



Machine Learning Apps. Feature engineering

regression -

train_df, test_df

- Linear Regression
- $\text{linreg.fit}(X_{\text{train}}, y_{\text{train}})$
- Predict
- $y_{\text{pred}} = \text{linreg.predict}(X_{\text{test}})$

Model Selection tools/Imports

- $\text{train_test_split}(data, test_size, rs)$

Regression

- Mean Squared Error
- Mean Absolute Error $(y_{\text{test}}, y_{\text{pred}})$
- R² Score
- Weights
- Intercept

Feature Improvement + data cleaning

Scaling:

Linear Regression: **InSensitive**

KNN: **Sensitive**

Regularization LR: **Sensitive**

Trees: **InSensitive**

Usefulness of data normalization and other data trans.

The usual preprocessing in machine learning: Standardize each variable $\left\{ \begin{array}{l} \text{Mean } 0 \\ \text{variance } 1 \end{array} \right.$ **StandardScaler**

- Minmax Scaler: Scale to a given range
- Robust Scaler: Center data to median for outliers.
- MaxabsScaler: Max absolute value - Sparse data

Discretization:

- Uniform: ^{Non-linear} Separate most frequent and Concentrate the less frequent. Reduce impact of outliers.
- Gaussian: Concentrate around mean.

* Expand several orders but most data is around first order \rightarrow apply log

One-hot Encoding

Aims to use categorical values in a model that only accepts continuous variables.

* has Colinearity problems

2.1 Mathematical formulation

$$\mathbf{U} = \underset{\mathbf{U}}{\operatorname{argmax}} \operatorname{Tr} \{ \mathbf{U}^T \mathbf{X}^T \mathbf{X} \mathbf{U} \} = \underset{\mathbf{U}}{\operatorname{argmax}} \operatorname{Tr} \{ \mathbf{U}^T \mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{U} \}$$

$$\text{s.t. } \mathbf{U}^T \mathbf{U} = \mathbf{I}$$

Solution:

Let's start considering the projection onto a one-dimensional space ($K = 1$). So, we want to solve the following optimization problem:

$$\mathbf{u}_1 = \underset{\mathbf{u}_1}{\operatorname{argmax}} \mathbf{u}_1^T \mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u}_1$$

$$\text{s.t. } \mathbf{u}_1^T \mathbf{u}_1 = 1$$

If we apply Lagrange multipliers, we can include the constraint into functional by means of the lagrange multiplier λ_1 and we arrive to an unconstrained problem

$$\mathbf{u}_1 = \underset{\mathbf{u}_1}{\operatorname{argmax}} \mathbf{u}_1^T \mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

by making its derivate equal to zero, we obtain that the optimum solution has to satisfy

$$\mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

which indicates that \mathbf{u}_1 must be an eigenvector of the covariance matrix $\mathbf{C}_{\mathbf{X}\mathbf{X}}$ and λ_1 is its associated eigenvalue.

Besides, if we left-multiply this expression by \mathbf{u}_1^T and make use of $\mathbf{u}_1^T \mathbf{u}_1 = 1$, we obtain that

$$\mathbf{u}_1^T \mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u}_1 = \lambda_1$$

that is, the variance of the projected data by the eigenvector \mathbf{u}_1 is equal to its associated eigenvalue λ_1 . So, we can get the maximum projected variance, with a single projection, if we set this first principal component \mathbf{u}_1 as the eigenvector with the largest eigenvalue λ_1 . Later, we can define additional principal components in an incremental fashion by choosing each new direction as the one eigenvector with the next highest eigenvalue so that the projected variance is maximized.

So, we can obtain the first K projections of the PCA algorithm by solving the following eigenvalue problem

$$\mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u} = \lambda \mathbf{u}$$

So, the projection matrix, \mathbf{U} , consists of the first eigenvectors of $\mathbf{C}_{\mathbf{X}\mathbf{X}}$ (i.e., those associated with largest eigenvalues)

$$\mathbf{U} = \operatorname{eigs}(\mathbf{C}_{\mathbf{X}\mathbf{X}})$$

Partial Least Squares

Find projections of Input and output data with maximum covariance.

$$\mathbf{U}, \mathbf{V} = \underset{\mathbf{U}, \mathbf{V}}{\operatorname{argmax}} \operatorname{Tr} \{ \mathbf{U}^T \mathbf{X}^T \mathbf{Y} \mathbf{V} \} = \underset{\mathbf{U}, \mathbf{V}}{\operatorname{argmax}} \operatorname{Tr} \{ \mathbf{U}^T \mathbf{C}_{\mathbf{X}\mathbf{Y}} \mathbf{V} \}$$

Max # of projections = Output Classes - 1

Canonical Correlation Analysis

Find directions of max correlation between Input and output data.

$$\mathbf{u}, \mathbf{v} = \underset{\mathbf{u}, \mathbf{v}}{\operatorname{argmax}} \frac{(\mathbf{u}^T \mathbf{C}_{\mathbf{X}\mathbf{Y}} \mathbf{v})^2}{\mathbf{u}^T \mathbf{C}_{\mathbf{X}\mathbf{X}} \mathbf{u} \mathbf{v}^T \mathbf{C}_{\mathbf{Y}\mathbf{Y}} \mathbf{v}} \quad \text{and} \quad \mathbf{U}, \mathbf{V} = \underset{\mathbf{U}, \mathbf{V}}{\operatorname{argmax}} \operatorname{Tr} \{ \mathbf{U}^T \mathbf{C}_{\mathbf{X}\mathbf{Y}} \mathbf{V} \}$$

Output Classes - 1.

LDA Linear discriminant analysis

Classification method that projects each piece of data onto the real line by a linear equation $\mathbf{z} = \mathbf{W}^T \mathbf{x}$

Feature Selection

Filtering methods (Correlation, mutual info...)

Wrappers (Classification/Regression performance over CV to find features)

Embedded (Integrated into model) e.g. L1 regressor

Filtering

- Independent from learning stage
 - Variance (unsupervised)
 - Correlation coefficient (regression)
 - Statistical tests (binary)
 - Mutual Information

Do we think correlation is a good reference criteria?

Mutual Info:

Mutual Information is capable of capturing any type of relationship between two random variables, unlike correlation.

$$MI(X, Y) = \iint p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} = [0, \infty] \text{ but } MI > 2 \text{ is rare.}$$

MI between two quantities is a measure of the degree to which knowledge of one quantity reduces the uncertainty about the other. * We don't know these distributions

→ Most use MI estimators based on histograms. But sklearn's MI implementation uses K-NN based estimation.

Summary

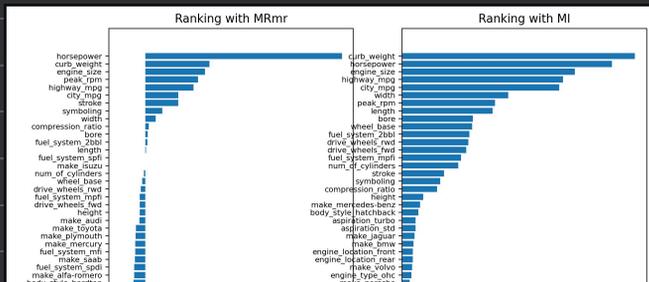
- Easy to use/interpret
- Computationally efficient
- able to detect any rel.
- May be important when interacting with another variable but not on it's own. MI cannot detect
- Usefulness depends on model being used

Multivariate Criteria

- Selecting subsets rather than individual criteria.

Minimum Redundancy Maximal Relevance (MRMR)

- Based on forward search, at each step we add feature.



$$\text{Relevance}(i) = R_{\text{REL}}(X_i, Y)$$

$$\text{Redundancy}(i) = \sum_{i' \in \text{var}_{\text{rel}}} R_{\text{RED}}(X_i, X_{i'})$$

$$\text{mrMR}(i) = \text{Relevance}(i) - \text{Redundancy}(i)$$

$$i^* = \underset{i}{\text{argmax}} \text{mrMR}(i)$$

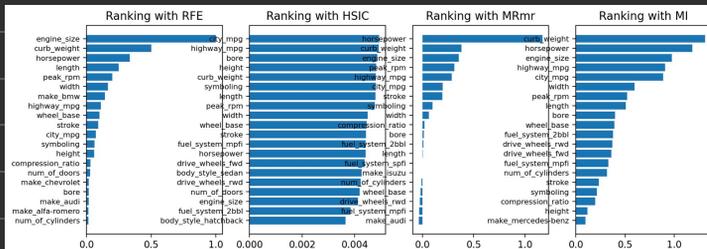
* MRMR Subselects redundancy (Correlation)

Search Strategies: $2^D - 1$ Subsets \rightarrow greedy search.

Forward, Backward

Wrapper methods

Directly use performance of final classifier/regressor to evaluate features. \rightarrow n-fold, CV process



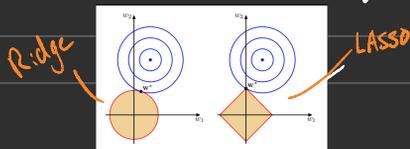
* RFE (Recursive Feature Elim.)

Embedded methods

Carry out selection during train process $\left\{ \begin{array}{l} \text{Random Forests} \\ L_2 \text{ regularization} \end{array} \right. \rightarrow$ * RF Sensitivity due to averaging between models.

L_1 regularization is used for overfitting but causes model Coef. to drop to zero

LASSO (least absolute shrinkage selection)



Let's consider a regularized linear regression model, with L_2 (Ridge Regression) and L_1 (Lasso) regularizations, to analyze this property. To do this, let's reformulate their formulations as a least squares optimization problem plus a constraint, that is

• Ridge regression:

$$w^* = \underset{w}{\text{argmin}} \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2$$

$$\text{s.t. } \sum_{j=1}^D |w_j| \leq t$$

• Lasso regression:

$$w^* = \underset{w}{\text{argmin}} \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2$$

$$\text{s.t. } \sum_{j=1}^D |w_j| \leq t$$