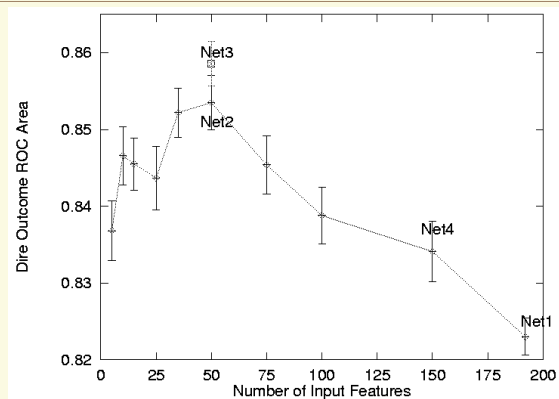


Special Topic: Feature Selection

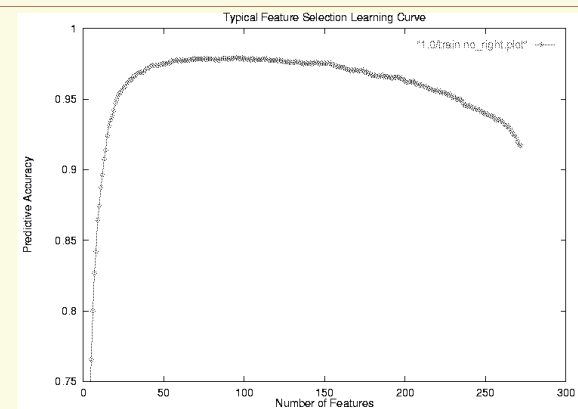
Anti-Motivation

- Most learning methods implicitly do feature selection:
 - decision trees: use info gain or gain ratio to decide what attributes to use as tests. many features don't get used.
 - neural nets: backprop learns strong connections to some inputs, and near-zero connections to other inputs.
 - kNN, MBL: weights in Weighted Euclidean Distance determine how important each feature is. weights near zero mean feature is not used.
 - SVMs: maximum margin hyperplane may focus on important features, ignore irrelevant features.
- So why do we need feature selection?

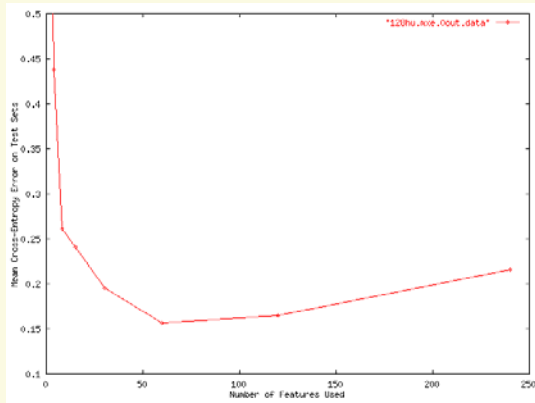
Motivation



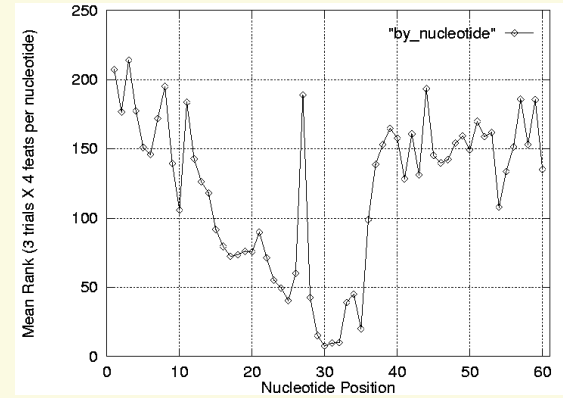
Motivation



Motivation



Motivation



Motivation: Ornithology



23,948 observations
- 12,849 presence
- 11,099 absence

What is important in this domain?

yearseason	nearby_feeders	snow_cov_atleast	ned_age_m	gcrb2311	gcnosw2910
dayofweek	nearby_feeders	snow_cov_atleast	ned_age_f	gcrb2312	gcnosw2911
halfdays	cat	snow_cruty	households_sqm	gcrb2313	gcnosw2912
latitude	dogs_precip_leatleast	ave_hh_sq	hshld_1_m_hh	gcrb2501	gcnosw4
longitude	humans_effort_hrs_atleast	hshld_1_m_hh	hshld_1_f_hh	gcrb2502	gcnosw513
yard_type_garden	count_area_size	nlod_water	marbh_cnd_hh	gcrb2503	gcnosw301
yard_type_landsc	housing_density	nlod_developed	marbh_cnd_hh	gcrb2504	gcnosw302
yard_type_woods	high_feeders	nlod_barren	marbh_no_c_hh	gcrb2505	gcnosw303
yard_type_desert	fed_yr_round	nlod_forested	whh_chlid_hh	gcrb2506	gcnosw304
yard_type_pavement	fed_in_jan	nlod_shrubland	whh_chlid_hh	gcrb2507	gcnosw311
hab_dcid_woods	fed_in_feb	nlod_orchards	families	gcrb2508	gcnosw312
hab_evgr_woods	fed_in_mar	nlod_grasslands	ave_fam_sq	gcrb2509	gcnosw313
hab_mixed_woods	fed_in_apr	nlod_cultivated	vacant_hu	gcrb2510	gcnosw301
hab_orchard	fed_in_may	nlod_wetlands	owner_occ_hu	gcrb2511	gcnosw302
hab_park	fed_in_jun	nlod_gro_10	renter_occ_hu	gcrb2512	gcnosw303
hab_water_fresh	fed_in_jul	pop00_sqm_av	housden	gcrb2513	gcnosw304
hab_water_salt	fed_in_aug	white_pop00	nlod_red	gcrb2514	gcnosw305
hab_residential	fed_in_sep	black_pop00	nlod_black	gcrb2515	gcnosw306
hab_industrial	fed_in_oct	amer1_sq_pop00	nlod_black	gcrb2516	gcnosw307
hab_agricultural	fed_in_nov	hawn_pi_pop00	nlod_black	gcrb2517	gcnosw308
hab_desert_scrub	fed_in_dec	other_pop00	nlod_black	gcrb2518	gcnosw309
hab_young_woods	munfeeders_suat	nlod_black	nlod_black	gcrb2519	gcnosw310
hab_swamp	munfeeders_ground	nlod_black	nlod_black	gcrb2520	gcnosw311
hab_marsh	munfeeders_hanging	nlod_black	nlod_black	gcrb2521	gcnosw312
hab_other	munfeeders_platfrm	nlod_black	nlod_black	gcrb2522	gcnosw313
evgr_trees_atleast	munfeeders_humming	nlod_black	nlod_black	gcrb2523	gcnosw314
evgr_shrubs_atleast	munfeeders_water	nlod_black	nlod_black	gcrb2524	gcnosw315
dcid_trees_atleast	munfeeders_thistle	nlod_black	nlod_black	gcrb2525	gcnosw316
dcid_shrubs_atleast	munfeeders_fruit	nlod_black	nlod_black	gcrb2526	gcnosw317
fru_trees_atleast	day1_am	nlod_black	nlod_black	gcrb2527	gcnosw318
cacti_atleast	day1_pm	nlod_black	nlod_black	gcrb2528	gcnosw319
		nlod_black	nlod_black	gcrb2529	gcnosw320

After feature selection ...



Yearseason,	nlcd_barren,
Dayselapsed,	nlcd_shrubland,
Longitude,	nlcd_grasslands,
hab_young_woods,	pop00_sqmi_mv,
hab_swamp,	asian_pop00,
nearby_feeders,	Gcprec0101,
Cats,	Gcrh2313,
Dogs,	Gcskyc5002,
fed_in_feb,	Gcwnd60b01,
Numfeeders,	Gcslvp6101,
Hanging,	Gcsnow1411,
day1_am,	Gcsnow3513,
snow_dep_atleast,	gcsun5312
effort_hrs_atleast,	

Brute-Force Approach

- Try all possible combinations of features
- Given N features, 2^N subsets of features
 - usually too many to try
 - danger of overfitting
- Train on train set, evaluate on test set (or use cross-validation)
- Use set of features that performs best on test set(s)

Two Basic Approaches

- Wrapper Methods:
 - give different sets of features to the learning algorithm and see which works better
 - *algorithm dependent*
- Proxy Methods (relevance determination methods)
 - determine what features are important or not important for the prediction problem without knowing/using what learning algorithm will be employed
 - *algorithm independent*

Wrapper Methods

- Wrapper methods find features that work best with some particular learning algorithm:
 - best features for kNN and neural nets may not be best features for decision trees
 - can eliminate features learning algorithm “has trouble with”
- Forward stepwise selection
- Backwards elimination
- Bi-directional stepwise selection and elimination

Relevance Determination Methods

- Rank features by information gain
 - Info Gain = reduction in entropy due to attribute

$$Entropy = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

- Try first 10, 20, 30, ..., N features with learner
- Evaluate on test set (or use cross validation)
- May be only practical method if thousands of attributes

Ornithology Example

0.9948	houfin
0.6259	+longitude
0.6863	+dayselapsd
0.7369	+yearseason
0.7506	+numfeeders_hanging
0.7594	+asian_pop00
0.7619	+nearby_feeders
0.7642	+hab_young_woods
0.7656	+pop00_sqmi_mv
0.7682	+effort_hrs_atleast
0.7691	+gcsnow3513
0.7699	+gcsnow1411
0.7703	+gcsun5312
0.7716	+hab_swamp
0.7722	+nlcd_barren
0.1083	latitude
0.0896	gctemp0302
0.0846	gctemp0313
0.0841	gcsnow34
0.0840	gctemp0312
0.0822	gcvis68a01
0.0788	gcsnow1401
0.0736	gctmp07a13
0.0714	gctmp07b13
0.0701	gcfog4403
0.0679	gctemp0301
0.0678	gcfog4401
0.0672	gcrh2303
0.0661	gcsnow1413

Advantages of Feature Selection

- Improved accuracy!
- Less complex models:
 - run faster
 - easier to understand, verify, explain
- Feature selection points you to most important features
- Don't need to collect/process features not used in models

Limitations of Feature Selection

- Given many features, feature selection can overfit
 - consider 10 relevant features, and 10^9 random irrelevant features
- Wrapper methods require running base learning algorithm many times, which can be expensive!
- Just because feature selection doesn't select a feature, doesn't mean that feature isn't a strong predictor
 - redundant features
- May throw away features domain experts want in model
- Most feature selection methods are greedy and won't find optimal feature set

Current Research in Feature Selection

- Speeding-up feature selection (1000's of features)
- Preventing overfitting (1000's of features)
- Better proxy methods
 - would be nice to know what the good/relevant features are independent of the learning algorithm
- Irrelevance detection:
 - truly irrelevant attributes can be ignored
 - better algorithms
 - better definition(s)

Bottom Line

- Feature selection almost always improves accuracy on real problems
- Plus:
 - simpler, more intelligible models
 - features selected can tell you about problem
 - less features to collect when using model in future

Feature selection usually is a win.

Special Topic:

Missing Values

Missing Values Common in Real Data

- Pneumonia:
 - 6.3% of attribute values are missing
 - one attribute is missing in 61% of cases
- C-Section:
 - only about 1/2% of attribute values are missing
 - but 27.9% of cases have at least 1 missing value
- Ornithology:
 - 8% of attribute values are missing
 - 94.5% of cases have at least 1 missing value
- UCI machine learning repository:
 - 31 of 68 data sets reported to have missing values

“Missing” Can Mean Many Things

- MAR: "Missing at Random":
 - usually best case
 - usually not true
- Non-randomly missing
- Presumed normal, so not measured
- Causally missing:
 - attribute value is missing because of other attribute values (or because of the outcome value!)

Dealing With Missing Data

- Some learning methods can handle missing values
- Throw away cases with missing values
 - in some data sets, most cases get thrown away
 - if not missing at random, throwing away cases can bias sample towards certain kinds of cases
- Treat “missing” as a new attribute value
 - what value should we use to code for missing with continuous or ordinal attributes?
 - if missing causally related to what is being predicted?
- Impute (fill-in) missing values
 - once filled in, data set is easy to use
 - if missing values poorly predicted, may hurt performance of subsequent uses of data set

Imputing Missing Values

- Fill-in with mean, median, or most common value
- Predict missing values using machine learning
- Expectation Maximization (EM):
 - Build model of data values (ignore missing vals)
 - Use model to estimate missing values
 - Build new model of data values (including estimated values from previous step)
 - Use new model to re-estimate missing values
 - Re-estimate model
 - Repeat until convergence

Potential Problems

- Imputed values may be inappropriate:
 - in medical databases, if missing values not imputed separately for male and female patients, may end up with male patients with 1.3 prior pregnancies, and female patients with low sperm counts
 - many of these situations will not be so humorous/obvious!
- If some attributes are difficult to predict, filled-in values may be random (or worse)
- Some of the best performing machine learning methods are impractical to use for filling in missing values (neural nets)
- Beware of coding - reliably detect missing cases can be difficult

Research in Handling Missing Values

- Lazy learning:
 - don't train a model until you know test case
 - missing in test case may "shadow" missing values in train set
- Better algorithms:
 - Expectation maximization (EM)
 - Non-parametric methods (since parametric methods often work poorly when assumptions are violated)
- Faster Algorithms:
 - apply to very large datasets