



Wrap-up quiz 5

This quiz requires some programming to be answered.

Open the dataset `ames_housing_no_missing.csv` with the following command:

```
import pandas as pd

ames_housing = pd.read_csv("../datasets/ames_housing_no_missing.csv")
target_name = "SalePrice"
data = ames_housing.drop(columns=target_name)
target = ames_housing[target_name]
```

`ames_housing` is a pandas dataframe. The column "SalePrice" contains the target variable.

To simplify this exercise, we will only used the numerical features defined below:

```
numerical_features = [
    "LotFrontage", "LotArea", "MasVnrArea", "BsmtFinSF1", "BsmtFinSF2",
    "BsmtUnfSF", "TotalBsmtSF", "1stFlrSF", "2ndFlrSF", "LowQualFinSF",
    "GrLivArea", "BedroomAbvGr", "KitchenAbvGr", "TotRmsAbvGrd", "Fireplaces",
    "GarageCars", "GarageArea", "WoodDeckSF", "OpenPorchSF", "EnclosedPorch",
    "3SsnPorch", "ScreenPorch", "PoolArea", "MiscVal",
]

data_numerical = data[numerical_features]
```

We will compare the generalization performance of a decision tree and a linear regression. For this purpose, we will create two separate predictive models and evaluate them by 10-fold cross-validation.

Thus, use `sklearn.linear_model.LinearRegression` and `sklearn.tree.DecisionTreeRegressor` to create the models. Use the default parameters for the linear regression and set `random_state=0` for the decision tree.

Be aware that a linear model requires to scale numerical features. Please use `sklearn.preprocessing.StandardScaler` so that your linear regression model behaves the same way as the quiz author intended ;)

Question

By comparing the cross-validation test scores for both models fold-to-fold, count the number of times the linear model has a better test score than the decision tree model. Select the range which this number belongs to:

- a) [0, 3]: the linear model is substantially worse than the decision tree
- b) [4, 6]: both models are almost equivalent
- c) [7, 10]: the linear model is substantially better than the decision tree

Select a single answer

Instead of using the default parameters for the decision tree regressor, we will optimize the `max_depth` of the tree. Vary the `max_depth` from 1 level up to 15 levels. Use nested cross-validation to evaluate a grid-search (`sklearn.model_selection.GridSearchCV`). Set `cv=10` for both the inner and outer cross-validations, then answer the questions below.

Question

What is the optimal tree depth for the current problem?

- a) The optimal depth is ranging from 3 to 5
- b) The optimal depth is ranging from 5 to 8
- c) The optimal depth is ranging from 8 to 11
- d) The optimal depth is ranging from 11 to 15

Select a single answer

Now, we want to evaluate the generalization performance of the decision tree while taking into account the fact that we tune the depth for this specific dataset. Use the grid-search as an estimator inside a `cross_validate` to automatically tune the `max_depth` parameter on each cross-validation fold.

Question

A tree with tuned depth

- a) is always worse than the linear models on all CV folds
- b) is often but not always worse than the linear model
- c) is often but not always better than the linear model
- d) is always better than the linear models on all CV folds

Select a single answer

Note: Try to set the `random_state` of the decision tree to different values e.g. `random_state=1` or `random_state=2` and re-run the nested cross-validation to check that your answer is stable enough.

Instead of using only the numerical features you will now use the entire dataset available in the variable `data`.

Create a preprocessor by dealing separately with the numerical and categorical columns. For the sake of simplicity, we will assume the following:

- categorical columns can be selected if they have an `object` data type;
- use an `OrdinalEncoder` to encode the categorical columns;
- numerical columns should correspond to the `numerical_features` as defined above. This is a subset of the features that are not an `object` data type.

In addition, set the `max_depth` of the decision tree to 7 (fixed, no need to tune it with a grid-search).

Evaluate this model using `cross_validate` as in the previous questions.

Question

A tree model trained with both numerical and categorical features

- a) is most often worse than the tree model using only the numerical features
- b) is most often better than the tree model using only the numerical features

Select a single answer

Note: Try to set the `random_state` of the decision tree to different values e.g. `random_state=1` or `random_state=2` and re-run the (this time single) cross-validation to check that your answer is stable enough.

By scikit-learn developers
© Copyright 2022.

[Join the full MOOC for better learning!](#)

Brought to you under a [CC-BY License](#) by [Inria Learning Lab](#), [scikit-learn @ La Fondation Inria](#), [Inria Academy](#), with many thanks to the [scikit-learn](#) community as a whole!