UVA CS 6316: Machine Learning

Lecture 7: Feature Selection

Dr. Yanjun Qi

University of Virginia

Department of Computer Science

Course Content Plan Six major sections of this course

□ Regression (supervised)
 □ Classification (supervised)
 □ Unsupervised models
 □ Learning theory

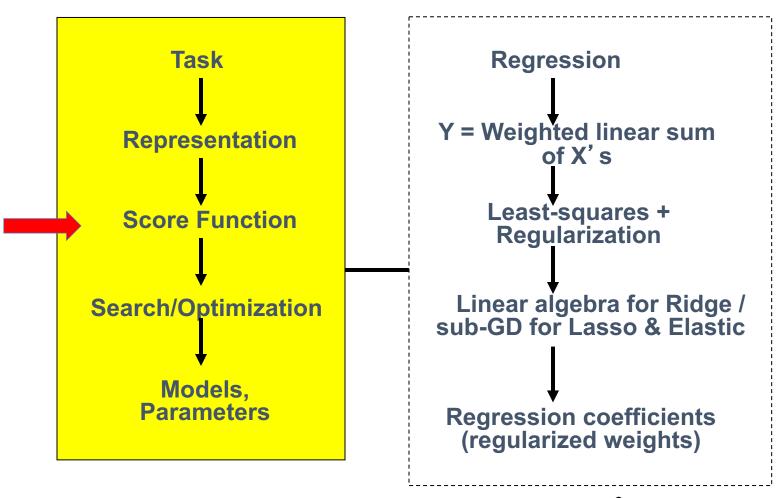
- ☐ Graphical models
- ☐ Reinforcement Learning

About interactions among X1,... Xp

About f()

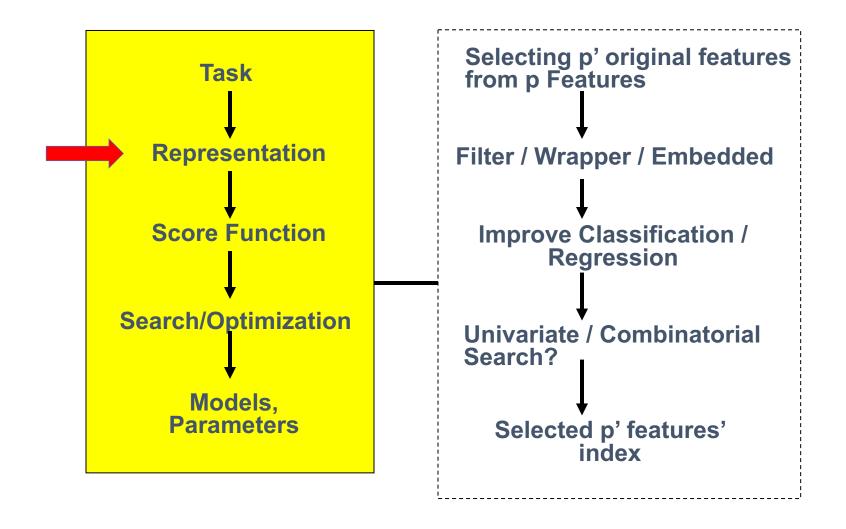
Learn program to Interact with its environment

Last: Regularized multivariate linear regression



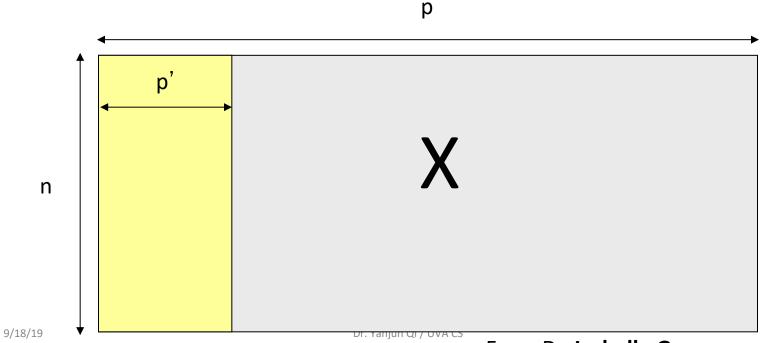
$$\min_{\text{Dr. Yanjun Qi / UVA CS}} \sum_{i=1}^{n} \left(Y - \hat{Y} \right)^2 + \lambda (\sum_{j=1}^{p} \beta_j^q)^{1/q}$$

Today: Feature Selection



Feature Selection

• Thousands to millions of low level features: select the most relevant ones to build better, faster, and easier to understand learning models.

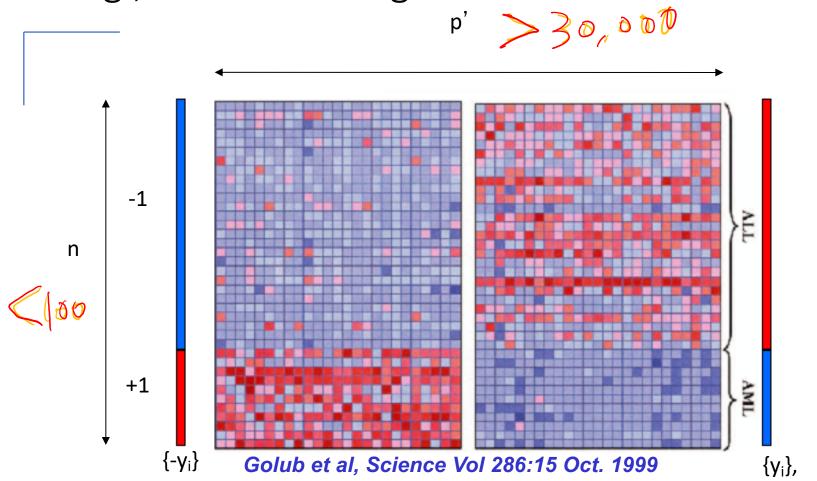


From Dr. Isabelle Guyon

e.g., Movie Reviews and Revenues: An Experiment in Text Regression, Proceedings of HLT '10 (1.7k n / >3k features)

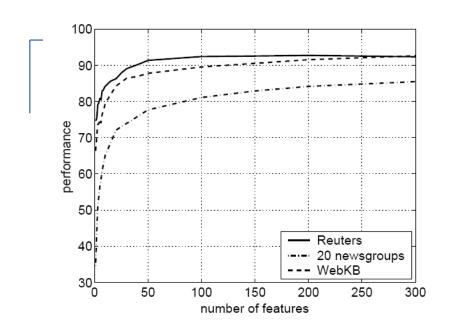
IV. Features e.g. counts of a ngram in			
1	Lexical n-grams (1,2,3)		
11	Part-of-speech n-grams (1,2,3)		
111	Dependency relations (nsubj,advmod,)		
Meta	U.S. origin, running time, budget (log), # of opening screens, genre, MPAA rating, holiday release (summer, Christmas, Memorial day,), star power (Oscar winners, high-grossing actors)		

e.g., Leukemia Diagnosis



7

e.g., Text Categorization with feature Filtering



Reuters: 21578 news wire, 114 semantic categories.

20 newsgroups: 19997 articles, 20 categories.

WebKB: 8282 web pages, 7 categories.

Bag-of-words: >100,000 features.

Top 3 words of some output Y categories:

- Alt.atheism: atheism, atheists, morality
- Comp.graphics: image, jpeg, graphics
- Sci.space: space, nasa, orbit
- Soc.religion.christian: god, church, sin
- Talk.politics.mideast: israel, armenian, turkish
- Talk.religion.misc: jesus, god, jehovah

Bekkerman et al, JMLR, 2003

We aim to make the learned model: Feature Selection → Simpler models

- 1. Generalize Well
 - Less sensitive to noise
 - Lower variance Occam's razor) --- More later!
- 2. Computationally Scalable and Efficient
 - Easier to train (to need less labeled examples)
 - Simpler to use (computationally)
- 3. Robust / Trustworthy / Interpretable
 - Especially for some domains, this is about trust!
 - Easier to explain (more interpretable!)

Occam's razor: law of parsimony

The principle of Occam's razor

states that the explanation of any phenomenon should make as few assumptions as possible, eliminating those that make no difference to any observable predictions of the theory

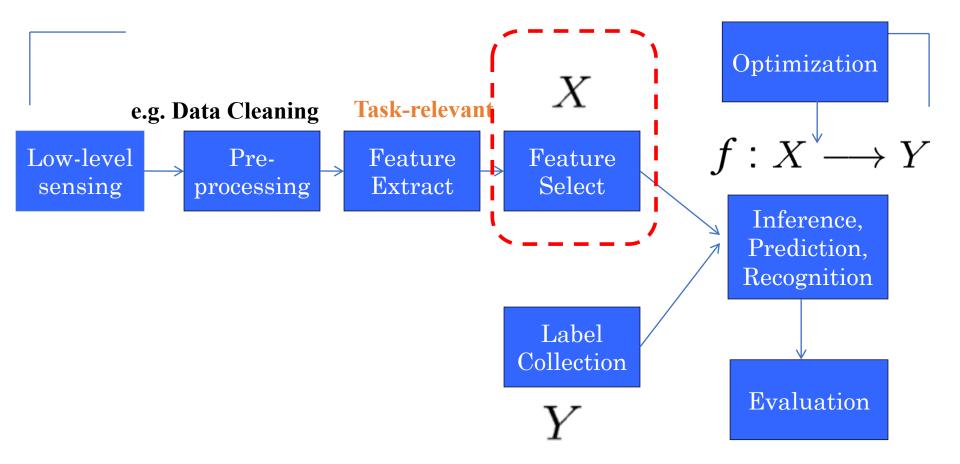
image at:

ww.butterflyeffect.ca/.../OccamsRaz
or.htmlRemove frame

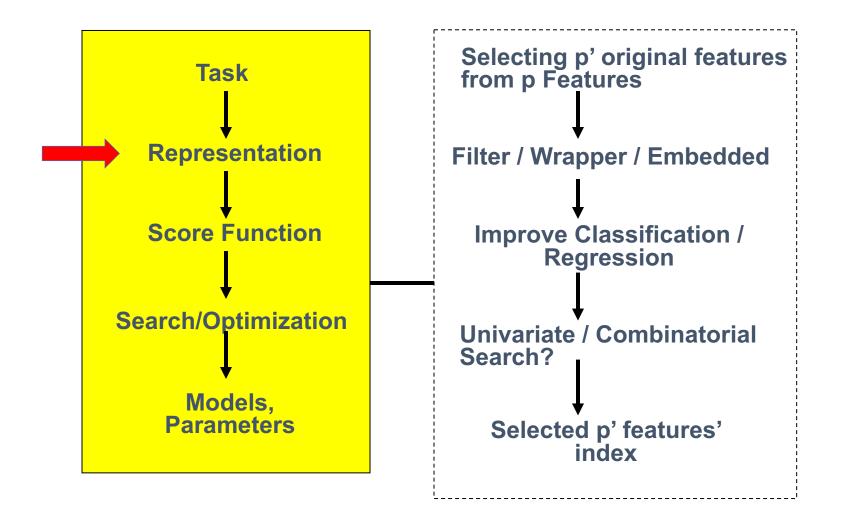


parsimony: extreme unwillingness to spend money or use resources.

A Typical Machine Learning Pipeline



Today: Feature Selection



Summary of Feature Selection Methods:

- Filtering approach:
 - ranks features or feature subsets independently of the predictor.
 - ...using univariate methods: consider one variable at a time
 - ...using multivariate methods: consider more than one variables at a time
- Wrapper approach:
 uses a predictor to assess features or feature subsets.
- Embedding approach:

uses a predictor to build a (single) model with a subset of features that are internally selected.

Nomenclature

- Univariate method: considers one variable (feature) at a time.
- Multivariate method: considers subsets of variables (features) together.
- Filter method: ranks features or feature subsets independently of the predictor.
- Wrapper method: uses a predictor to assess features or feature subsets.

Summary of Feature Selection Methods:

Filtering approach:

ranks features or feature subsets independently of the predictor.

- ...using univariate methods: consider one variable at a time
- ...using multivariate methods: consider more than one variables at a time
- Wrapper approach:
 uses a predictor to assess features or feature subsets.
- Embedding approach:
 uses a predictor to build a (single) model with a subset of
 features that are internally selected.

(I) Filtering: Univariate:

e.g., Pearson Correlation

Pearson correlation coefficient

$$r(x,y) = \frac{\sum_{i=1}^{n} (x_i - x)(y_i - y)}{\sqrt{\sum_{i=1}^{n} (x_i - x)^2 \times \sum_{i=1}^{n} (y_i - y)^2}}$$

- Measuring the linear correlation between two variables: x and y,
- giving a value between +1 and -1 inclusive, where 1 is total positive correlation, 0 is no correlation, and -1 is total negative correlation.

where
$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
 and $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$.

$$|r(x,y)| \leq 1$$

(I) Filtering: Univariate:

e.g., Pearson Correlation

Pearson correlation coefficient

$$r(x,y) = \frac{\sum_{i=1}^{n} (x_i - x)(y_i - y)}{\sqrt{\sum_{i=1}^{n} (x_i - x)^2 \times \sum_{i=1}^{n} (y_i - y)^2}}$$

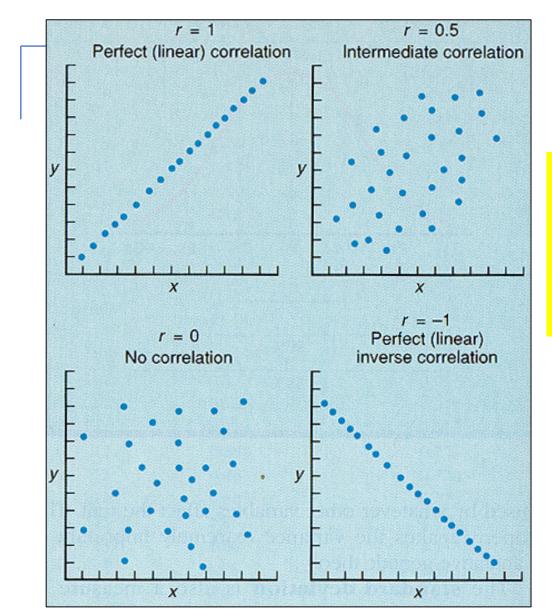
$$|r(x,y)| \leq 1$$

where
$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
 and $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$.

• Special case: cosine distance

$$s(x,y) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|}$$

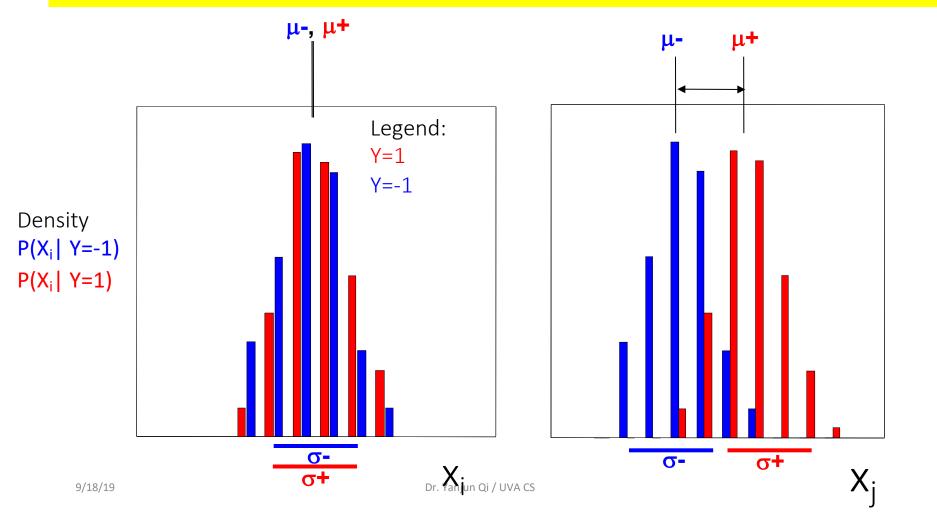
(I) Filtering: Univariate: e.g., Pearson Correlation



- can only detect linear dependencies between two variables
- (e.g. between one feature vs. target)

(I) Filtering: univariate filtering e.g. T-test

Goal: determine the relevance of a given single feature for two classes of samples.



(I) Filtering: univariate filtering e.g. T-test

significant?

T-test

 Assumption: Two Normally distributed classes with equal variance σ^2 unknown; estimated from data as σ^2_{within} .

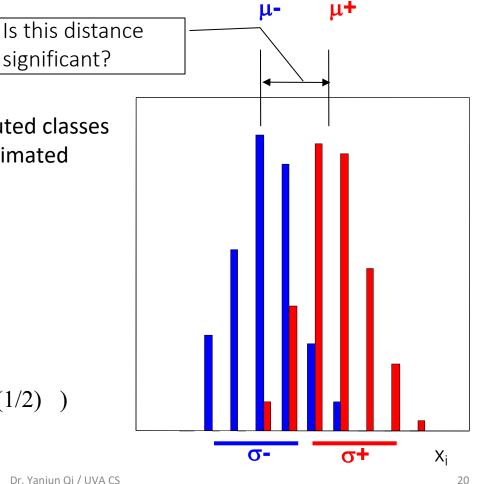
• Null hypothesis H_0 : μ + = μ -

T statistic:

If H₀ is true, then

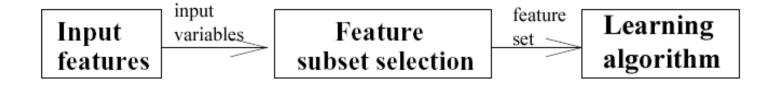
t=
$$(\mu + - \mu -)/(\sigma_{within}(1/|m^{+}| + 1/|m^{-}|) \perp (1/2)$$

Student(m⁺+m⁻-2 d.f.)



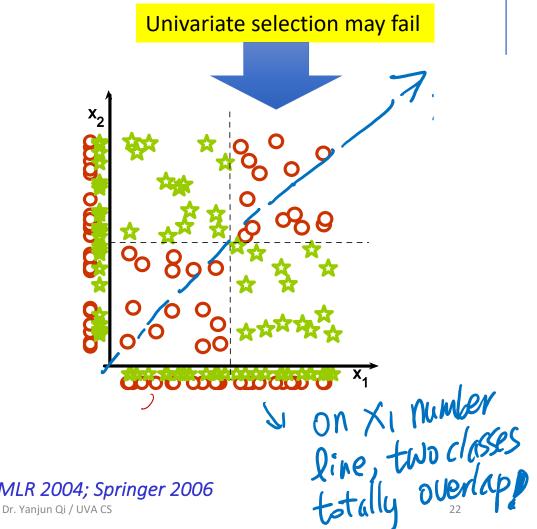
(I) Filtering: multi-variate: Feature Subset Selection

- Filter Methods
 - Select subsets of variables as a pre-processing step, independently of the used classifier!!



- E.g. Group correlation
- E.g. Information theoretic filtering methods such as Markov blanket

(I) Filtering: multi-variate: **Feature Subset Selection**



(I) Filtering: multi-variate: Feature Subset Selection

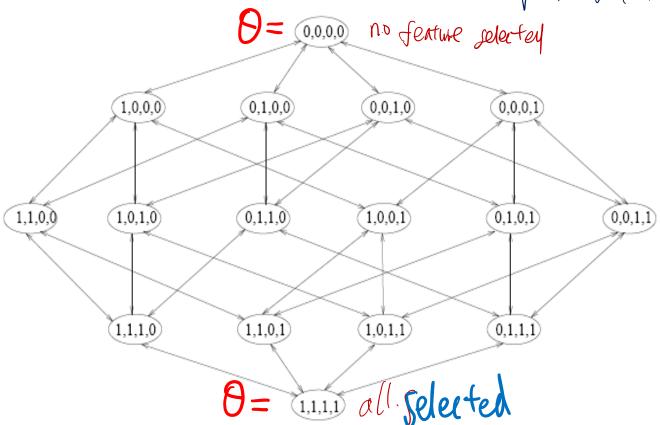
good, not, boring, many possible words
2 grain
Setthers 7 grains Very good, very very good, not very boring,

(I) Filtering: multi-variate:

Feature Subset Selection

- You need:
 - a measure for assessing the goodness of a feature subset (scoring function)
 - a strategy to search the space of possible feature subsets
- Finding a minimal optimal feature set for an arbitrary target is NP-hard
 - => Good heuristics are needed!

each feature subset can be described by 30 = [6/1, 0/1, 0/1, 0/1, 0/1] A = [6/1, 0/1, 0/1, 0/1, 0/1] A = [6/1, 0/1, 0/1, 0/1, 0/1] A = [6/1, 0/1, 0/1, 0/1, 0/1]



p features, 2^p possible feature subsets!

)

(I) Filtering: Summary

- usually fast
- provide generic selection of features, not tuned by given learner (universal, learner-agnostic)
- this is also often criticised (feature set not optimized for used learner)
- Often used as a pre-processing step for other methods

(I) Filtering: (many other choices)

Method	X	Y	Comments
Name	$ \mathrm{Formula} \mathrm{B} \mathrm{M} \mathrm{C} \mathrm{B} $	M C	
Bayesian accuracy Balanced accuracy		s s	Theoretically the golden standard, rescaled Bayesian relevance Eq. 3.2. Average of sensitivity and specificity; used for unbalanced dataset, same as AUC for binary targets.
Bi-normal separation F-measure Odds ratio	Eq. 3.5 + s + Eq. 3.7 + s + Eq. 3.6 + s +	s s	Used in information retrieval. Harmonic of recall and precision, popular in information retrieval. Popular in information retrieval.
Means separation T-statistics Pearson correlation Group correlation χ^2 Relief Separability Split Value	Eq. 3.13 + i + + Eq. 3.8 + s +	i + s s +	Based on two class means, related to Fisher's criterion. Based also on the means separation. Linear correlation, significance test Eq. 3.12, or a permutation test. Pearson's coefficient for subset of features. Results depend on the number of samples m. Family of methods, the formula is for a simplified version ReliefX, captures local correlations and feature interactions. Decision tree index.
Kolmogorov distance Bayesian measure Kullback-Leibler divergence Jeffreys-Matusita distance Value Difference Metric	Eq. 3.16 + s + + Eq. 3.20 + s + +	s + s + s +	Difference between joint and product probabilities. Same as Vajda entropy Eq. 3.23 and Gini Eq. 3.39. Equivalent to mutual information. Rarely used but worth trying. Used for symbolic data in similarity-based methods, and symbolic feature-feature correlations.
Mutual Information V Information Gain Ratio V Symmetrical Uncertainty J-measure Weight of evidence MDL 9/18/19	Eq. 3.32 + s + + Eq. 3.35 + s + + Eq. 3.36 + s + +	s + s + s +	Equivalent to information gain Eq. 3.30. Information gain divided by feature entropy, stable evaluation. Low bias for multivalued features. Measures information provided by a logical rule. So far rarely used. So far rarely used.

Summary of Feature Selection Methods:

- Filtering approach:
 - ranks features or feature subsets independently of the predictor.
 - ...using univariate methods: consider one variable at a time
 - ...using multivariate methods: consider more than one variables at a time
- Wrapper approach:
 uses a predictor to assess features or feature subsets.
- Embedding approach:
 uses a predictor to build a (single) model with a subset of
 features that are internally selected.

(2) Wrapper: Feature Subset Selection

- Learner is considered a black-box
- Interface of the black-box is used to score subsets of variables according to the predictive power of the learner when using the subsets.
- Results vary for different learners

(2) Wrapper: Feature Subset Selection

Two major questions to answer:

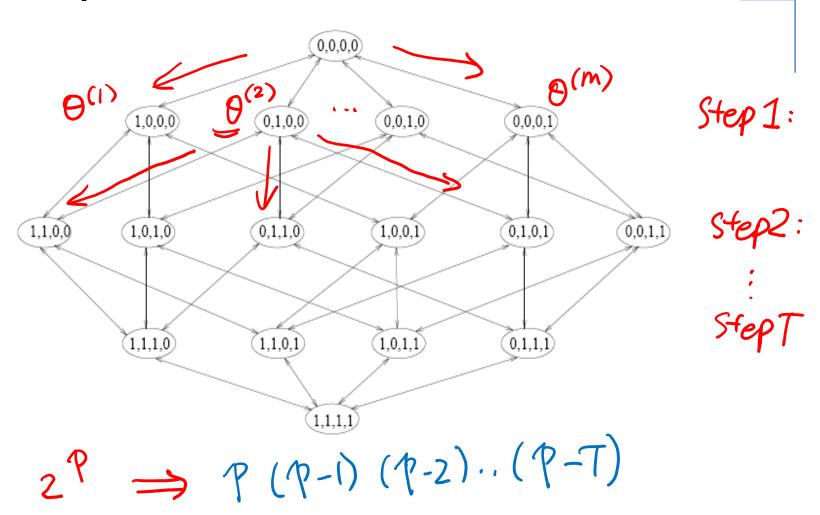
 (a). Assessment: How to measure performance of a learner that uses a particular feature subset?

 (b). Search: How to search in the space of all feature subsets?

(b). Search: How to search the space of all feature subsets?

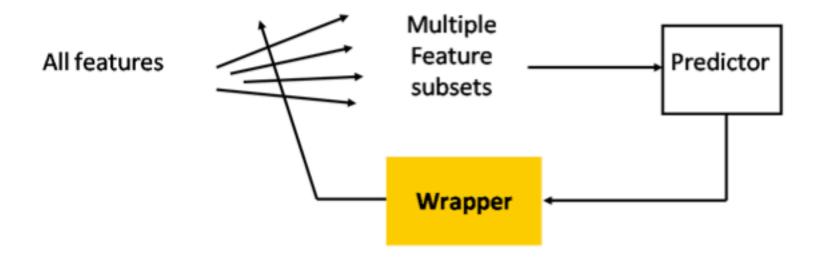
- The problem of finding the optimal subset is NP-hard!
- A wide range of heuristic search strategies can be used.
 Two different classes:
 - Forward selection (start with empty feature set and add features at each step)
 - Backward elimination (start with full feature set and discard features at each step)
- predictive power is usually measured on a validation set or by cross-validation
- By using the learner as a black box wrappers are universal and simple!
- Criticism: a large amount of computation is required.

(b). Search: How to search the space of all feature subsets?



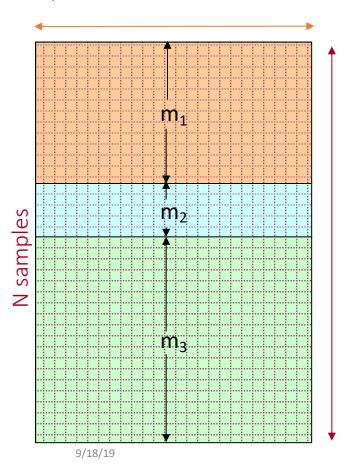
(a). Assessment: How to access multiple candidates of feature subsets

Wrapper Methods



(a). Assessment: feature subset assessment (for wrapper approach)

p variables/features



Split data into 3 sets:

training, validation, and test set.

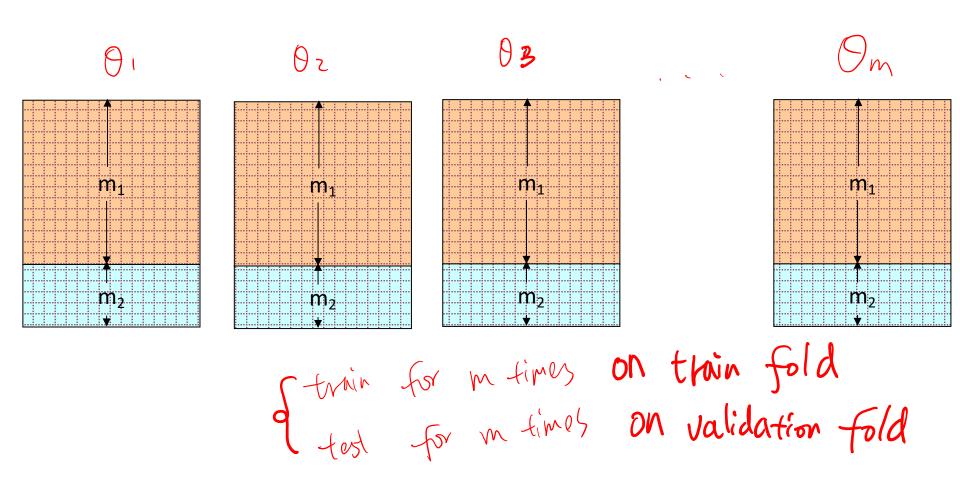
- 1) For each feature subset, train predictor on training data.
- 2) Select the feature subset, which performs best on validation data.
 - Repeat and average if you want to reduce variance (cross-validation).
- 3) Test on test data.

Danger of over-fitting with intensive search!

Dr. Yanjun Qi / UVA CS

From Dr. Isabelle Guyon

(a). Assessment: How to access multiple candidates of feature subsets



train data: argmin $J(\beta_{0}(i)) \Rightarrow \beta^{*}(i)$ $\beta_{0}(i)$ validation argmin Predict Loss (BX)

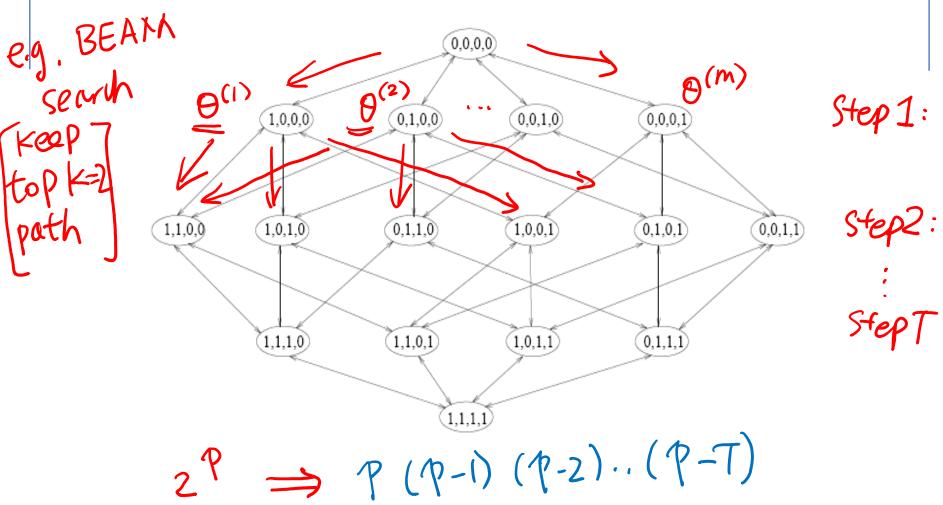
data {0(1), 0(2), ..., 0(m)} Predict Loss (Bx)

9/18/19

(b). Search: even more search strategies for selecting feature subset

- Forward selection or backward elimination.
- Beam search: keep k best path at each step.
- GSFS: generalized sequential forward selection when (n-k) features are left try all subsets of g features. More trainings at each step, but fewer steps.
- PTA(I,r): plus I, take away r at each step, run SFS I times then SBS r times.
- Floating search: One step of SFS (resp. SBS), then SBS (resp. SFS) as long as we find better subsets than those of the same size obtained so far.

(b). Search: How to search the space of all feature subsets?



9/18/19 Dr. Yanjun Qi / UVA CS 38

Summary of Feature Selection Methods:

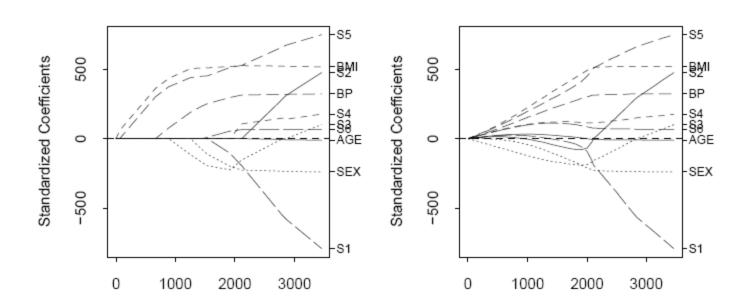
- Filtering approach:
 - ranks features or feature subsets independently of the predictor.
 - ...using univariate methods: consider one variable at a time
 - ...using multivariate methods: consider more than one variables at a time
- Wrapper approach:
 uses a predictor to assess features or feature subsets.
- Embedding approach:
 - uses a predictor to build a (single) model with a subset of features that are internally selected.

(3) Embedded: Feature Subset Selection

- Specific to a given learning machine!
- Performs variable selection (implicitly) in the process of training
- Just train a (single) model

(3) Embedded: e.g. Feature Selection via Embedded Methods: e.g., L₁-regularization

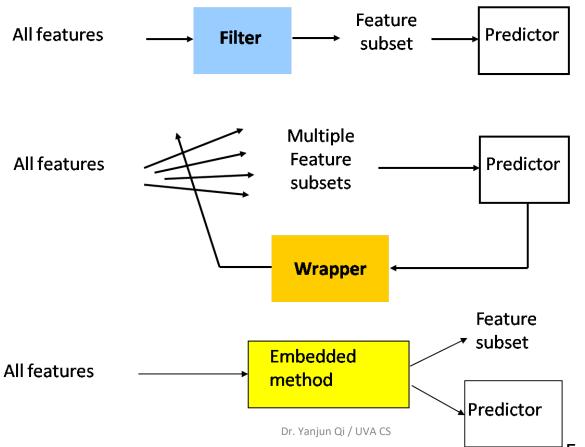
 l_1 penalty: $y \sim Model(X\beta) + \lambda \sum |\beta_i|$ (lasso) l_2 penalty: $y \sim Model(X\beta) + \lambda \sum \beta_i^2$ (ridge regression) LASSO Ridge Regression



9/18/19 Dr. Yanjun Qi / UVA CS 41

Summary: filters vs. wrappers vs. embedding

Main goal: rank subsets of useful features



42

In practice...

- No method is universally better:
 - wide variety of types of variables, data distributions, learning machines, and objectives.
- Feature selection is not always necessary to achieve good performance.

NIPS 2003 and WCCI 2006 challenges: http://clopinet.com/challenges

Later: Dimensionality Reduction,

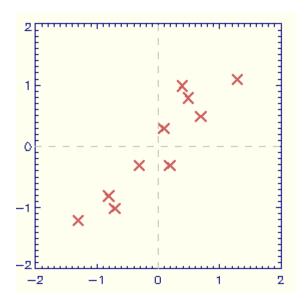
In the presence of many of features, select the most relevant subset of (weighted) combinations of features.

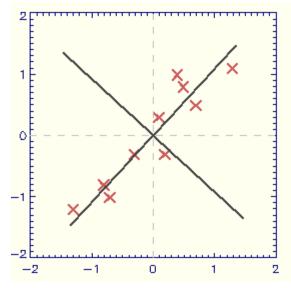
Feature Selection:
$$X_1, ..., X_p \rightarrow X_{k1}, ..., X_{kp'}$$

Dimensionality Reduction:
$$X_1, ..., X_m \rightarrow g_1(X_1, ..., X_m), ..., g_{p'}(X_1, ..., X_m)$$

Later: Dimensionality Reduction, e.g., (Linear) Principal Components Analysis

■ **PCA** finds a *linear* mapping of dataset X to a dataset X' of lower dimensionality. The variance of X that is remained in X' is maximal.





Dataset X is mapped to dataset X', here of the same dimensionality. The first dimension in X' (= the first principal component) is the direction of maximal variance. The second principal component is orthogonal to the first.

References

- ☐ Prof. Andrew Moore's slides
- ☐ Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.
- ☐ Dr. Isabelle Guyon's feature selection tutorials