniana*
(serrated boundary)

Globally similar leaves with distinct local features.

All data is from a database of 7481 leaves native to the Washington DC/Baltimore area.

Acknowledgements

I. Gaurav Agarwal, Haibin Ling, David Jacobs, Sameer Shirdhonkar, W. John Kress, Rusty Russell, Peter Belhumeur, Nandan Dixit, Steve Feiner, Dhruv Mahajan, Kalyan Sunkavalli, Ravi Ramamoorthi, Sean White. "First Steps Toward an Electronic Field Guide for Plants". *Taxon*, vol. 55, no. 3, Aug. 2006.

Wavelet Representation

The discrete wavelet transform decomposes a curve of length n into two curves of length $n/2$:

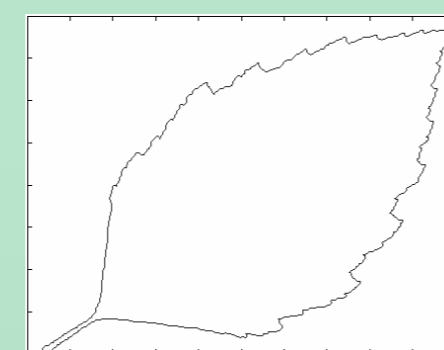
- ϕ : approximation coefficients generate the closest representation to the curve using half as many points
- ψ : detail coefficients store the local information required to regenerate the original curve from the approximation coefficients

$$\begin{aligned} f(t) &= \sum_{n=-\infty}^{\infty} [d_n \psi_{1n}(t) + c_n \phi_{1n}(t)] \\ &= \sum_{n=-\infty}^{\infty} [d_n \psi_{1n}(t) + c_n \sum_{m=-\infty}^{\infty} [d_m \psi_{2m}(t) + c_m \phi_{2m}(t)]] \\ &= \dots \end{aligned}$$

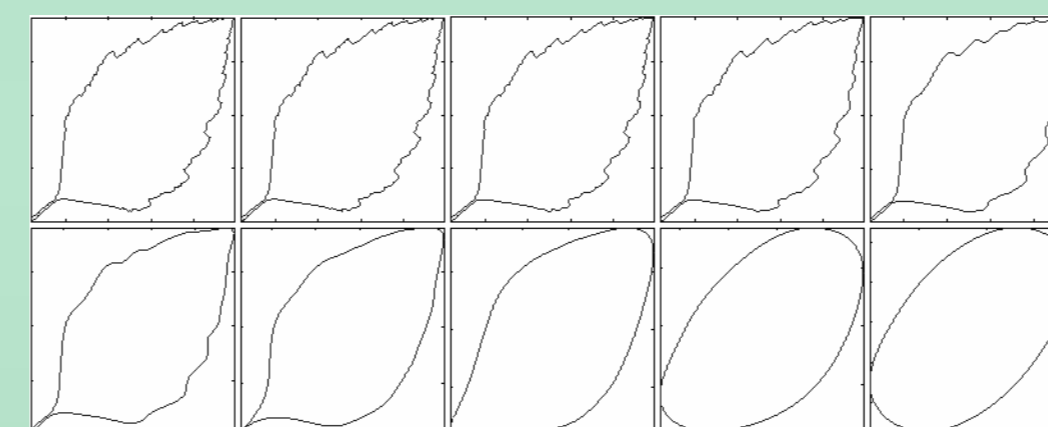
where $\psi_{in}(t) := n^{\text{th}}$ point of i^{th} scale

Apply the wavelet transform repeatedly to approximation coefficients:

- Generates detail coefficients over many scales (empirically determined that 3 was optimal)
- Reduces the approximation curves to an increasingly oval-like form



Left: Original boundary curve
Right: First 10 scales of wavelet approximations

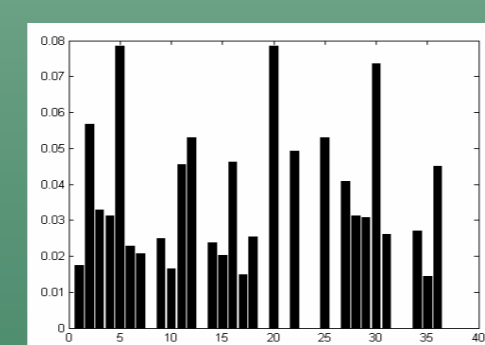
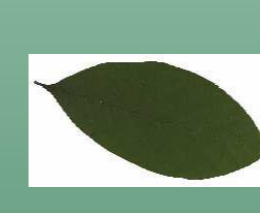
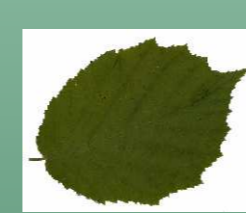


Each boundary point: $(x, y) \rightarrow [x_{d1}, y_{d1}, x_{d2}, y_{d2}, 0, y_{d3}]$

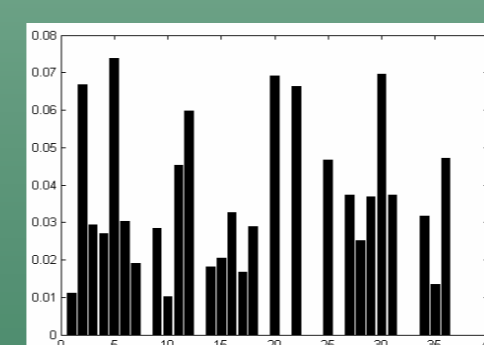
- Detail coefficients over 3 scales
- Rotated to the x-axis at the coarsest scale to preserve rotation invariance of original image

Clustering

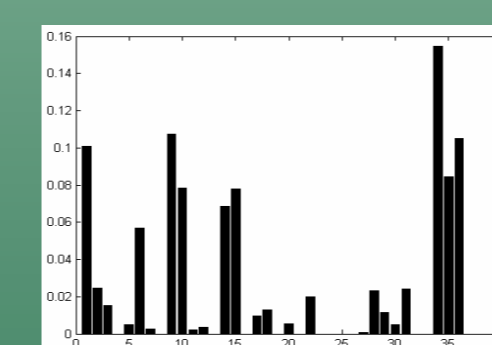
- *K-Means clustering*: find 36 cluster centers from all boundary points from all training leaves
- Individual leaf \rightarrow normalized distribution of its boundary points across cluster centers



(a)



(b)



(c)

Leaf image and corresponding histogram for (a) *Corylus americana*, (b) *Corylus americana*, different example, (c) *Asimina triloba*

Naïve Bayes Classification

Chi-squared distance between any pair of leaf distributions:

$$d(\ell_1, \ell_2) = \sum_{n=1}^{36} [\ell_1(n) - \ell_2(n)]^2$$

The distance from a leaf to a species is the shortest distance to any leaf known to be of that species:

$$d(\ell_{new}, S_k) = \min_{\ell_k} d(\ell_{new}, \ell_k | \ell_k \in S_k)$$

Define the probability of every combination of test data, given

- dW : distance using wavelet model
- dI : distance using IDSC model

Use Bayes' Rule to predict the true species of any unclassified leaf.

$$\text{Bayes' Rule: } P[A|B] = \frac{P[B|A] \cdot P[A]}{P[B]}$$

$$\begin{aligned} \text{Species}(\ell) &= \underset{S_k}{\text{argmax}} P[\ell \in S_k | dW(\ell, S_k), dI(\ell, S_k)] \\ &= \underset{S_k}{\text{argmax}} \frac{P[dW(\ell, S_k), dI(\ell, S_k) | \ell \in S_k] \cdot P[\ell \in S_k]}{P[dW(\ell, S_k), dI(\ell, S_k)]} \end{aligned}$$

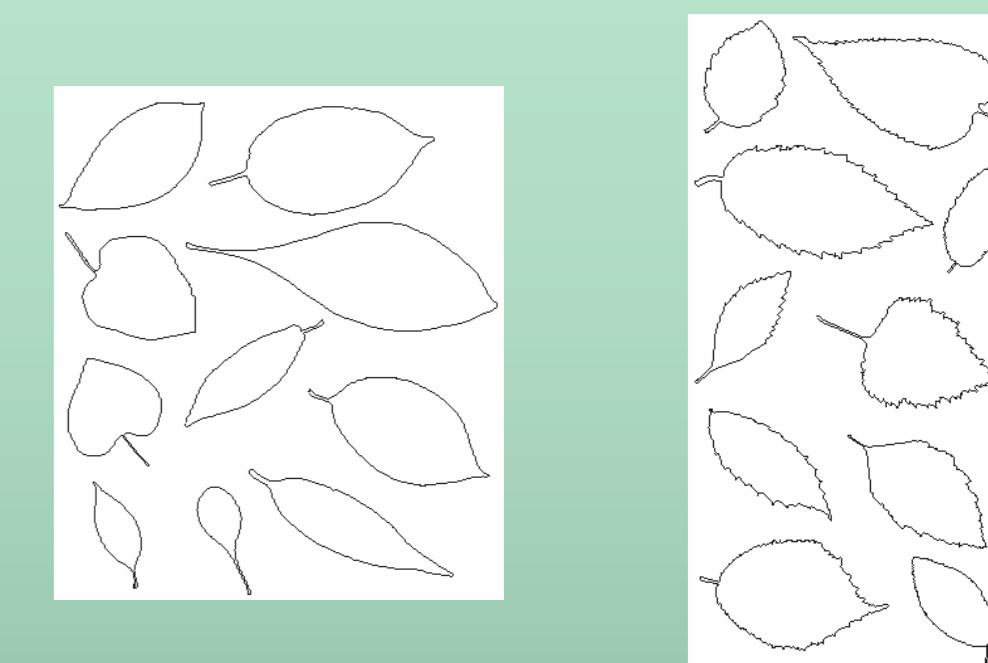
Preliminary Results

Training data: 200 leaves

- 10 smooth species
- 10 serrated species
- 10 examples of each

Testing data: 100 leaves

- 5 new examples of each of the trained species



(a) Smooth species (b) Serrated species

	IDSC alone	Wavelet alone	Wavelet + IDSC
Identified correct species	62%	46%	71%
Identified incorrect species with correct serration	53%	100%	100%

Comments:

- IDSC predicts serration no better than chance if it fails to identify the species.
- Wavelets captures the missing local information, but cannot make useful classification decisions on its own.
- Combining local and global information succeeds in producing a better overall species prediction!