

Supervised learning

Aprendizaje Automático para la Robótica
Máster Universitario en Ingeniería Industrial

Departamento de Automática

Objectives

1. Extend supervised learning algorithms
2. Apply supervised learning to real-world problems

Bibliography

- Müller, Andreas C., Guido, Sarah. Introduction to Machine Learning with Python. O'Reilly. 2016

All figures have been taken from

https://github.com/amueller/introduction_to_ml_with_python/blob/master/02-supervised-learning.ipynb

Table of Contents

1. Generalization, overfitting and underfitting

2. k-Nearest Neighbors

- k-NN classification
- Scikit-Learn
- kNN regression
- Scikit-Learn
- Summary

3. Linear models

- Ordinary least squares
- Linear regression
- Regularized linear models
- Ridge regression
- Lasso regression
- ElasticNet
- Regularized linear models comparison
- Scikit-Learn

- Linear models for classification

- Scikit-Learn

- Summary

4. Decision Trees

- Scikit-Learn

- Summary

5. Ensembles of Decision Trees

- Scikit-Learn

- Summary

6. Support Vector Machines

- Kernelized Support Vector Machines

- Support Vector Machines

- Summary

7. A

- b

- A: Scikit-Learn

- A: Summary

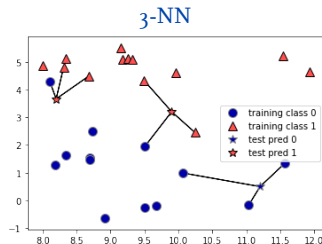
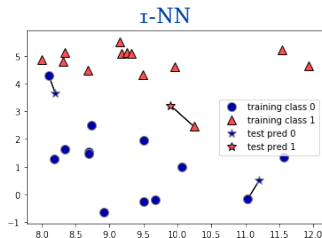
- ARIMA

k-Nearest Neighbors

k-NN classification (I)

k-NN (k-Nearest Neighbors): Likely, the simplest classifier

- Given a data point, it takes its k closest neighbors
- Same prediction than its neighbors



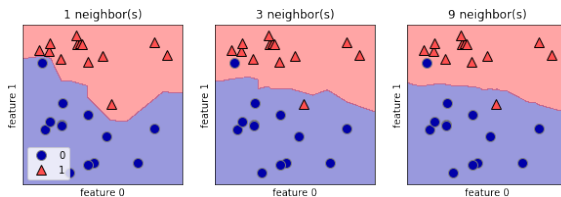
k-NN does not generate a model

- The whole dataset must be stored

k uses to be an odd number (1-NN, 3-NN, 5-NN, ...)

k-Nearest Neighbors

k-NN classification (II)



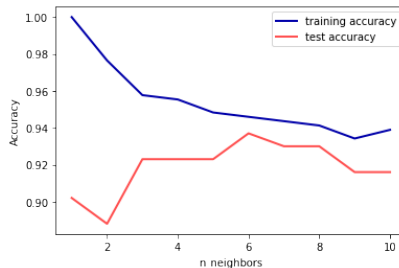
k determines the model complexity

- Smoother boundaries in larger k values
- Model complexity decreases with k
- If k equals the number of samples, k -NN always predicts the most frequent class

How to figure out the best k ?

k-Nearest Neighbors

k-NN classification (III)



k-Nearest Neighbors classifier

Scikit-learn

```
sklearn.neighbors.KNeighborsClassifier
```

Constructor arguments:

- `n_neighbors`: int, default=5
- `metric`: string, default='minkowski'
- `p`: int, default=2 ($p = 1$ Manhattan distance, $p = 2$ euclidean distance)

Attributes:

- `classes_`: ndarray(`n_samples`)

Methods: `fit()`, `predict()`

(Scikit-Learn reference)

k-Nearest Neighbors

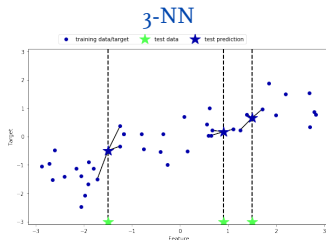
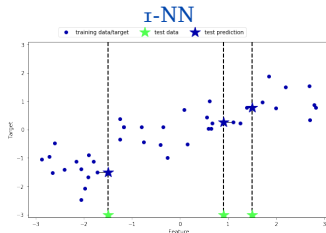
kNN regression (I)

k-NN regression

Given a data point

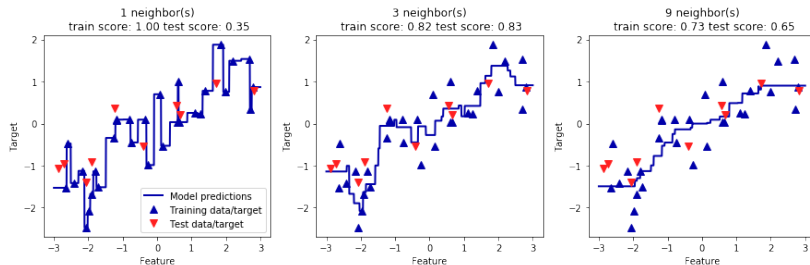
1. Take the k closest data points
2. Predict same target value (1-NN) or average target value (k-NN)

Performance is measured with a regression metric, by default, R^2



k-Nearest Neighbors

kNN regression (II)



k determines boundary smoothness

1. With $k = 1$, prediction visits all data points
2. With large k values, fit is worse

k-Nearest Neighbors regressor

Scikit-learn

`sklearn.neighbors.KNeighborsRegressor`

Constructor arguments:

- `n_neighbors`: int, default=5
- `metric`: string, default='minkowski'
- `p`: int, default=2 ($p = 1$ Manhattan distance, $p = 2$ euclidean distance)

Attributes:

Methods: `fit()`, `predict()`

(Scikit-Learn reference)

k-Nearest Neighbors

Summary

Hyperparameters	Advantages	Disadvantages
k	Simple	Slow with large datasets
Distance	Baseline	Bad performance with hundreds or more attributes
		No model
		Dataset must be stored in memory

Linear models

Linear regression

Different linear models for regression

- The difference lies in how β_i parameters are learned

Ordinary Least Squares (OLS): Minimizes mean squared error

- OLS does not have any hyperparameter
- No complexity control

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2$$

Linear regression can be used to fit non-linear models

- Just adding new attributes

Linear models

Regularized linear models

Regularization: Term that penalizes complexity

- Added to the cost function
- Linear models remain the same
- Train to minimize cost function and coefficients
- Intercepts are not part of regularization

Three regularizations

- L_1 (Lasso regression), L_2 (Ridge regression) and ElasticNet (L_1 and L_2)

Lasso (L_1)

$$\alpha \sum_j^n |\beta_j|$$

Ridge (L_2)

$$\frac{\alpha}{2} \sum_j^n \beta_j^2$$

ElasticNet

$$\alpha \left(\frac{\lambda}{2} \sum_j^n \beta_j^2 + (1 - \lambda) \sum_j^n |\beta_j| \right)$$

Linear models

Ridge regression

Ridge regression (or L2 regularization) adds a new term to cost function

$$\text{MSE} + \alpha \sum_{i=1}^n \beta_i^2$$

α controls the model complexity

- If $\alpha = 0$ Ridge becomes a regular linear regression
- Optimal α depends on the problem

Ridge by default

Linear models

Lasso regression (I)

Lasso regression (or L_1 regularization) adds a new term to cost function

$$\text{MSE} + \alpha \frac{1}{2} \sum_{i=1}^n |\beta_i|$$

α controls the model complexity

- If $\alpha = 0$ Ridge becomes a regular linear regression
- Optimal α depends on the problem

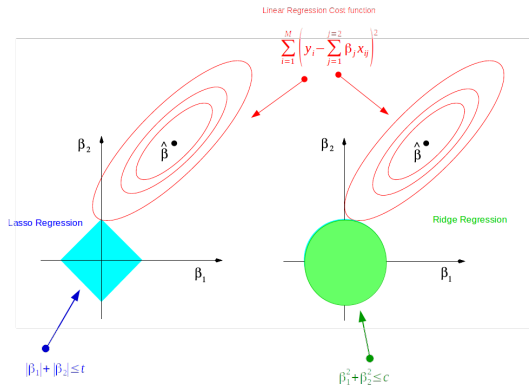
Some coefficients may be exactly zero

- Implicit feature selection
- Easier interpretation
- Better with large number of attributes

Linear models

Lasso regression (II)

Dimension Reduction of Feature Space with LASSO



(Source)

Linear models

ElasticNet

Lasso and Ridge can be combined

$$\text{MSE} + \alpha \left(\lambda \frac{1}{2} \sum_{i=1}^n |\beta_i| + (1 - \lambda) \sum_{i=1}^n \beta_i^2 \right)$$

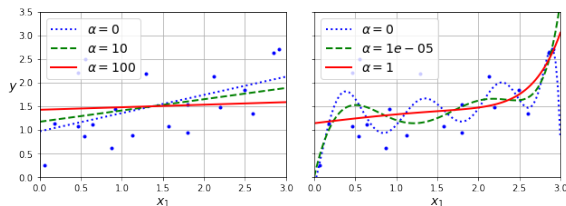
Two hyperparameters

- α controls the model complexity
- λ balances between L1 and L2

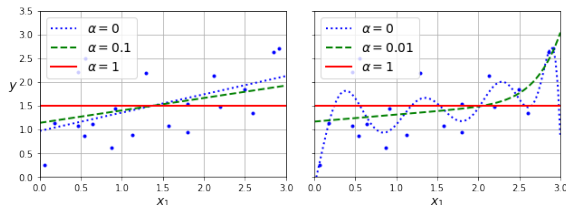
Linear models

Regularized linear models comparison

Ridge - L2



Lasso - L1



Linear models

Scikit-learn (I)

```
sklearn.linear_model.LinearRegression
```

Constructor arguments:

- `fit_intercept`: boolean, default=True

Attributes:

- `coef_`: ndarray (n_features,)
- `intercept_`: ndarray (n_targets,)
- `n_features_in_`: int

Methods: `fit()`, `predict()`

Linear models

Scikit-learn (II)

sklearn.linear_model.Ridge

Constructor arguments:

- `fit_intercept`: boolean, default=True
- `alpha`: float, default=1.0

Attributes:

- `coef_`: ndarray (n_features,)
- `intercept_`: ndarray (n_targets,)
- `n_features_in_`: int

Methods: `fit()`, `predict()`

Linear models

Scikit-learn (IV)

`sklearn.linear_model.ElasticNet`

Constructor arguments:

- `fit_intercept`: boolean, default=True
- `alpha`: float, default=1.0
- `l1_ratio`: float, default=0.5

Attributes:

- `coef_`: ndarray (n_features,)
- `intercept_`: ndarray (n_targets,)
- `n_features_in_`: int

Methods: `fit()`, `predict()`

Linear models

Scikit-learn (III)

`sklearn.linear_model.Lasso`

Constructor arguments:

- `fit_intercept`: boolean, default=True
- `alpha`: float, default=1.0

Attributes:

- `coef_`: ndarray (n_features,)
- `intercept_`: ndarray (n_targets,)
- `n_features_in_`: int

Methods: `fit()`, `predict()`

Linear models

Linear models for classification (I)

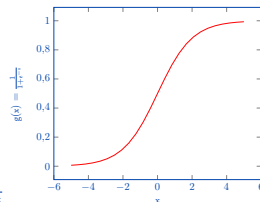
A linear regression can be used as classifier

- Just compare the prediction with a threshold (0, for instance)
 - If $\hat{y} > 0$, assign class 1
 - If $\hat{y} \leq 0$, assign class -1
- The decision boundary for any binary linear classifier is a line, plane or hyperplane

A **logistic regression** is a generalization of a linear regression

- It is a binary classifier
- Its output is a probability

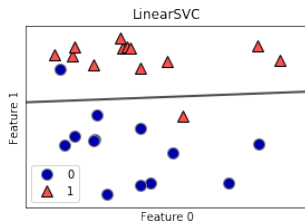
$$\hat{p} = \sigma \left(\beta_0 + \sum_{i=1}^n \beta_i x_i \right), \quad \sigma(x) = \frac{1}{1+e^{-x}}$$



where $\sigma(t)$ is the logistic function, defined as $\sigma(t) = \frac{1}{1+e^t}$

Linear models

Linear models for classification (II)

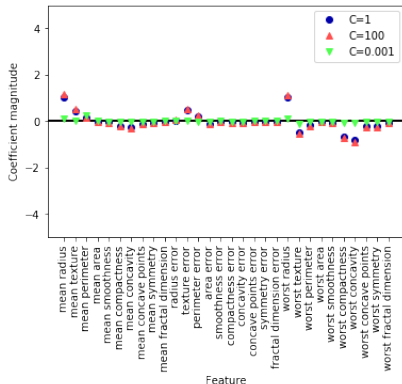


Linear models

Linear models for classification (III)

The model can be regularized with L_1 , L_2 and ElasticNet

- In Scikit-Learn, regularization strength is given by C
- Lower values of C correspond to smaller regularization strength



Linear models

Scikit-learn

`sklearn.linear_model.ElasticNet`

Constructor arguments:

- `penalty`: 'l1', 'l2', 'elasticnet', 'none', default='l2'
- `fit_intercept`: boolean, default=True
- `alpha`: float, default=1.0
- `l1_ratio`: float, default=0.5

Attributes:

- `coef_`: ndarray (n_features,)
- `intercept_`: ndarray (n_targets,)
- `n_features_in_`: int

Methods: `fit()`, `predict()`

Linear models

Summary

Hyperparameters	Advantages	Disadvantages
-	Fast train and predict	No complexity tuning
α (L1, L2, ElasticNet)	Scales well to large data-sets	Limited in low dimensional spaces
l1_ratio (ElasticNet)	Better in high dimensional spaces Few hyperparameters Interpretable	

Better when the number of features is large compared to the number of samples

Decission Trees

TODO

Decission Trees

Scikit-learn

sklearn.cluster.AgglomerativeClustering

Constructor arguments:

- linkage: 'ward', 'complete', 'average', 'single'

Attributes:

- n_clusters: int
- labels_: ndarray (n_samples)

Methods: fit(), fit_predict()

(Scikit-Learn reference)

Decission Trees

Summary

Hyperparameters	Advantages	Disadvantages

Ensembles of Decision Trees

TODO

Ensembles of Decision Trees

Ensembles of Decision Trees : Scikit-learn

```
sklearn.cluster.AgglomerativeClustering
```

Constructor arguments:

- `linkage`: 'ward', 'complete', 'average', 'single'

Attributes:

- `n_clusters`: int
- `labels_`: ndarray (n_samples)

Methods: `fit()`, `fit_predict()`

(Scikit-Learn reference)

Ensembles of Decision Trees

Summary

Hyperparameters	Advantages	Disadvantages

Support Vector Machines

TODO

Support Vector Machines

Kernelized Support Vector Machines

TODO

Scikit-Learn

Support Vector Machines

Scikit-learn

`sklearn.cluster.AgglomerativeClustering`

Constructor arguments:

- `linkage`: 'ward', 'complete', 'average', 'single'

Attributes:

- `n_clusters`: int
- `labels_`: ndarray (n_samples)

Methods: `fit()`, `fit_predict()`

(Scikit-Learn reference)

Support Vector Machines

Summary

Hyperparameters	Advantages	Disadvantages

A

B

TODO

A

B: Scikit-learn

```
sklearn.cluster.AgglomerativeClustering
```

Constructor arguments:

- `linkage`: 'ward', 'complete', 'average', 'single'

Attributes:

- `n_clusters`: int
- `labels_`: ndarray (n_samples)

Methods: `fit()`, `fit_predict()`

(Scikit-Learn reference)

A

B: Summary

Hyperparameters	Advantages	Disadvantages

Algorithms

ARIMA (I)

AR: Autoregressive model

- Current observation depends on the last p observations
- Long term memory

AR(p)

$$X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon_t$$

MA: Moving Average model

- Current observation linearly depends on the last q innovations
- Short term memory

MA(q)

$$X_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$$

ARMA model = AR + MA

- ARMA(p, q): Two hyperparameters, p and q

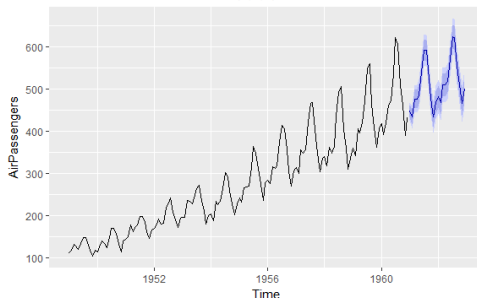
Algorithms

ARIMA (II)

ARIMA = AR + i + MA (AR integrated MA)

- ARIMA(p, d, q)
- Three integer parameters: p, q and d (in practice, low order models)

Forecasts from STL + ARIMA(1,1,1) with drift



(Source)

autoarima: search over p, q and d