

Kaggle Master Week-2 Q&A

Q1- Which of the following statements are true about the intended use of cross-validation?

- I - To reduce randomness while measuring model performance.
- II - To get a better measure of model performance.
- III - To increase model's training performance.
- IV - To increase MAE (mean absolute error) or MSE (mean squared error).

- I, II, IV
- II, III
- I, II ✓
- All of them

A1- Cross-validation kullanmamızdaki amaç modelimizde kullandığımız metrikleri daha doğru bir şekilde gözlemleyebilmektir. Dolayısı ile modelimizin hatasını düşürmesi veya modeli daha iyi eğitmemiz üzerinde doğrudan bir etkisi yoktur.

Q2- Which of the following statements are true about LabelEncoder and OneHotEncoder?

- I-They help us to deal with categorical values.
- II-Label Encoding assigns each value to a different integer whether it is unique or not.
- III-One Hot Encoding creates new column for every possible value in the original data.
- IV-For large number of categorical variable count value (such as 15 different values) it is not good to use One Hot Encoder generally.

- I, II, IV
- I, III, IV ✓
- I, II, III
- All of them

A2- Label Encoder ve One Hot Encoding kategorik verilerin üstesinden gelmek için kullanılırlar. Label Encoding her eşsiz (unique) değer için bir değer üretip atama yapar. One Hot Encoding ise her değer için yeni bir kolon oluşturur. Bu değerler eşsiz (unique) değerlerdir. Kolon sayısı arttıkça, genelde One Hot Encoding iyi bir performans sergilemez. Bu yüzden One Hot Encoding genelde fazla sayıdaki unique kolon içeren kategorik verilerle kullanıldığında iyi sonuç vermez.

Q3- Which of the following statement is inconsistent with pipelines?

- With pipelines, there is less probability to forget a preprocessing step.
 - It's hard to productionize a model with pipelines. ✓
 - You won't need to manually keep track of your training and validation data at each step with a pipeline.
- With a pipeline, we can use the cross-validation technique easily.

A3- Pipelinesi, modelimize input olarak verilecek datanın her zaman aynı işlemlerden geçirilmesi, ön-işlemede meydana gelebilecek hata ve eksiklik risklerinin azaltılması ve cross-validation gibi model değerlendirmesi yaptığımız işlemleri kolayca yapabilmek için kullanıyoruz. Modelleri pipeline ile oluşturmak zor değil ve model deployment aşamasında hata yapmanızı büyük ölçüde engelleyeceğinden dolayı oldukça kullanışlılar.

Q4- print(df(head).method())

Assume that you want to print locations of the missing values in the top 10 rows. Which method is suitable for this?

- dropna(how='any')
- isnan
- notnull
- isnull ✓

A4- Bir veri çerçevesinde NULL değerlerini denetlemek ve yönetmek için isnull () ve notnull () yöntemleri kullanılır. isnull () yöntemi, NaN değeri için True ve boş olmayan değer için False döndürür. notnull() bu durumun tam tersidir.

Q5- Which of the following is not a Booster parameter of XGBoost?

- min_child_weight
- objective ✓
- max_leaf_nodes
- colsample_bylevel

A5- “objective” parametresi bir learning task parametresidir. Bunun gibi parametreler, her adımda hesaplanacak

metriğin optimizasyon hedefini tanımlamak için kullanılır.

Q6- What do the highlighted code pieces mean?

```
x_train_plus = x_train.copy()
x_valid_plus = x_valid.copy()
for col in cols_with_missing:
    x_train_plus[col + '_was_missing'] = x_train_plus[col].isnull()
    x_valid_plus[col + '_was_missing'] = x_valid_plus[col].isnull()
my_imputer = SimpleImputer()
imputed_x_train_plus = pd.DataFrame(my_imputer.fit_transform(x_train_plus))
imputed_x_valid_plus = pd.DataFrame(my_imputer.transform(x_valid_plus))
imputed_x_train_plus.columns = x_train_plus.columns
imputed_x_valid_plus.columns = x_valid_plus.columns
```

- To make new columns indicating what will be imputed
- For imputation
- To make copy to avoid changing original data
- To put removed column names back ✓

A6- Yukarıdaki code parçası kayıp verileri işlemede kullanılan bir yöntem olan Imputation adımlarını ifade etmektedir. İşaretli satırlar da, imputing işlemi sırasında kayıp verileri çıkarılmış kolonları, temizlenmiş olarak geri almamızı sağlar.

Q7- Which of the below is/are nominal variable(s)?

- I - Gender
- II - Genotype
- III - Religious preference
- IV- IQ
- V - Income earned in a week.

- I, II
- I, II, III ✓
- II, III, IV
- All of them

A7- Nominal değişkenler aralarında sıralama yapılamayan kategorik değişkenlerdir. Cinsiyet, genotip ve dini tercihler değerlerinin birbirlerine herhangi bir üstünlüğü bulunmayan değişkenlerdir. IQ ve haftalık kazanç kategorik değişkenler olmadığından nominal değişken olarak değerlendirilemezler.

Q8- Which of the following statements are true about “max_depth” hyperparameter in Random Forest?

- I- Lower is better parameter in case of same validation accuracy
- II- Higher is better parameter in case of same validation accuracy
- III- Increase the value of max_depth may overfit the data
- IV- Increase the value of max_depth may underfit the data

- I, IV
- II, IV
- I, III ✓
- II, III

A8- Çünkü maksimum derinliği gereğinden fazla artırmamız modelimizin veriyi ezberlemesine ve overfit olmasına yol açar. Farklı derinlikler ile oluşturduğumuz modellerden aynı skoru alırsak modelimiz karmaşıklığını azaltmak için düşük derinlikli olanı tercih etmemiz gerekir.

Q9- You will build a model to predict housing prices. The model will be deployed on an ongoing basis, to predict the price of a new house when a description is added to a website. Here are four features that could be used as predictors. Which of the features is most likely to be a source of leakage?

- Size of the house (in square meters)
- Average sales price of homes in the same neighborhood ✓
- Latitude and longitude of the house
- Whether the house has a basement

A9- Data leakage (veri sızıntısı), eğitim verileri hedef hakkında bilgi içerdiğinde gerçekleşir, ancak model tahmini için kullanıldığında benzer veriler kullanılamaz. Bu, eğitim setinde (ve hatta muhtemelen doğrulama verilerinde) yüksek performansa yol açar, ancak model üretimde kötü performans gösterecektir.

Başka bir deyişle, karar verme mekanizması başlayana kadar o model çok doğru görünür fakat en sonunda modelin çok yanlış kurulduğu ortaya çıkar.

- 1- Bir evin büyüklüğünün satıldıktan sonra değiştirilmesi olası değildir (teknik olarak mümkün olsa da). Ancak tipik olarak bu bir tahmin yapmamız gerektiğinde kullanılabilir ve veriler ev satıldıktan sonra değiştirilmez. Bu yüzden oldukça güvenlidir.
- 2- Bunun ne zaman güncellendiğini bilmiyoruz. Bir ev satıldıktan sonra ham verilerde alan güncellenirse ve ortalamanın hesaplanması için evin satışı kullanılırsa, bu veri sızıntısı anlamına gelir. Bir uçta, mahallede sadece bir ev satılıyorsa ve tahmin etmeye çalıştığımız ev ise, o zaman ortalama tahmin etmeye çalıştığımız değere tam olarak eşit olacaktır. Genel olarak, az satış yapılan mahalleler için model, eğitim verileri üzerinde çok iyi performans gösterecektir. Ancak modeli uyguladığınızda, tahmin ettiğiniz ev henüz satılmayacaktır, bu nedenle bu özellik eğitim verilerinde olduğu gibi çalışmaz.
- 3- Bunlar değişmez ve bir tahmin yapmak istediğimiz zaman hazır olur. Yani burada veri sızıntı riski yoktur.
- 4- Bu da değişmez ve bir tahmin yapmak istediğimiz anda kullanılabilir. Yani burada veri sızıntı riski yoktur.

Q10- How is the Gradient Boosting cycle proceed? Please choose the correct order from the mixed statements below.

- I- We add the new model to ensemble.
- II- We use the current ensemble to generate predictions for each observation in the dataset.
- III- We use the loss function to fit a new model that will be added to the ensemble.

- I-II-III
- I-III-II
- II-I-III
- II-III-I ✓

A10- Gradient boosting döngüsünde ilk olarak, veri grubundaki her bir gözlem için tahminler oluşturmak üzere mevcut topluluğu (ensemble) kullanıyoruz. Bir tahmin yapmak için, topluluktaki tüm modellerden tahminleri ekliyoruz. Bu tahminler bir kayıp fonksiyonunu (loss function) hesaplamak için kullanılır.



Daha sonra loss function ı, topluluğa (ensemble) eklenecek yeni bir modele uyacak şekilde kullanıyoruz. Özellikle, model parametrelerini belirlemede kullanıyoruz ki böylece bu yeni modeli topluluğa eklemekle olası zaman kayıplarını azaltıyoruz. Son olarak da topluluğa yeni modeli ekliyoruz.