

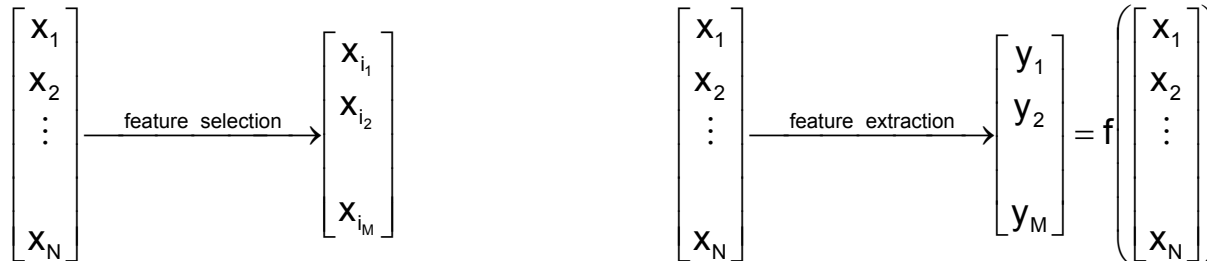
LECTURE 11: Sequential Feature Selection

- **Feature extraction vs. Feature selection**
- **Search strategy and objective functions**
- **Objective functions**
 - Filters
 - Wrappers
- **Sequential search strategies**
 - Sequential Forward Selection
 - Sequential Backward Selection
 - Plus-I Minus-r Selection
 - Bidirectional Search
 - Floating Search



Feature extraction vs. Feature selection

- As we discussed in Lecture 9, there are two general approaches for performing dimensionality reduction
 - **Feature extraction:** Transforming the existing features into a lower dimensional space
 - **Feature selection:** Selecting a subset of the existing features without a transformation



- **Feature extraction was covered in lectures 9 and 10**
 - We derived the “optimal” linear features for two objective functions
 - **Signal representation: PCA**
 - **Signal classification: LDA**
- **Feature selection, also called Feature Subset Selection (FSS) in the literature, will be the subject of the last two lectures**
 - Although **FSS can be thought of as a special case of feature extraction** (think of a sparse projection matrix with a few ones), in practice it is a quite different problem
 - FSS looks at the issue of dimensionality reduction from a different perspective
 - FSS has a unique set of methodologies



Feature Subset Selection

■ Definition

- Given a feature set $X = \{x_i \mid i=1 \dots N\}$, find a subset $Y_M = \{x_{i_1}, x_{i_2}, \dots, x_{i_M}\}$, with $M < N$, that optimizes an objective function $J(Y)$, ideally the probability of correct classification

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \xrightarrow{\text{feature selection}} \begin{bmatrix} x_{i_1} \\ x_{i_2} \\ \vdots \\ x_{i_M} \end{bmatrix} \quad \{x_{i_1}, x_{i_2}, \dots, x_{i_M}\} = \underset{M, i_m}{\operatorname{argmax}} [J\{x_i \mid i = 1 \dots N\}]$$

■ Why Feature Subset Selection?

- Why not use the more general feature extraction methods, and simply project a high-dimensional feature vector onto a low-dimensional space?

■ Feature Subset Selection is necessary in a number of situations

- Features may be expensive to obtain
 - You evaluate a large number of features (sensors) in the test bed and select only a few for the final implementation
- You may want to extract meaningful rules from your classifier
 - When you transform or project, the measurement units of your features (length, weight, etc.) are lost
- Features may not be numeric
 - A typical situation in the machine learning domain

■ In addition, fewer features means fewer parameters for pattern recognition

- Improved the generalization capabilities
- Reduced complexity and run-time



Search strategy and objective function

■ Feature Subset Selection requires

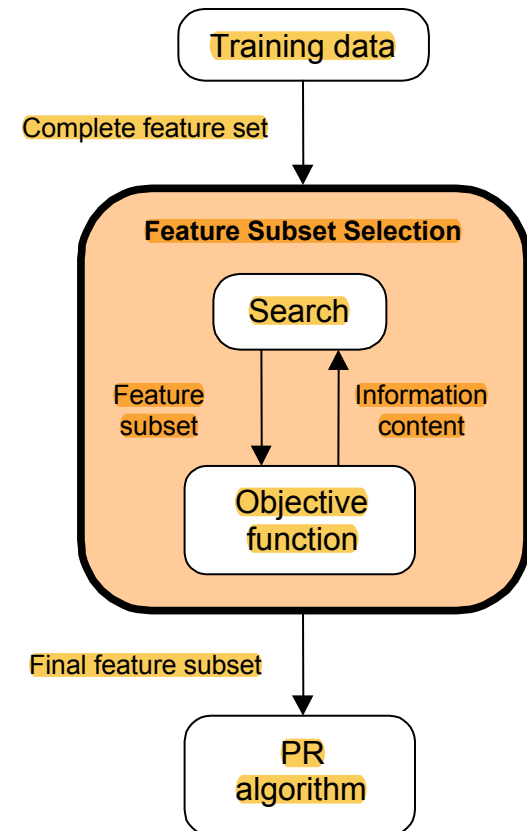
- A search strategy to select candidate subsets
- An objective function to evaluate these candidates

■ Search Strategy

- Exhaustive evaluation of feature subsets involves $\binom{N}{M}$ combinations for a fixed value of M, and 2^N combinations if M must be optimized as well
 - This number of combinations is unfeasible, even for moderate values of M and N, so a search procedure must be used in practice
 - For example, exhaustive evaluation of 10 out of 20 features involves 184,756 feature subsets; exhaustive evaluation of 10 out of 100 involves more than 10^{13} feature subsets [Devijver and Kittler, 1982]
- A search strategy is therefore needed to direct the FSS process as it explores the space of all possible combination of features

■ Objective Function

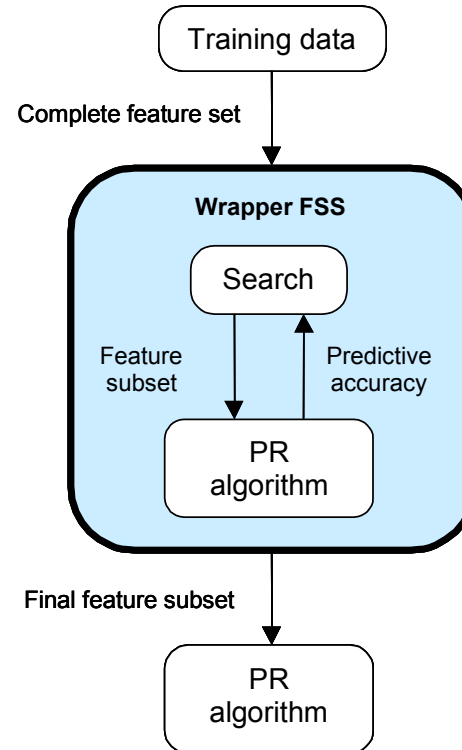
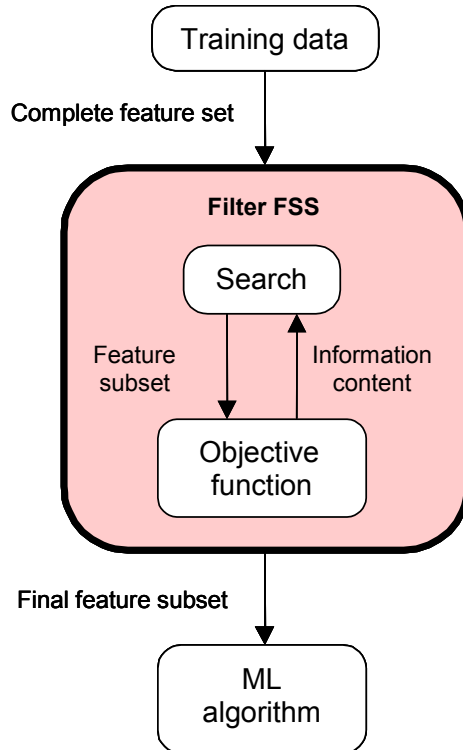
- The objective function evaluates candidate subsets and returns a measure of their “goodness”, a feedback signal used by the search strategy to select new candidates



Objective function

■ Objective functions are divided in two groups

- **Filters:** The objective function evaluates feature subsets by their information content, typically interclass distance, statistical dependence or information-theoretic measures
- **Wrappers:** The objective function is a pattern classifier, which evaluates feature subsets by their predictive accuracy (recognition rate on test data) by statistical resampling or cross-validation



Filter types

■ Distance or separability measures

- These methods use distance metrics to measure class separability, such as
 - Distance between classes: Euclidean, Mahalanobis, etc.
 - Determinant of $S_W^{-1}S_B$ (LDA eigenvalues)

■ Correlation and information-theoretic measures

- These methods are based on the rationale that good feature subsets contain features highly correlated with (predictive of) the class, yet uncorrelated with (not predictive of) each other

• Linear relation measures

- Linear relationship between variables can be measured using the correlation coefficient

$$J(Y_M) = \frac{\sum_{i=1}^M \rho_{ic}}{\sum_{i=1}^M \sum_{j=i+1}^M \rho_{ij}}$$

- Where ρ_{ic} is the correlation coefficient between feature 'i' and the class label and ρ_{ij} is the correlation coefficient between features 'i' and 'j'

• Non-Linear relation measures

- Correlation is only capable of measuring linear dependence. A more powerful measure is the mutual information $I(Y_k; C)$

$$J(Y_M) = I(Y_M; C) = H(C) - H(C | Y_M) = \sum_{c=1}^C \int_{Y_M} P(Y_M, \omega_c) \lg \frac{P(Y_M, \omega_c)}{P(Y_M) P(\omega_c)} dx$$

- The mutual information between the feature vector and the class label $I(Y_M; C)$ measures the amount by which the uncertainty in the class $H(C)$ is decreased by knowledge of the feature vector $H(C|Y_M)$, where $H(\cdot)$ is the entropy function

- Note that mutual information requires the computation of the multivariate densities $P(Y_M)$ and $P(Y_M, \omega_c)$, which is an ill-posed problem for high-dimensional spaces. In practice [Battiti, 1994], mutual information is replaced by a heuristic like

$$J(Y_M) = \sum_{m=1}^M I(x_{i_m}; C) - \beta \sum_{m=1}^M \sum_{n=m+1}^M I(x_{i_m}; x_{i_n})$$



Filters vs. Wrappers

■ Filters

• Advantages

- **Fast execution:** Filters generally involve a non-iterative computation on the dataset, which can execute much faster than a classifier training session
- **Generality:** Since filters evaluate the intrinsic properties of the data, rather than their interactions with a particular classifier, their results exhibit more generality: the solution will be “good” for a larger family of classifiers

• Disadvantages

- **Tendency to select large subsets:** Since the filter objective functions are generally monotonic, the filter tends to select the full feature set as the optimal solution. This forces the user to select an arbitrary cutoff on the number of features to be selected

■ Wrappers

• Advantages

- **Accuracy:** wrappers generally achieve better recognition rates than filters since they are tuned to the specific interactions between the classifier and the dataset
- **Ability to generalize:** wrappers have a mechanism to avoid overfitting, since they typically use cross-validation measures of predictive accuracy

• Disadvantages

- **Slow execution:** since the wrapper must train a classifier for each feature subset (or several classifiers if cross-validation is used), the method can become unfeasible for computationally intensive methods
- **Lack of generality:** the solution lacks generality since it is tied to the bias of the classifier used in the evaluation function. The “optimal” feature subset will be specific to the classifier under consideration



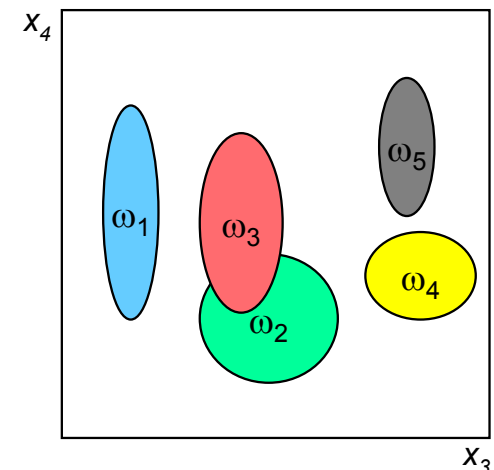
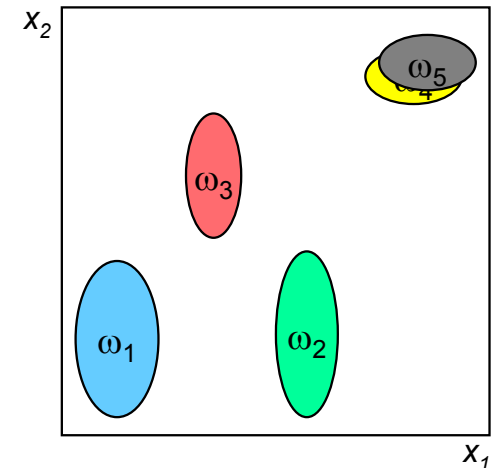
Search strategies

- **There is a large number of search strategies, which can be grouped in three categories**
 - Exponential algorithms (Lecture 12)
 - These algorithms evaluate a number of subsets that grows exponentially with the dimensionality of the search space
 - The most representative algorithms under this class are
 - Exhaustive Search (already discussed)
 - Branch and Bound
 - Approximate Monotonicity with Branch and Bound
 - Beam Search
 - Sequential algorithms (Lecture 11)
 - These algorithms add or remove features sequentially, but have a tendency to become trapped in local minima
 - Representative examples of sequential search include
 - Sequential Forward Selection
 - Sequential Backward Selection
 - Plus-I Minus-r Selection
 - Bidirectional Search
 - Sequential Floating Selection
 - Randomized algorithms (Lecture 12)
 - These algorithms incorporating randomness into their search procedure to escape local minima
 - Representative examples are
 - Random Generation plus Sequential Selection
 - Simulated Annealing
 - Genetic Algorithms



Naïve sequential feature selection

- One may be tempted to evaluate each individual feature separately and select those M features with the highest scores
 - Unfortunately, this strategy will VERY RARELY work since it does not account for feature dependence
- An example will help illustrate the poor performance that can be expected from this naïve approach
 - The figures show a 4-dimensional pattern recognition problem with 5 classes. Features are shown in pairs of 2D scatter plots
 - The objective is to select the best subset of 2 features using the naïve sequential feature selection procedure
 - Any reasonable objective function will rank features according to this sequence: $J(x_1) > J(x_2) \approx J(x_3) > J(x_4)$
 - x_1 is, without a doubt, the best feature. It clearly separates $\omega_1, \omega_2, \omega_3$ and $\{\omega_4, \omega_5\}$
 - x_2 and x_3 have similar performance, separating classes in three groups
 - x_4 is the worst feature since it can only separate ω_4 from ω_5 , the rest of the classes having a heavy overlap
 - The optimal feature subset turns out to be $\{x_1, x_4\}$, because x_4 provides the only information that x_1 needs: discrimination between classes ω_4 and ω_5
 - However, if we were to choose features according to the individual scores $J(x_k)$, we would certainly pick x_1 and either x_2 or x_3 , leaving classes ω_4 and ω_5 non separable
 - This naïve strategy fails because it cannot consider features with complementary information



Sequential Forward Selection (SFS)

■ Sequential Forward Selection is the simplest greedy search algorithm

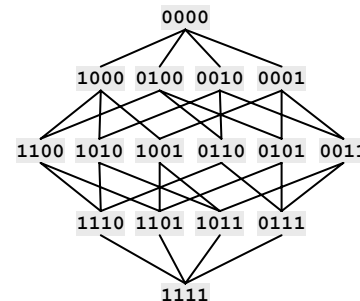
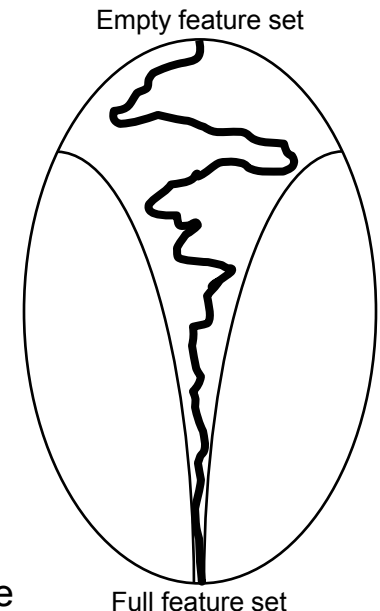
- Starting from the empty set, sequentially add the feature x^+ that results in the highest objective function $J(Y_k + x^+)$ when combined with the features Y_k that have already been selected

■ Algorithm

1. Start with the empty set $Y_0 = \{\emptyset\}$
2. Select the next best feature $x^+ = \underset{x \notin Y_k}{\operatorname{argmax}} [J(Y_k + x)]$
3. Update $Y_{k+1} = Y_k + x^+$; $k = k + 1$
4. Go to 2

■ Notes

- SFS performs best when the optimal subset has a small number of features
 - When the search is near the empty set, a large number of states can be potentially evaluated
 - Towards the full set, the region examined by SFS is narrower since most of the features have already been selected
- The search space is drawn like an ellipse to emphasize the fact that there are fewer states towards the full or empty sets
 - As an example, the state space for 4 features is shown. Notice that the number of states is larger in the middle of the search tree
 - The main disadvantage of SFS is that it is unable to remove features that become obsolete after the addition of other features



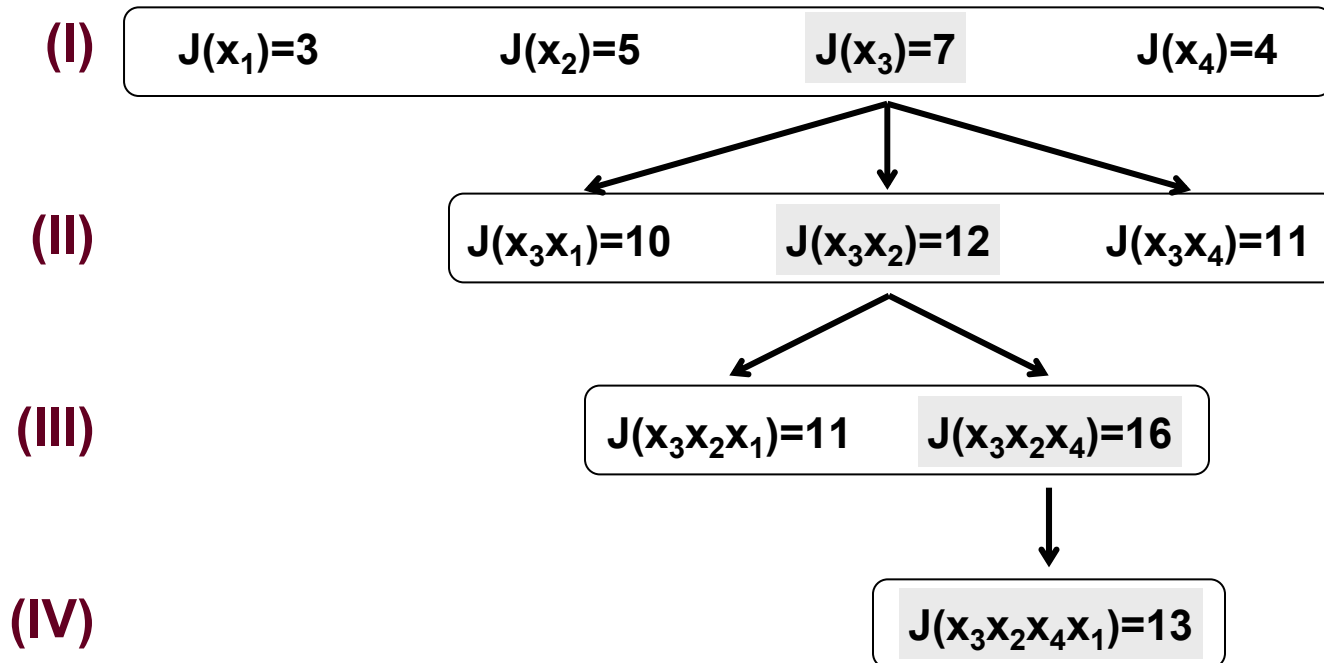
SFS example

- Assuming the objective function $J(\underline{X})$ below, perform a Sequential Forward Selection to completion

$$J(\underline{X}) = -2x_1x_2 + 3x_1 + 5x_2 - 2x_1x_2x_3 + 7x_3 + 4x_4 - 2x_1x_2x_3x_4$$

- where x_k are *indicator variables* that determine if the k-th feature has been selected ($x_k=1$) or not ($x_k=0$)

Solution



Sequential Backward Selection (SBS)

■ Sequential Backward Selection works in the opposite direction of SFS

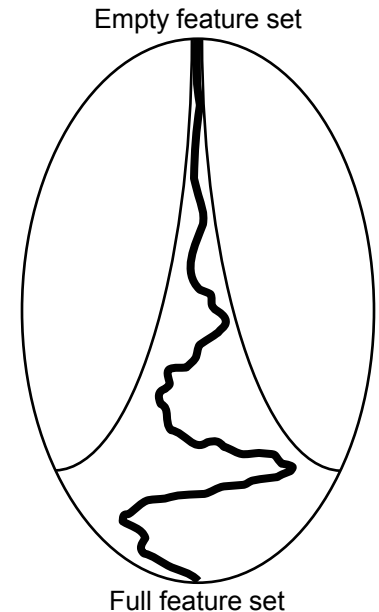
- Starting from the full set, sequentially remove the feature x^- that results in the smallest decrease in the value of the objective function $J(Y-x^-)$
 - Notice that removal of a feature may actually lead to an increase in the objective function $J(Y_k-x^-) > J(Y_k)$. Such functions are said to be *non-monotonic* (more on this when we cover Branch and Bound)

■ Algorithm

1. Start with the full set $Y_0=X$
2. Remove the worst feature $x^- = \arg\max_{x \in Y_k} [J(Y_k - x)]$
3. Update $Y_{k+1}=Y_k-x^-$; $k=k+1$
4. Go to 2

■ Notes

- SBS works best when the optimal feature subset has a large number of features, since SBS spends most of its time visiting large subsets
- The main limitation of SBS is its inability to reevaluate the usefulness of a feature after it has been discarded



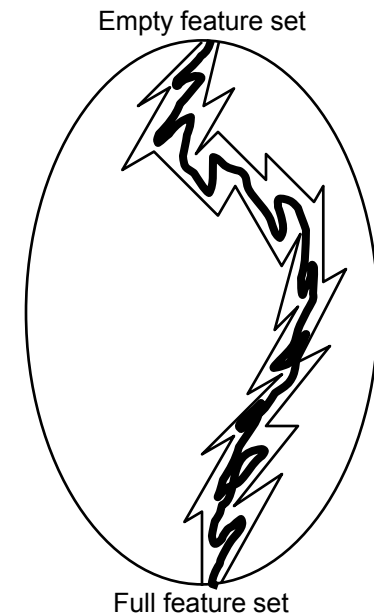
Plus-L Minus-R Selection (LRS)

■ Plus-L Minus-R is a generalization of SFS and SBS

- If $L > R$, LRS starts from the empty set and repeatedly adds 'L' features and removes 'R' features
- If $L < R$, LRS starts from the full set and repeatedly removes 'R' features followed by 'L' feature additions

■ Algorithm

1. If $L > R$ then
start with the empty set $Y = \{\emptyset\}$
else
start with the full set $Y = X$
go to step 3
2. Repeat L times
$$x^+ = \operatorname{argmax}_{x \notin Y_k} [J(Y_k + x)]$$
$$Y_{k+1} = Y_k + x^+; \quad k = k + 1$$
3. Repeat R times
$$x^- = \operatorname{argmax}_{x \in Y_k} [J(Y_k - x)]$$
$$Y_{k+1} = Y_k - x^-; \quad k = k + 1$$
4. Go to 2



■ Notes

- LRS attempts to compensate for the weaknesses of SFS and SBS with some backtracking capabilities
- Its main limitation is the lack of a theory to help predict the optimal values of L and R



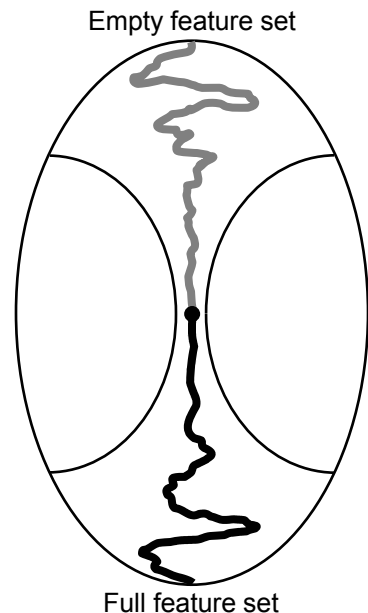
Bidirectional Search (BDS)

■ Bidirectional Search is a parallel implementation of SFS and SBS

- SFS is performed from the empty set
- SBS is performed from the full set
- To guarantee that SFS and SBS converge to the same solution, we must ensure that
 - Features already selected by SFS are not removed by SBS
 - Features already removed by SBS are not selected by SFS
 - For example, before SFS attempts to add a new feature, it checks if it has been removed by SBS and, if it has, attempts to add the second best feature, and so on. SBS operates in a similar fashion.

■ Algorithm

1. Start SFS with the empty set $Y_F = \{\emptyset\}$
2. Start SBS with the full set $Y_B = X$
3. Select the best feature
$$x^+ = \underset{\substack{x \notin Y_{F_k} \\ x \in Y_{B_k}}}{\operatorname{argmax}} [J(Y_{F_k} + x)]$$
$$Y_{F_{k+1}} = Y_{F_k} + x^+$$
3. Remove the worst feature
$$x^- = \underset{\substack{x \in Y_{B_k} \\ x \notin Y_{F_{k+1}}}}{\operatorname{argmax}} [J(Y_{B_k} - x)]$$
$$Y_{B_{k+1}} = Y_{B_k} - x^-; \quad k = k + 1$$
4. Go to 2



Sequential Floating Selection (SFFS and SFBS)

- **Sequential Floating Selection methods are an extension to the LRS algorithms with flexible backtracking capabilities**
 - Rather than fixing the values of 'L' and 'R', these *floating* methods allow those values to be determined from the data:
 - The dimensionality of the subset during the search can be thought to be “floating” up and down
- **There are two floating methods**
 - Sequential Floating Forward Selection (SFFS) starts from the empty set
 - After each forward step, SFFS performs backward steps as long as the objective function increases
 - Sequential Floating Backward Selection (SFBS) starts from the full set
 - After each backward step, SFBS performs forward steps as long as the objective function increases
- **SFFS Algorithm (SFBS is analogous)**

```
1. Start with the empty set  $Y=\{\emptyset\}$ 
2. Select the best feature

$$x^+ = \operatorname{argmax}_{x \notin Y_k} [J(Y_k + x)]$$


$$Y_k = Y_k + x^+; \quad k = k + 1$$

3. Select the worst feature*

$$x^- = \operatorname{argmax}_{x \in Y_k} [J(Y_k - x)]$$

4. If  $J(Y_k - x^-) > J(Y_k)$  then
   
$$Y_{k+1} = Y_k - x^-; \quad k = k + 1$$

   go to Step 3
else
   go to Step 2
```

**Notice that you'll need to do some book-keeping to avoid infinite loops*

