

**Applied Machine Learning**

# Imputation and Feature Selection

Yasin Ceran

# Dealing with missing values

- Missing values can be encoded in many ways
- Numpy has no standard format for it (often `np.NaN`) - pandas does
- Sometimes: 999, ???, ?, `np.inf`, “N/A”, “Unknown“ ...
- Not discussing “missing output” - that’s semi-supervised learning.
- Often missingness is informative!

```

array([[ nan, 3.2, 5.7, 2.3],
       [ nan, 2.8, 4.9, 2. ],
       [ nan, 2.8, 6.7, 2. ],
       [ nan, 2.7, 4.9, 1.8],
       [ 6.7, 3.3, 5.7, 2.1],
       [ nan, 3.2, 6. , 1.8],
       [ nan, 2.8, 4.8, 1.8],
       [ nan, 3. , 4.9, 1.8],
       [ nan, 2.8, 5.6, 2.1],
       [ nan, 3. , 5.8, 1.6],
       [ 7.4, 2.8, 6.1, 1.9],
       [ nan, 3.8, 6.4, 2. ],
       [ 6.4, 2.8, 5.6, 2.2],
       [ nan, 2.8, 5.1, 1.5],
       [ nan, 2.6, 5.6, 1.4],
       [ nan, 3. , 6.1, 2.3],
       [ nan, 3.4, 5.6, 2.4],
       [ nan, 3.1, 5.5, 1.8],
       [ nan, 3. , 4.8, 1.8],
       [ 6.9, 3.1, 5.4, 2.1],
       [ 6.7, 3.1, 5.6, 2.4],
       [ nan, 3.1, 5.1, 2.3],
       [ nan, 2.7, 5.1, 1.9],
       [ nan, 3.2, 5.9, 2.3],
       [ nan, 3.3, 5.7, 2.5],
       [ nan, 3. , 5.2, 2.3],
       [ nan, 2.5, 5. , 1.9],
       [ nan, 3. , 5.2, 2. ],
       [ 6.2, 3.4, 5.4, 2.3],
       [ nan, 3. , 5.1, 1.8]])

```

```

array([[ 5.1, 3.5, 1.4, 0.2],
       [ nan, nan, 1.4, 0.2],
       [ 4.7, 3.2, 1.3, 0.2],
       [ 4.6, 3.1, 1.5, 0.2],
       [ 5. , 3.6, 1.4, 0.2],
       [ nan, nan, nan, nan],
       [ 4.6, 3.4, 1.4, 0.3],
       [ 5. , 3.4, 1.5, 0.2],
       [ 4.4, 2.9, 1.4, 0.2],
       [ 4.9, 3.1, 1.5, 0.1],
       [ 5.4, 3.7, 1.5, 0.2],
       [ 4.8, 3.4, 1.6, 0.2],
       [ 4.8, 3. , 1.4, 0.1],
       [ 4.3, 3. , 1.1, 0.1],
       [ nan, nan, nan, nan],
       [ 5.7, 4.4, 1.5, 0.4],
       [ 5.4, 3.9, 1.3, 0.4],
       [ 5.1, 3.5, 1.4, 0.3],
       [ 5.7, 3.8, 1.7, 0.3],
       [ 5.1, 3.8, 1.5, 0.3],
       [ 5.4, 3.4, 1.7, 0.2],
       [ 5.1, 3.7, 1.5, 0.4],
       [ 4.6, 3.6, 1. , 0.2],
       [ 5.1, nan, nan, nan],
       [ 4.8, 3.4, 1.9, 0.2],
       [ 5. , 3. , 1.6, 0.2],
       [ nan, nan, nan, 0.4],
       [ 5.2, 3.5, 1.5, 0.2],
       [ 5.2, 3.4, 1.4, 0.2],
       [ 4.7, 3.2, 1.6, 0.2]])

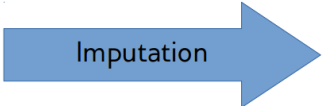
```

```

array([[ 6. ,  3.4,  4.5, nan],
       [ 6.9,  3.1,  5.1,  2.3],
       [ 4.6,  3.2,  1.4, nan],
       [ 5.1,  3.8,  1.5, nan],
       [ 4.4,  2.9, nan, nan],
       [ 6.6,  2.9,  4.6,  1.3],
       [ 6.7,  3. ,  5.2,  2.3],
       [ 6.3,  3.3,  6. ,  2.5],
       [ 7.2,  3. ,  5.8,  1.6],
       [ 4.6,  3.4,  1.4, nan],
       [ 5.2,  3.5,  1.5, nan],
       [ 5.4,  3.4, nan, nan],
       [ 5.9,  3.2,  4.8,  1.8],
       [ 4.9,  3.1, nan, nan],
       [ 6.9,  3.2,  5.7,  2.3],
       [ 5.7,  3.8, nan,  0.3],
       [ 5.3,  3.7,  1.5, nan],
       [ 4.5,  2.3,  1.3, nan],
       [ 6.5,  3. ,  5.5,  1.8],
       [ 6.2,  2.9,  4.3,  1.3],
       [ 6.4,  2.8,  5.6,  2.2],
       [ 6.1,  3. ,  4.6,  1.4],
       [ 6.2,  2.8,  4.8,  1.8],
       [ 4.9,  2.5,  4.5,  1.7],
       [ 6. ,  2.7,  5.1, nan],
       [ 6.8,  3.2,  5.9,  2.3],
       [ 6. ,  2.9,  4.5,  1.5],
       [ 4.9,  2.4,  3.3,  1. ],
       [ 5.8,  2.7,  5.1, nan],
       [ 5.5,  2.4,  3.8, nan]])

```

Imputation



```

array([[ 6. ,  3.4,  4.5,  1.6],
       [ 6.9,  3.1,  5.1,  2.3],
       [ 4.6,  3.2,  1.4,  0.2],
       [ 5.1,  3.8,  1.5,  0.3],
       [ 4.4,  2.9,  1.4,  0.2],
       [ 6.6,  2.9,  4.6,  1.3],
       [ 6.7,  3. ,  5.2,  2.3],
       [ 6.3,  3.3,  6. ,  2.5],
       [ 7.2,  3. ,  5.8,  1.6],
       [ 4.6,  3.4,  1.4,  0.3],
       [ 5.2,  3.5,  1.5,  0.2],
       [ 5.4,  3.4,  1.5,  0.4],
       [ 5.9,  3.2,  4.8,  1.8],
       [ 4.9,  3.1,  1.5,  0.1],
       [ 6.9,  3.2,  5.7,  2.3],
       [ 5.7,  3.8,  1.7,  0.3],
       [ 5.3,  3.7,  1.5,  0.2],
       [ 4.5,  2.3,  1.3,  0.3],
       [ 6.5,  3. ,  5.5,  1.8],
       [ 6.2,  2.9,  4.3,  1.3],
       [ 6.4,  2.8,  5.6,  2.2],
       [ 6.1,  3. ,  4.6,  1.4],
       [ 6.2,  2.8,  4.8,  1.8],
       [ 4.9,  2.5,  4.5,  1.7],
       [ 6. ,  2.7,  5.1,  1.6],
       [ 6.8,  3.2,  5.9,  2.3],
       [ 6. ,  2.9,  4.5,  1.5],
       [ 4.9,  2.4,  3.3,  1. ],
       [ 5.8,  2.7,  5.1,  1.9],
       [ 5.5,  2.4,  3.8,  1.1]])

```

# Imputation Methods

- Mean / Median
- kNN
- Regression models
- Probabilistic models

# Baseline: Dropping Columns

```
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler

nan_columns = np.any(np.isnan(X_train), axis=0)
X_drop_columns = X_train[:, ~nan_columns]
logreg = make_pipeline(StandardScaler(), LogisticRegression(solver='lbfgs', multi_class='multinomial'))
scores = cross_val_score(logreg, X_drop_columns, y_train, cv=10)
np.mean(scores)
```

0.767

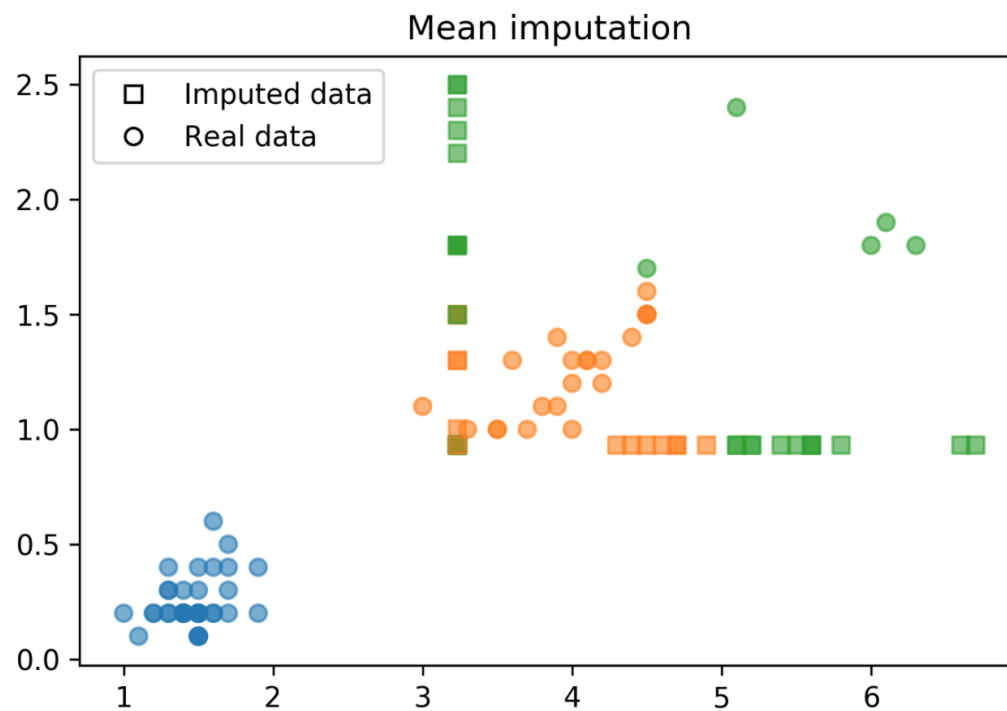
# Mean and Median

```
[[ 6.  2.9  4.5  1.5]
 [ 5.9  3.  5.1  1.8]
 [ 4.4  3.  1.3  0.2]
 [ 5.1  3.3  nan  nan]
 [ 5.  3.5  1.6  0.6]
 [ 5.4  3.4  nan  nan]
 [ 5.7  3.8  nan  0.3]
 [ 5.6  2.5  3.9  nan]
 [ 7.7  2.6  6.9  2.3]
 [ 5.8  2.7  5.1  1.9]
 [ 6.7  3.1  5.6  2.4]
 [ 4.8  3.4  1.9  nan]
 [ 7.2  3.2  6.  1.8]
 [ 4.4  2.9  nan  nan]
 [ 6.9  3.2  5.7  2.3]
 [ 5.5  4.2  1.4  nan]
 [ 6.3  2.3  4.4  1.3]
 [ 7.  3.2  4.7  1.4]
 [ 5.8  2.7  nan  nan]
 [ 6.8  2.8  4.8  1.4]
 [ 5.4  3.9  1.7  nan]
 [ 7.6  3.  6.6  2.1]
 [ 7.7  2.8  6.7  2. ]
 [ 5.  3.3  nan  0.2]
 [ 5.9  3.  4.2  1.5]
 [ 6.1  2.8  4.  1.3]
 [ 5.  3.6  1.4  0.2]
 [ 7.4  2.8  6.1  1.9]
 [ 6.3  2.5  5.  1.9]
 [ 6.7  3.3  5.7  2.5]]
```

```
from sklearn.preprocessing import Imputer
imp = Imputer(strategy="mean").fit(X_train)
imp.transform(X_train)[-30:]
```

Imputation

```
array([[ 6.  ,  2.9 ,  4.5 ,  1.5 ],
 [ 5.9 ,  3.  ,  5.1 ,  1.8 ],
 [ 4.4 ,  3.  ,  1.3 ,  0.2 ],
 [ 5.1 ,  3.3 ,  4.116,  1.462],
 [ 5.  ,  3.5 ,  1.6 ,  0.6 ],
 [ 5.4 ,  3.4 ,  4.116,  1.462],
 [ 5.7 ,  3.8 ,  4.116,  0.3 ],
 [ 5.6 ,  2.5 ,  3.9 ,  1.462],
 [ 7.7 ,  2.6 ,  6.9 ,  2.3 ],
 [ 5.8 ,  2.7 ,  5.1 ,  1.9 ],
 [ 6.7 ,  3.1 ,  5.6 ,  2.4 ],
 [ 4.8 ,  3.4 ,  1.9 ,  1.462],
 [ 7.2 ,  3.2 ,  6.  ,  1.8 ],
 [ 4.4 ,  2.9 ,  4.116,  1.462],
 [ 6.9 ,  3.2 ,  5.7 ,  2.3 ],
 [ 5.5 ,  4.2 ,  1.4 ,  1.462],
 [ 6.3 ,  2.3 ,  4.4 ,  1.3 ],
 [ 7.  ,  3.2 ,  4.7 ,  1.4 ],
 [ 5.8 ,  2.7 ,  4.116,  1.462],
 [ 6.8 ,  2.8 ,  4.8 ,  1.4 ],
 [ 5.4 ,  3.9 ,  1.7 ,  1.462],
 [ 7.6 ,  3.  ,  6.6 ,  2.1 ],
 [ 7.7 ,  2.8 ,  6.7 ,  2.  ],
 [ 5.  ,  3.3 ,  4.116,  0.2 ],
 [ 5.9 ,  3.  ,  4.2 ,  1.5 ],
 [ 6.1 ,  2.8 ,  4.  ,  1.3 ],
 [ 5.  ,  3.6 ,  1.4 ,  0.2 ],
 [ 7.4 ,  2.8 ,  6.1 ,  1.9 ],
 [ 6.3 ,  2.5 ,  5.  ,  1.9 ],
 [ 6.7 ,  3.3 ,  5.7 ,  2.5 ]])
```





```
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScalar

nan_columns = np.any(np.isnan(X_train), axis = 0)
X_drop_columns = X_train[:,~nan_columns]
logreg = make_pipeline(StandardScalar(), LogisticRegression())
scores = cross_val_score(logreg, X_drop_columns, y_train, cv = 10)
np.mean(scores)
```

0.794

```
mean_pipe = make_pipeline(Imputer(), StandardScalar(),
                           LogisticRegression())
scores = cross_val_score(mean_pipe, X_train, y_train, cv = 10)
np.mean(scores)
```

0.729

# KNN Imputation

- Find  $k$  nearest neighbors that have non-missing values.
- Fill in all missing values using the average of the neighbors.

# KNN Imputation

```
# Very inefficient didactic implementation

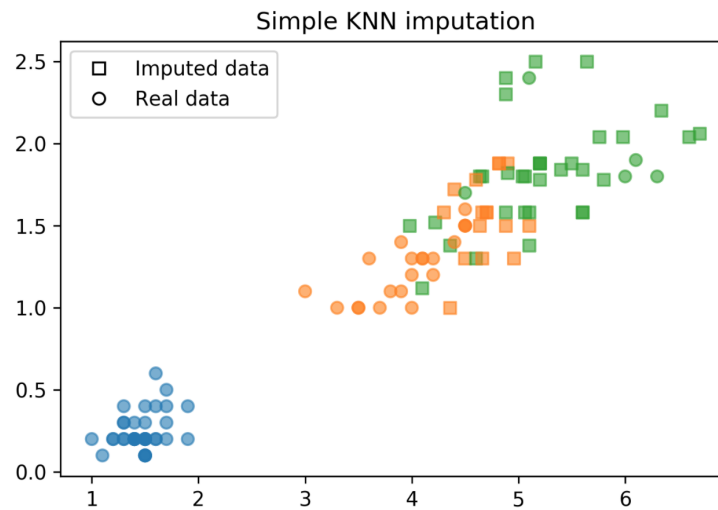
distances = np.zeros((X_train.shape[0], X_train.shape[0]))
for i, x1 in enumerate(X_train):
    for j, x2 in enumerate(X_train):
        dist = (x1 - x2) ** 2
        nan_mask = np.isnan(dist)
        distances[i, j] = dist[~nan_mask].mean() * X_train.shape[1]

neighbors = np.argsort(distances, axis=1)[: , 1:]
n_neighbors = 3

X_train_knn = X_train.copy()
for feature in range(X_train.shape[1]):
    has_missing_value = np.isnan(X_train[:, feature])
    for row in np.where(has_missing_value)[0]:
        neighbor_features = X_train[neighbors[row], feature]
        non_nan_neighbors = neighbor_features[~np.isnan(neighbor_features)]
        X_train_knn[row, feature] = non_nan_neighbors[:n_neighbors].mean()
```

```
scores = cross_val_score(logreg, X_train_knn, y_train, cv=10)  
np.mean(scores)
```

0.849



# Model-Driven Imputation

- Train regression model for missing values
- Possibly iterate: retrain after filling in
- Very flexible!

# Model-driven Imputation w RF

```
rf = RandomForestRegressor(n_estimators=100)
X_imputed = X_train.copy()

for i in range(10):
    last = X_imputed.copy()
    for feature in range(X_train.shape[1]):
        inds_not_f = np.arange(X_train.shape[1])
        inds_not_f = inds_not_f[inds_not_f != feature]
        f_missing = np.isnan(X_train[:, feature])
        rf.fit(X_imputed[~f_missing][:, inds_not_f], X_train[~f_missing, feature])

        X_imputed[f_missing, feature] = rf.predict(
            X_imputed[f_missing][:, inds_not_f])

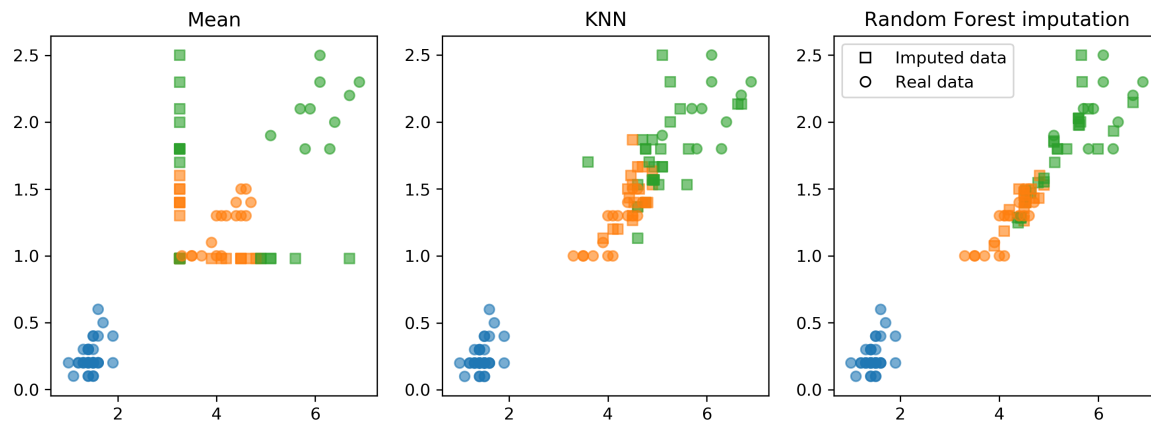
    if (np.linalg.norm(last - X_imputed)) < .5:
        break

scores = cross_val_score(logreg, X_imputed, y_train, cv=10)
np.mean(scores)
```

0.855

14 / 40

# Comparision of Imputation Methods



# Feature Selection



# Why Select Features?

- Faster prediction and training
- Less storage for model and dataset
- More interpretable model

# Types of Feature Selection

- Unsupervised vs Supervised
- Univariate vs Multivariate
- Model based or not

# Unsupervised Feature Selection

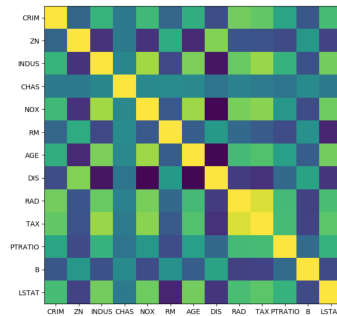
- May discard important information
- Variance-based: 0 variance or few unique values
- Covariance-based: remove correlated features
- PCA: remove linear subspaces

# Covariance

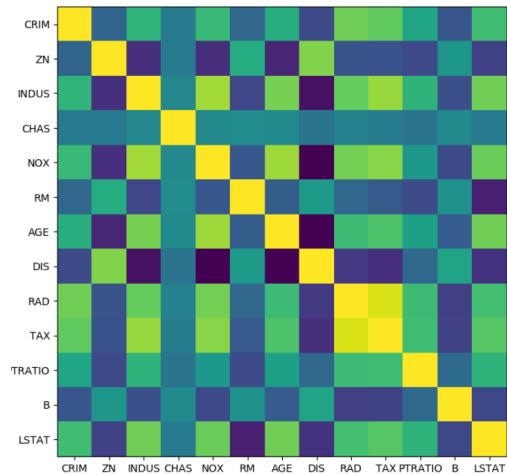
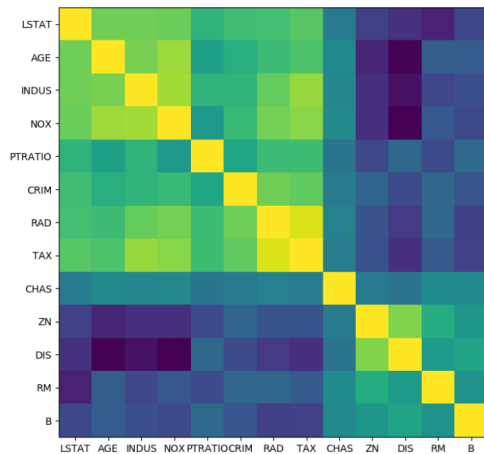
```
from sklearn.preprocessing import scale

boston = load_boston()
X, y = boston.data, boston.target
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
X_train_scaled = scale(X_train)

cov = np.cov(X_train_scaled, rowvar=False)
```



```
from scipy.cluster import hierarchy
order = np.array(hierarchy.dendrogram(
    hierarchy.ward(cov),no_plot=True)['ivl'], dtype="int")
```

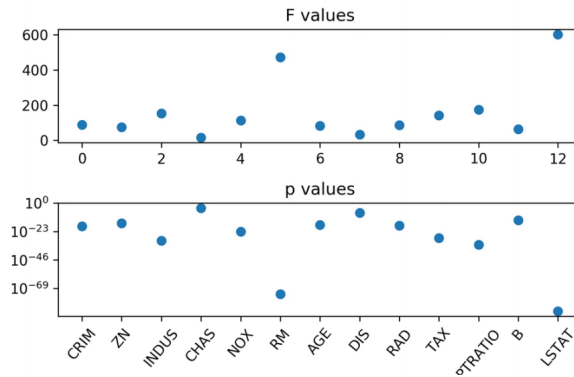


# Supervised Feature Selection

# Univariate Statistics

- Pick statistic, check p-values !
- `f_regression`, `f_classif`, `chi2` in scikit-learn

```
from sklearn.feature_selection import f_regression  
f_values, p_values = f_regression(X, y)
```



```
from sklearn.feature_selection import SelectKBest, SelectPercentile, SelectFpr  
from sklearn.linear_model import RidgeCV
```

```
select = SelectKBest(k=2, score_func=f_regression)  
select.fit(X_train, y_train)  
print(X_train.shape)  
print(select.transform(X_train).shape)
```

```
(379, 13)
```

```
(379, 2)
```



```
from sklearn.feature_selection import SelectKBest, SelectPercentile, SelectFpr
from sklearn.linear_model import RidgeCV
```

```
select = SelectKBest(k=2, score_func=f_regression)
select.fit(X_train, y_train)
print(X_train.shape)
print(select.transform(X_train).shape)
```

(379, 13)

(379, 2)

```
all_features = make_pipeline(StandardScaler(), RidgeCV())
np.mean(cross_val_score(all_features, X_train, y_train, cv=10))
```

0.718

```
from sklearn.feature_selection import SelectKBest, SelectPercentile, SelectFpr
from sklearn.linear_model import RidgeCV
```

```
select = SelectKBest(k=2, score_func=f_regression)
select.fit(X_train, y_train)
print(X_train.shape)
print(select.transform(X_train).shape)
```

```
(379, 13)
(379, 2)
```

```
all_features = make_pipeline(StandardScaler(), RidgeCV())
np.mean(cross_val_score(all_features, X_train, y_train, cv=10))
```

0.718

```
select_2 = make_pipeline(StandardScaler(),
                        SelectKBest(k=2, score_func=f_regression), RidgeCV())
np.mean(cross_val_score(select_2, X_train, y_train, cv=10))
```

0.624

# Model-Based Feature Selection

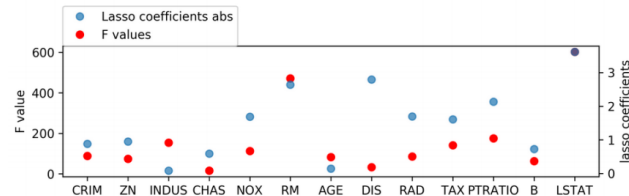
- Get best fit for a particular model
- Ideally: exhaustive search over all possible combinations
- Exhaustive is infeasible (and has multiple testing issues)
- Use heuristics in practice.

# Model based (single fit)

- Build a model, select features important to model
- Lasso, other linear models, tree-based Models
- Multivariate - linear models assume linear relation

```
from sklearn.linear_model import LassoCV
X_train_scaled = scale(X_train)
lasso = LassoCV().fit(X_train_scaled, y_train)
print(lasso.coef_)
```

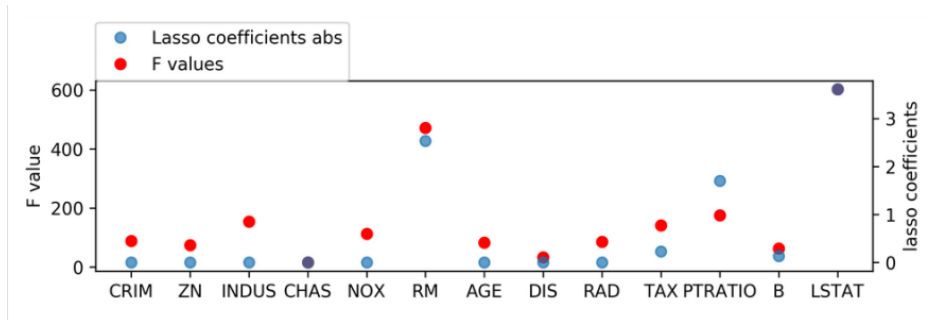
[-0.881 0.951 -0.082 0.59 -1.69 2.639 -0.146 -2.796 1.695 -1.614 -2.133 0.729 -3.615]



# Changing Lasso alpha

```
from sklearn.linear_model import Lasso
X_train_scaled = scale(X_train)
lasso = Lasso().fit(X_train_scaled, y_train)
print(lasso.coef_)
```

[-0. 0. -0. 0. -0. 2.529 -0. -0. -0. -0.228 -1.701 0.132 -3.606]



# SelectFromModel

```
from sklearn.feature_selection import SelectFromModel
select_lassocv = SelectFromModel(LassoCV(), threshold=1e-5)
select_lassocv.fit(X_train, y_train)
print(select_lassocv.transform(X_train).shape)
```

(379,11)

# SelectFromModel

```
from sklearn.feature_selection import SelectFromModel
select_lassocv = SelectFromModel(LassoCV(), threshold=1e-5)
select_lassocv.fit(X_train, y_train)
print(select_lassocv.transform(X_train).shape)
```

(379,11)

```
pipe_lassocv = make_pipeline(StandardScaler(), select_lassocv, RidgeCV())
np.mean(cross_val_score(pipe_lassocv, X_train, y_train, cv=10))
np.mean(cross_val_score(all_features, X_train, y_train, cv=10))
```

0.717

0.718

# SelectFromModel

```
from sklearn.feature_selection import SelectFromModel
select_lassocv = SelectFromModel(LassoCV(), threshold=1e-5)
select_lassocv.fit(X_train, y_train)
print(select_lassocv.transform(X_train).shape)
```

(379,11)

```
pipe_lassocv = make_pipeline(StandardScaler(), select_lassocv, RidgeCV())
np.mean(cross_val_score(pipe_lassocv, X_train, y_train, cv=10))
np.mean(cross_val_score(all_features, X_train, y_train, cv=10))
```

0.717

0.718

```
# could grid-search alpha in lasso
select_lasso = SelectFromModel(Lasso())
pipe_lasso = make_pipeline(StandardScaler(), select_lasso, RidgeCV())
np.mean(cross_val_score(pipe_lasso, X_train, y_train, cv=10))
```

0.671



# Iterative Model-Based Selection

- Fit model, find least important feature, remove, iterate.
- Or: Start with single feature, find most important feature, add, iterate.

# Recursive Feature Elimination

- Uses feature importances / coefficients, similar to “SelectFromModel”
- Iteratively removes features (one by one or in groups)
- Runtime:  $(n_{\text{features}} - n_{\text{feature\_to\_keep}}) / \text{stepsize}$

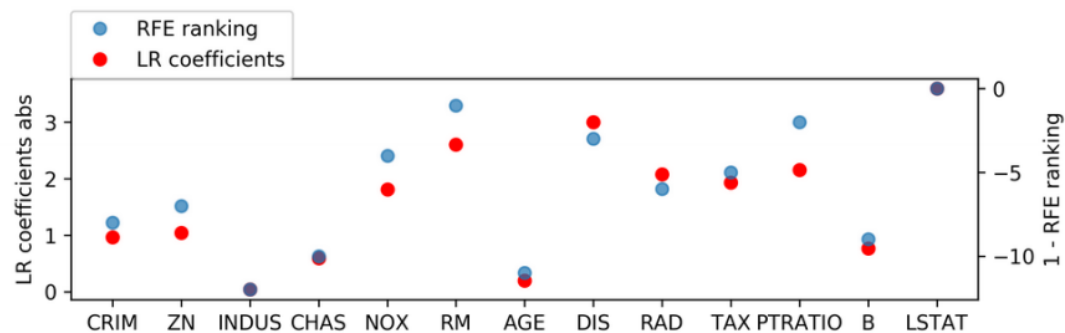
```

from sklearn.linear_model import LinearRegression
from sklearn.feature_selection import RFE

# create ranking among all features by selecting only one
rfe = RFE(LinearRegression(), n_features_to_select=1)
rfe.fit(X_train_scaled, y_train)
rfe.ranking_

```

array([ 9, 8, 13, 11, 5, 2, 12, 4, 7, 6, 3, 10, 1])



# RFECV

```
from sklearn.linear_model import LinearRegression
from sklearn.feature_selection import RFECV
rfe = RFECV(LinearRegression(), cv=10)
rfe.fit(X_train_scaled, y_train)
print(rfe.support_)
print(boston.feature_names[rfe.support_])
```

```
[ True  True False  True  True  True False  True  True  True  True  True
  True]
['CRIM' 'ZN' 'CHAS' 'NOX' 'RM' 'DIS' 'RAD' 'TAX' 'PTRATIO' 'B' 'LSTAT']
```

```
pipe_rfe_ridgecv = make_pipeline(StandardScaler(),
                                RFECV(LinearRegression(), cv=10), RidgeCV())
np.mean(cross_val_score(pipe_rfe_ridgecv, X_train, y_train, cv=10))
```

0.710

```
pipe_rfe_ridgecv = make_pipeline(StandardScaler(),  
                                RFECV(LinearRegression(), cv=10), RidgeCV())  
np.mean(cross_val_score(pipe_rfe_ridgecv, X_train, y_train, cv=10))
```

0.710

```
from sklearn.preprocessing import PolynomialFeatures  
pipe_rfe_ridgecv = make_pipeline(StandardScaler(), PolynomialFeatures(),  
                                RFECV(LinearRegression(), cv=10), RidgeCV())  
np.mean(cross_val_score(pipe_rfe_ridgecv, X_train, y_train, cv=10))
```

0.820

# Wrapper Methods

- Can be applied for ANY model!
- Shrink / grow feature set by greedy search
- Called Forward or Backward selection
- Run CV / train-val split per feature
- Complexity:  $n_{\text{features}} * (n_{\text{features}} + 1) / 2$
- Implemented in mlxtend

# SequentialFeatureSelector

```
from mlxtend.feature_selection import SequentialFeatureSelector  
sfs = SequentialFeatureSelector(LinearRegression(), forward=False, k_features=7)  
sfs.fit(X_train_scaled, y_train)
```

Features: 7/7

```
print(sfs.k_feature_idx_)  
print(boston.feature_names[np.array(sfs.k_feature_idx_)])
```

```
(1, 4, 5, 7, 9, 10, 12)  
['ZN' 'NOX' 'RM' 'DIS' 'TAX' 'PTRATIO' 'LSTAT']
```

```
sfs.k_score_
```

0.725

Questions ?