



# **HOUSE PRICE DATASET'İ**

## **SEZGİSEL ANALİZ RAPORU**

### **GRUP 1**

Batuhan Özbay

Melisa Gündüz

Aybüke Akçay

Çiğdem Taş

Ali Şenyurt

Ercan Tuncay

Fuat Akdemir

Uğur Can Kıvanç

İzel Çelikkaya

Yiğit Hakverdi

Serenay Ardahanlı

Mehmet Özmen

Kodluyoruz ve İstanbul Büyükşehir Belediyesi, İstanbul Veri Bilimi Bootcamp kapsamında House Price Dataset üzerinde Sale Price değerinin tahminlenmesi için bir makine öğrenmesi modeli oluşturulması planlanmaktadır. Bu kapsamda, veri setinin sezgisel analizi gerçekleştirilmiştir.

Veri seti incelenirken ilk olarak kayıp değer olup olmadığı kontrol edilmiştir. Kategorik değişkenlerdeki kayıp olarak görülen değerler aslında bir kayıp değer değil, mevcut kategorik değişkenin özelliklerinin taşınıp taşınmadığını göstermektedir. Örneğin, PoolQC değişkenindeki kayıp değer olarak görülen değerler aslında o ev için havuzun bulunmadığını dolayısıyla da bir kalite olmayacağı anlamına gelmektedir.

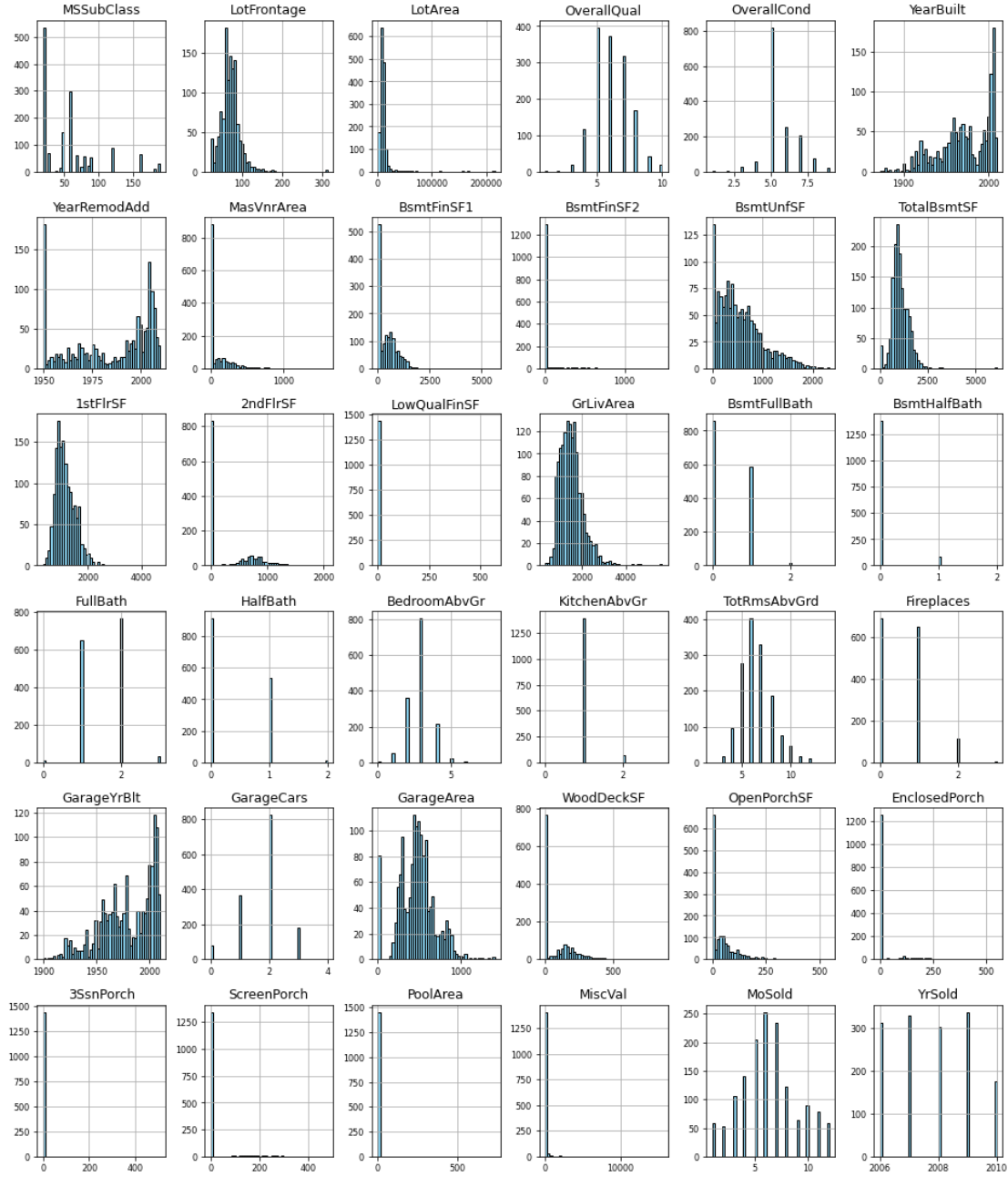
Aşağıda verilen değişkenlerde kayıp değer olarak görülen değerler kayıp değer değildir. Bu değerler 'None' olarak atanarak kayıp değerler yok edilebilir.

PoolQC	1453
MiscFeature	1406
Alley	1369
Fence	1179
FireplaceQu	690
GarageType	81
GarageCond	81
GarageQual	81
GarageFinish	81
GarageYrBlt	81
BsmtFinType2	38
BsmtExposure	38
BsmtFinType1	37
BsmtQual	37
BsmtCond	37
MasVnrArea	8
LotFrontage	259

Aşağıda verilen değerlerin karşısında kayıp değer sayısı verilmiştir.

MasVnrType	8
Electrical	1

Numeric değerlerin histogramları, buradan değerlerin discrete valuelarının dağılımına göre hangi kolonların kendi içinde baskın davrandığı ve düşürülmesi gerektiği görülebilir. Örnek olarak 'LowQualFinSF' histogramına bakılabilir.



Veri seti detaylı bir şekilde incelendiğinde aşağıdaki bulgulara ulaşılmıştır:

MSSubClass değişkeninin türü numerik olarak gösterilmiştir ancak değişken kategorik değişkendir. Evler yapıldıkları tarihlere göre gruplandırılmışlardır. MSSubClass değişkeninin değeri arttıkça evin yapıldığı tarih de günümüze yaklaşmıştır. Herhangi bir encoding gerekmeden numerik hali ile modelde kullanılabilir.

Bazı değişkenler çok yakın değerlere sahip olduğundan tek değer olarak ele alınabilir. Birlikte kullanılmak üzere birleştirilecek değerler seçilirken yüksek korelasyona sahip olup olmadıkları incelenecektir.

Overall Quality yükseldikçe ev fiyatları da yükselmiştir. Bu durum Overall Condition değişkeni için geçerli değildir.

Yearbuilt ve Yearremodeled değişkenleri SalePrice değişkeninde benzer bir artışa sebep olmuştur. 2000 ve sonrası için değerler yüksektir. 2000-2006 yılları arasında evlerde daha çok tadilat yapılmıştır. O dönem evlerde tadilat için bir talep olduğu söylenebilir.

Evler %99.0 ile Condition2 durumunda ve evlerin %49.7'si tek katlı; %30.5'i ise 2 katlı. Evlerin %83.6'sı tek ailelik evler / sonra %7.8 ile ikiz evler geliyor

Year remodeled — neighborhood değişkenleri birbirleri ile yüksek oranda ilişkilidir. İyi mahallelerde insanların evlerine tadilat yaptırmayı tercih ettiklerini söyleyebiliriz.

Saleprice'ın median değeri 163.000'dir.

KitchenQual değişkeni OvrQual ve Saleprice ile korelidir. Model içinde kullanılmaya uygun olduğu söylenebilir. Benzer şekilde TotRmsAbvGrd değişkeni BedroomAbvGrd ve Saleprice ile korelidir. Değişken unimodal bir dağılıma sahiptir. Model kurulurken kullanılması tercih edilebilir.

Functional değişkeninin %93'lük kısım typical ve herhangi bir değişkenle korele değil, kullanılması tercih edilmeyebilir. Ancak hasarlı evlerin fiyatlarının düşük olabileceği hipotezinden yola çıkılarak kullanılması tercih edilebilir.

FirePlaces değişkeninin modelde kullanılmaya uygun bir dağılımı vardır. Aynı zamanda SalePrice ile korele olduğu için kullanılabilir.

FireplaceQu: Missing value olmayan değerler kullanılabilir.

GarageType değişkeni Garagefinish değişkeni ile korelidir. Ancak kullanılması hakkında soru işaretleri mevcuttur.

GarageYrBlt değişkeni SalePrice'la koreledir. Bu değişken GarageCars ve GarageArea değişkenleri ile de koreledir. Kullanılması tercih edilebilir. Kullanıldığı takdirde yüksek korelasyona sahip olduğu diğer değişkenlerle birleştirilerek kullanılabilir.

GarageFinish değişkeni kategorik değişken olduğu için satış fiyatıyla olan korelasyonu değerlendirilememiştir. Ayrıca kategorideki 'No Garage' seçeneği missing value olarak işlenmiştir, bu hata düzeltilerek modellemeye ekleme yapılabilir.

GarageCars değişkeni nümerik değişken tipinde tanımlanmıştır. Satış fiyatıyla yüksek korelasyona sahiptir.

GarageArea değişkeni SalePrice ile yüksek korelasyona sahip nümerik değişkendir.

GarageQual ve GarageCond değişken tipleri çok yüksek korelasyona sahiptir bu sebeple beraber tek bir değişken olarak incelenebilirler.

PavedDrive, Paved değişkeninde yığılma olan kategorik değişken türü. Drop edilebilir.

Alley değişkeninde NA olarak adlandırılmış bir kategori de mevcuttur. Bu nedenle, eksik değer bulunmamaktadır. Değişken NA dışında iki tür kategori bilgisi içermektedir: paved ve gravel(çakıl).

LotFrontage değişkeni "Linear feet of street connected to property" olarak tanımlanmıştır. Bu tanımdan yola çıkılarak LotFrontage eksik değerleri bu tanıma göre LotFrontage eksik değerlerin bulunması mümkün değil gibi görünüyor

WoodDeckSF (Wood deck area in square feet) değişkeninin Phik korelasyonu sadece "PoolQC" ile bulunmaktadır. Wood Deck Area olmayan evler 0 değeri almıştır. 761 adet 0 değeri bulunmaktadır

OpenPorchSF: (Open porch area in square feet) değişkeninin Phik korelasyonu Condition2:0.51 BsmtFinSF1:0.59 GrLivArea:0.56 değerleri almaktadır. Open Porch Area bulunmayan evler 0 değeri almamıştır. (656 adet)

EnclosedPorch (Enclosed porch area in square feet) değişkeninin Phik korelasyonu PoolArea:0.51.Enclosed porch area olmayan evler 0 değeri almıştır. 1252 adet 0 değeri vardır.

PoolArea (Pool area in square feet) değişkeninin Phik korelasyon GrLivArea: 0.76, MiscFeature: 1.00. 1453 adet evde havuz yoktur.

PoolQC değişkenine ait bilgiler şu şekilde özetlenebilir: Pool quality Ex Excellent= 2 Gd Good :3 TA Average/Typical:0 Fa Fair : 2 adet NA No Pool: 1453 Null(NA) değeri yani havuz olmayan ev sayısı 1453 adettir. NA değerler eksik veri değil havuz olmamasıdır.

PoolQC Phik korelasyon deęerleri : LotConfig: 0.81 Neighborhood: 1.00 HouseStyle: 0.80  
YearBuilt: 0.63 YearRemodAdd: 0.83 MasVnrType: 0.92 ExterQual: 0.81 BsmtExposure:  
0.72 BsmtFinSF1: 0.63 BsmtFinSF2: 0.65 BsmtUnfSF: 0.82 GrLivArea: 0.83  
BedroomAbvGr: 0.81 FireplaceQu: 0.67 GarageYrBlt: 0.83 GarageFinish: 0.87 GarageArea:  
0.68 WoodDeckSF:0.67 SaleCondition: 0.68 SalePrice: 0.63

ExterQual (Dış cephedeki malzemenin kalitesini deęerlendirir): 17 farklı başlıkla yüksek korelasyonlu (ExterCond dahil deęil). %62.1'i Tipik, %33.4'ü İyi seviyede.

ExterCond (Malzemenin dış cephedeki mevcut durumunu deęerlendirir): OverallCond ve MiscFeature ile yüksek korelasyonlu. %87.8'i Tipik, %10'u İyi seviyede.

Foundation (Yapı Temeli: Yapı temeli türü): 12 alanla yüksek korelasyonlu. %44.3'ü Dökme Beton, %43.4'ü Cinder Bloęu (bir çeşit tuęla), %10'u ise Tuęla, Fayans'tan oluşmaktadır.

BsmtQual (Bodrumun yüksekliğini deęerlendirir): 16 alanla yüksek korelasyonlu. %44.5'i Tipik, %42.3'ü İyi, %8.3'ü Çok İyi seviyede. NA'ları (No basement) missing olarak görölmektedir. Ancak missing deęer deęildir.

BsmtCond (Bodrumun genel durumunu deęerlendirir): OverallQual ve OverallCond ile yüksek korelasyonlu. %89.8'i Tipik. NA'ları (No basement) missing olarak görölmektedir. Ancak missing deęer deęildir.

BsmtExposure (Walkout veya bahçe seviyesindeki duvarları ifade eder): PoolQC ile yüksek korelasyonlu. %65.3'ü No Exposure, %15.1'i Average Exposure, %9.2'si Good Exposure'dur. NA'ları (No basement) missing olarak görölmektedir. Ancak missing deęer deęildir.

BsmtFinType1 (Bodrum bitmiş alanın deęerlendirmesi): 7 alanla korelasyonu var. %29.5'i Unfinished, %28.6'sı Good Living Quarters, %15.1'i Average Living Quarters, %10.1'i Below Average Living Quarters, %9.1'i Average Rec Room'dan oluşmakta. NA'ları (No basement) missing olarak görölmektedir. Ancak Missing deęildir.

Evin içerisinde bulunan bodrum ile ilgili birden fazla deęişken mevcuttur. Bu deęişkenlerden SalePrice ile en yüksek korelasyona sahip olanının modelde kullanılması yeterli olacaktır. Modelle ilgili yapılacak geliştirmelerde birden fazla bsmt deęişkeni birleştirilerek modele eklenebilir.

HeatingQC, CentralAir, Electrical deęişkenleri kategorik deęişkenlerdir.

Electrical kolonunda sadece 1 adet null deęer bulunmaktadır. Bu null deęerin etkisi olmadığından en çok bulunan deęere atanmasında bir problem oluşmayacaktır.

Ancak CentralAir ve Electrical kolonlarındaki deęişkenler histogramda incelendiğinde ayrıık deęişkenlerin dağılımının %90 'nından fazlasının ayrıık deęişkenlerden sadece birine kaymış

olması, yani diğer değerler yüksek oranda baskılanması, bu kolonların SalesPrice üzerinde etkisinin çok az olacağını gösterir ve bu kolonların drop edilmesi gerekmektedir.

CentralAir Bool değerden oluşan, %93.5 'i 'True' olan bir kolon, SalePrice'a etkisi düşük bir değişkendir.

Electrical 5 ayırık değerden oluşan, %91.4 'ü 'SBrkr' olan bir kolon, SalePrice'a etkisi düşük bir değişkendir.

HeatingQC değişkeni 5 ayırık değerden oluşuyor. Ancak HeatingQC böyle bir durum olmadığında labeling yapılarak devam edilebilir.

1stFlrSF, 2ndFlrSF, LowQualFinSF, GrLivArea değişkenleri numerik değişkenlerdir.

1stFlrSF, 2ndFlrSF, GrLivArea kolonlarının SalesPrice ile yüksek korelasyon ilişkisi bulunan kolonlardır.

Ancak LowQualFinSF ise SalesPrice le korelasyon ilişkisi yoktur. Histogramında ise %98.2 oranında '0' değerinde bulunmakta ve bu yüzden drop edilmesi gerekmekte.

Full bathroom: "lavabo, tuvalet ve küvet veya duş içeren banyo"

Half bathroom: "lavabo ve tuvalet içeren ancak küvet veya duş içermeyen banyo" şeklinde tanımlanmaktadır.

BsmtFullBath değişkeni incelendiğinde, bodrum katındaki full banyo sayısı %58.6'sında 0, %40.3'ünde 1, kalanlarda 2 veya 3 olduğu görülür. BsmtFullBath SalePrice ile Pearsons correlation: 0.23'tür.

BsmtHalfBath değişkeni incelendiğinde, bodrum katındaki yarım banyo sayısı %94.4'ünde 0, kalanlarda 1 veya 2 olduğu görülür. BsmtHalfBath SalePrice ile Pearsons correlation: -0.02'dir.

FullBath değişkeni incelendiğinde, bodrum katı haricindeki full banyo sayısı %52.6'sında 2, %44.5'inde 1, kalanlarda 0 veya 3 olduğu görülür. FullBath değişkeni SalePrice ile yüksek korelasyona sahiptir. Pearsons korelasyonları 0.56'dır.

HalfBath değişkeni incelendiğinde, bodrum katı haricindeki yarım banyo sayısı %62.5'inde 0, %36.6'sında 1, kalanlarda 2 olduğu görülür. HalfBath SalePrice ile Pearsons correlation: 0.28'dir.

BedroomAbvGr değişkeni incelendiğinde Bodrum katı haricindeki yatak odası sayısı %55.1'inde 3, %24.5'inde 2, kalanlarda 0, 1, 4, 5, 6, 7, 8 değerleri görülür. BedroomAbvGr SalePrice ile Pearsons correlation: 0.17'dir.

KitchenAbvGr deęiřkeni incelendięinde bodrum katı haricindeki Mutfak sayısı % 95.3'ünde 1, kalanlarda 0, 2, 3 olduęu g r l r. KitchenAbvGr SalePrice ile Pearsons correlation: -0.14't r.

Fence degiskeni %80'nin ustunde N/A'dir, bu da 'No Fence' anlamina geliyor. Verilerin cogunlugunda N/A oldugu icin bu degiskten modelde tercih edilmemesi daha dogru olabilir.

MiscFeature, her bir verinin  eřitli  zelliklerini i eriyor; bunlar, asans r, ikinci bir garaj var mi veya tenis kortu var mi řeklinde ayrılmıřtır.

MiscVal, her bir miscellaneous deęerini i eren  zellik; bu  zellięin kullanılması  ok dogru olmayabilir   nk  bu  zellikte  ok fazla 0 deęer var.

MoSold, Hangi ay satildięi

YrSold, Hangi yıl satildięi

SaleType, Bu feature satisin tipini iceren veriler bulunuyor. (Kategorik) Verinin dagilimina bakildiginda en fazla Warranty Deed bulunmaktadır tahmin yapilmadan once veriler bu feature gore ayiklanip bunun WD uzerinden yapılabilir

SaleCondition Ayni sekilde satisin kosulu, en fazla Normal deęer icermektedir (Kategorik degisken) SaleType ile yuksek korelasyona sahiptir. Veriler preprocessing isleminde gecirilmedne once bu iki degisken goz onune alinarak bir ayiklama yapılabilir ve model bu iki kategorik deęerin uzerine kurulabilir

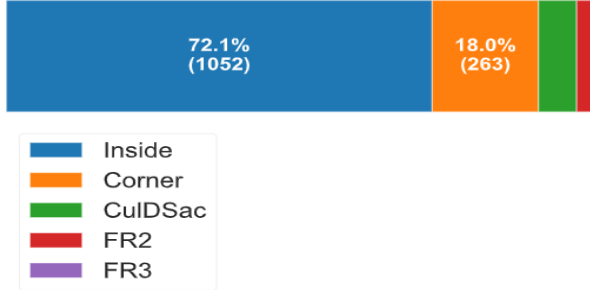
Landcontour: Binanın bulunduęu arazinin d zl ę n  ifade etmektedir. Missing value bulundurmamaktadır. Neighbourhood ile y ksek korelasyon i ermektedir. Landcontour i in d rt derece atanmıřtır lvl(zemin seviyesi) binaların y zde 89.8'inin toplandıęı alandır. Dięer dereceler banked,hillsed ve low'dan oluřmaktadır.Missing value yoktur.

Utilities: Binada ulařılabilir kamu hizmetlerini ifade etmektedir. Binaların %99.9'u t m kamu hizmetlerine ulařılabilmektedir. Y ksek korelasyonu olan hi bir deęiřken bulunmamaktadır. (Bence bu veriyi g z  n nde bulundurmamalıyız.)



### LotConfig: Parsel yapısı

Mülkün cephelerini tanımlayan bir değişkendir. Bu değişkenler inside,corner, cu-de-sac(çıkamaz), FR2 ( mülkün iki tarafı cephe ve FR3( mülkün üç tarafı cephe)den oluşmaktadır. LotConfig PoolQc ile yüksek korelasyon göstermektedir.Missing value yoktur



### LandSlope: Arazinin eğimi

LotArea, Roofstyle ve Neighbourhood ile yüksek korelasyon içindedir. (belki bu değişkenlerle birlikte incelenebilir) Gentle, moderate ve severe olmak üzere üç dereceye ayrılmıştır. Binaların büyük bir kısmı (%94.7) Gentle slope'tur. Missing value yoktur.

Neighbourhood: 36 değişkenle birden yüksek korelasyon içerisindedir. (Saleprice da buna dahildir.) Missing value yoktur.