https://www.sli.do/ #073374



Classical Machine Learning: Classification and Regression (V)

Learning Objectives

- Learn how to fine tune your machine learning model.
- Learn a few regression models.
- Learn the kernel method and regularization techniques.

Algorithm Tuning: AdaBoost, Gradient Boost

C-S David Chen, Department of Civil Engineering, National Taiwan University

Algorithm Tuning



- Algorithm tuning mainly involves tuning <u>hyperparameters</u>.
- Hyperparameters are parameters that are not directly learnt within estimators (algorithms).
- In scikit-learn they are passed as arguments to the constructor of the estimator classes. Typical examples include C, kernel and gamma for Support Vector Classifier, alpha for Lasso, etc.
- Two generic approaches to parameter search are provided in scikit-learn: for given values, GridSearchCV exhaustively considers all parameter combinations, while <u>RandomizedSearchCV</u> can sample a given number of candidates from a parameter space with a specified distribution.



sklearn.ensemble.AdaBoostClassifier

class sklearn.ensemble.AdaBoostClassifier(*base_estimator=None*, *, *n_estimators=50*, *learning_rate=1.0*, *algorithm='SAMME.R'*, *random_state=None*)

[source]

An AdaBoost classifier.

An AdaBoost [1] classifier is a meta-estimator that begins by fitting a classifier on the original dataset and then fits additional copies of the classifier on the same dataset but where the weights of incorrectly classified instances are adjusted such that subsequent classifiers focus more on difficult cases.

sklearn.ensemble.GradientBoostingClassifier

class sklearn.ensemble. GradientBoostingClassifier(*, loss='deviance', learning_rate=0.1, n_estimators=100, subsample=1.0, criterion='friedman_mse', min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_depth=3, min_impurity_decrease=0.0, min_impurity_split=None, init=None, random_state=None, max_features=None, verbose=0, max_leaf_nodes=None, warm_start=False, validation_fraction=0.1, n_iter_no_change=None, tol=0.0001, ccp_alpha=0.0)

Gradient Boosting for classification.

GB builds an additive model in a forward stage-wise fashion; it allows for the optimization of arbitrary differentiable loss functions. In each stage n_classes_ regression trees are fit on the negative gradient of the binomial or multinomial deviance loss function. Binary classification is a special case where only a single regression tree is induced.

Algo_Tuning.ipynb

Regression and Regression Algorithm Walkthrough

Supervised Learning: Regression

- J∈ R
 Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Extensively studied in statistics, neural network fields.
- Examples:
 - Predicting sales amounts of new product based on advertising expenditure.
 - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
 - Time series prediction of stock market indices.

Simple Linear Regression

 $y\approx w_0+w_1x$

 Predicting a quantitative response y on the basis of a single feature x (feature dimension = 1).



Multiple Linear Regression

 $y \approx w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$

 Predicting a quantitative response y on the basis of d distinct features.



Regression: Performance Metrics

 $TSS = \sum_{i=1}^{n} (y_i - \bar{y}_i)^2$ measures the total variance in the thought of as the amount of variability inherent in the response before the regression is performed. R2=0 3133

Fun Time: which R² gives us a better regression: (1) 1.0 (2) 0.5 (3) 0.0

An R^2 statistic measures the proportion of variability in y that can be explained using x. An R^2 statistic that is close to 1 indicates that a large proportion of the variability in the response has been explained by the regression.

Linear Regression.ipynb

Nonlinear Regression: Feature Transformation, Basis Functions, Feature Exploration and Kernel Method

Nonlinear Regression: Feature Transformation and Basis Function

• single feature x (feature dimension = 1).



 $y \approx w_0 + w_1 \phi(x)$ $\phi(x) = \sin(\alpha_0 + \alpha_1 x)$

- Important observation: $y \approx w_0 + w_1 \phi(x)$ remains linear with respect to w_1
- $\phi(x)$ is called **feature transformation**
- To apply feature transformation systematically, we often use some basis functions for φ(x).
 See Feature Transformation.pdf

Nonlinear Regression: Single vs. Multiple Features

• single feature x (feature dimension = 1).



$$y \approx w_0 + w_1 \phi(x)$$

 $\phi(x) = \sin(\alpha_0 + \alpha_1 x)$

• multiple features

 $y \approx \overline{w}_{0} + \overline{w}_{1}\phi(x_{1}) + \overline{w}_{2}\phi(x_{2}) \qquad \phi(x) = (\alpha_{0} + \alpha_{1}x + \alpha_{2}x^{2})$ $y \approx w_{0} + w_{1}x_{1} + w_{2}x_{2} + w_{3}x_{1}^{2} + w_{4}x_{1}x_{2} + w_{5}x_{2}^{2}$

we are now dealing with five-dimensional instead of two-dimensional feature space.

Nonlinear Regression: Feature Explosion

Fun Time: let the input dimension *d* is of moderate size, e.g., d = 100, then the associated polynomial degree n = 5, what would be the resulting features? (1) ~ 10^3 (2) ~ 10^5 (3) ~ 10^7 (4) ~ 10^9



#073374

Nonlinear Regression: Feature Explosion

Fun Time: let the input dimension *d* is of moderate size, e.g., d = 100, then the associated polynomial degree n = 5, what would be the resulting features? (1) ~ 10^3 (2) ~ 10^5 (3) ~ 10^7 (4) ~ 10^9



 The precise dimension of the resulting feature *m* after applying *nth*order polynomial transformation for *d* distinct input features is:

$$m = \frac{(d+n)!}{d!n!} - 1 \qquad m \approx O(d^n)$$

For the input dimension d = 100 and d = 1000, the associated polynomial degree n = 5 would have resulting features of m = 96,560,645 and 8,459,043,543,950, respectively.

Nonlinear Regression: Feature Explosion

- The exponential growth challenge of the feature dimension after feature transformation $\mathbf{x} \rightarrow \phi(\mathbf{x})$ can be cleverly solved by the kernel-based method.
- A kernel is just a form of generalized dot product $\phi(\mathbf{x}) \cdot \phi(\mathbf{z})$
- The key insight in kernel-based method is that you can rewrite many linear models in a way that doesn't require you to ever explicitly compute $\phi(\mathbf{x})$
- The topic is often called kernel tricks, a very important topic in ML but beyond the scope of this course. (refs: Kernel Trick01.pdf, Kernel Trick02.pdf)

Nonlinear Regression: Overfit and Regularization

C-S David Chen, Department of Civil Engineering, National Taiwan University

Regularization

• High-order nonlinear regression can easily lead to overfit.



- Regularization can help!
 - Ridge (L2) regularization
 - Lasso (L1) regularization

Regularization.pdf

Summary: Regression

- Regression is used to predict a continuous value based some given features.
- Feature transformation x → φ(x) allows us to extend linear models to nonlinear models.
- The exponential growth challenge of the feature dimension after feature transformation $x \rightarrow \phi(x)$ can be cleverly solved by the kernel-based method.
- The kernel method is also used in SVM.
- Nonlinear regression tends to overfit and ridge (L2) and lasso (L1) regularization techniques are often used to reduce overfit.
- Many classification methods such as SVM, decision tree, ensemble trees can also handle regression problems.

Cheat Sheet (備忘表)



https://medium.com/dataflair/beat-the-heat-with-machine-learning-cheat-sheet-365c25bd1c3

Midterm Competition Status

https://www.kaggle.com/t/b8d78043a0dc42149d 6bbbedbfaa67da