

# Regresyonlar

Fuat Can Beylunioğlu

December 21, 2017

# Giriş

- ▶ Regresyon bir bağımlı değişkenin  $Y_i$  ile bağımsız değişken  $X_i$  ya da değişkenlerle  $X$  arasındaki ilişkiyi ölçer.
- ▶ Bu ilişki tek yönlüdür, yani  $Y_i$ 'nin  $X_i$ 'deki değişikliklerden nasıl etkilendiği üzerine bir modeldir.
- ▶ Model lineerdir, yani ilişkiyi hata terimlerini dışarıda bırakarak doğrusal olarak modellemeyi hedefler.
- ▶ Temel regresyon denklemi şu şekildedir:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad (1)$$

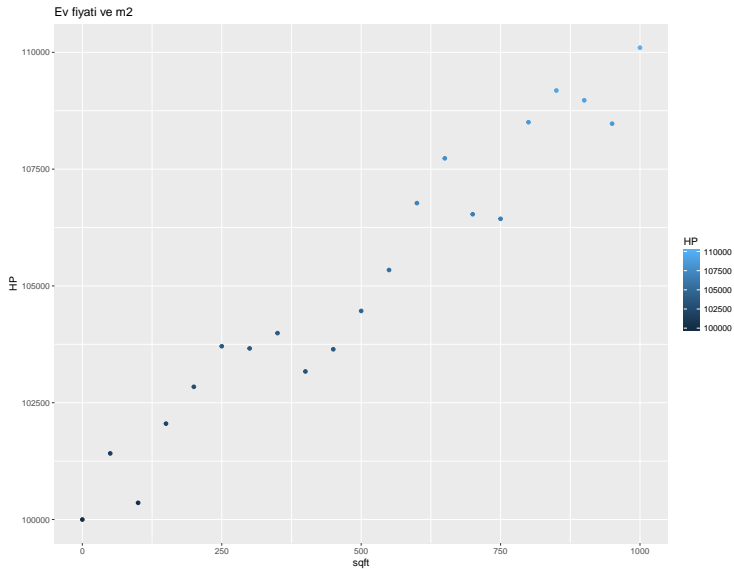
- ▶ Burada  $\epsilon_i$ 'nin tamamen rassal olduğu varsayılmaktadır. Matematiksel olarak ifade etmek gerekirse  $\epsilon_i \sim N(0, \sigma^2)$  olmalıdır. Aksi takdirde denklem taraflı (biased) olacaktır.

# Regresyon: Ev Fiyatı Örneği

- ▶ Örneğin ev fiyatının, evin  $m^2$  genişliği ile ilişkili olduğunu varsayalım. Bu durumda evin teorik değeri aşağıdaki gibidir:

$$HP_i = \beta_0 + \beta_1 \text{sqft}_i + \epsilon_i \quad (2)$$

$$\beta_0 = 100000, \beta_1 = 1000$$



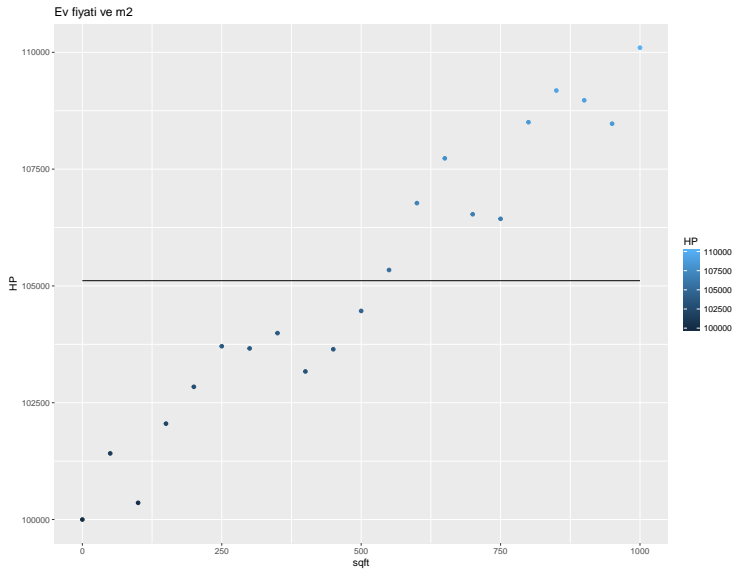
## Regresyon: Ev Fiyatı Örneği

- ▶ Örneğin ev fiyatını, evin  $m^2$  genişliği ile tahmin eden aşağıdaki modeli ele alalım:

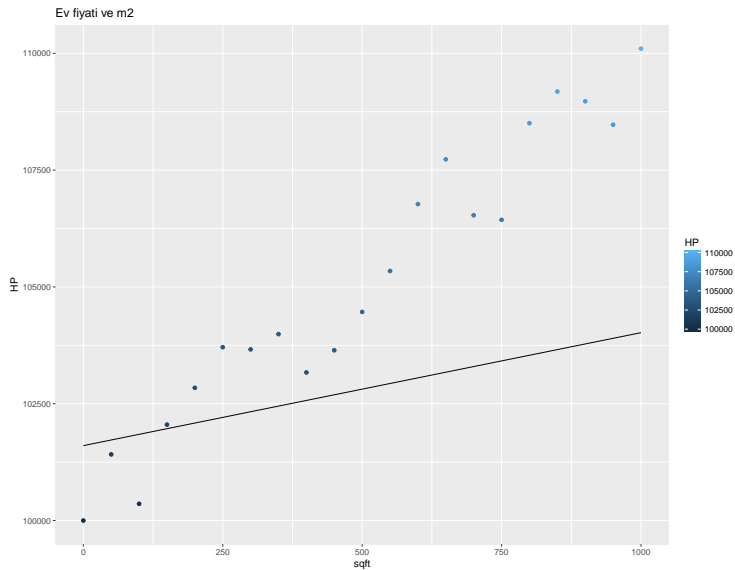
$$HP_i = \beta_0 + \beta_1 sqft_i \quad (3)$$

- ▶  $\beta_0 = 100000$  ve  $\beta_1 = 1000$  olduğu durumda, modele göre evin baz değerinin 100000 olduğu ve her bir  $m^2$ 'lik artışın evin değerine 1000TL değer eklediği söylenebilir.
- ▶ Yukarıdaki model regresyon çizgisinin formülüdür,  $\epsilon_i$  terimi yazılmaz.

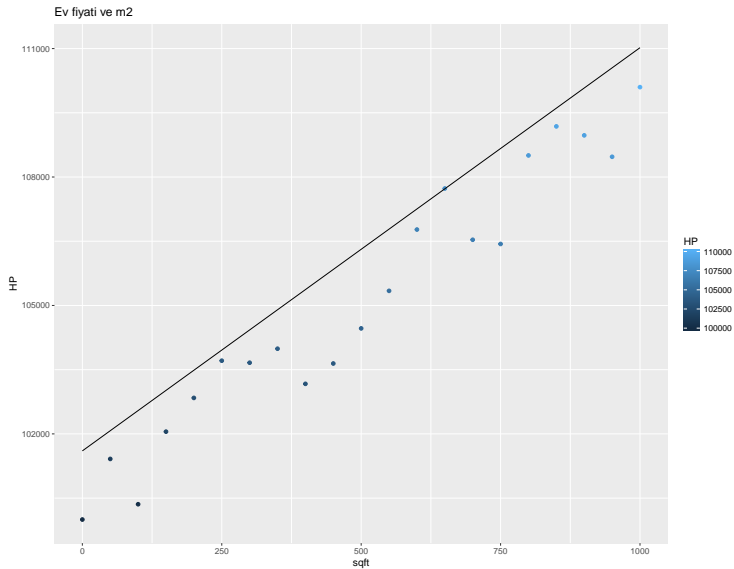
# Hangi Model



# Hangi Model

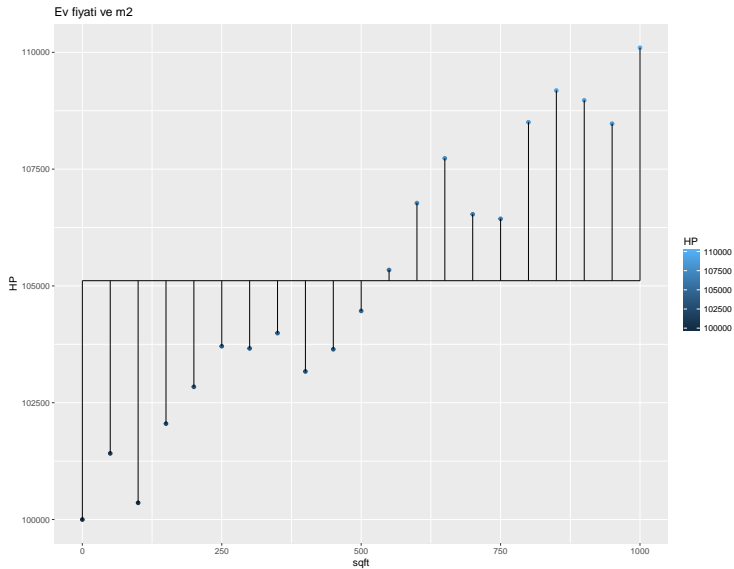


# Hangi Model

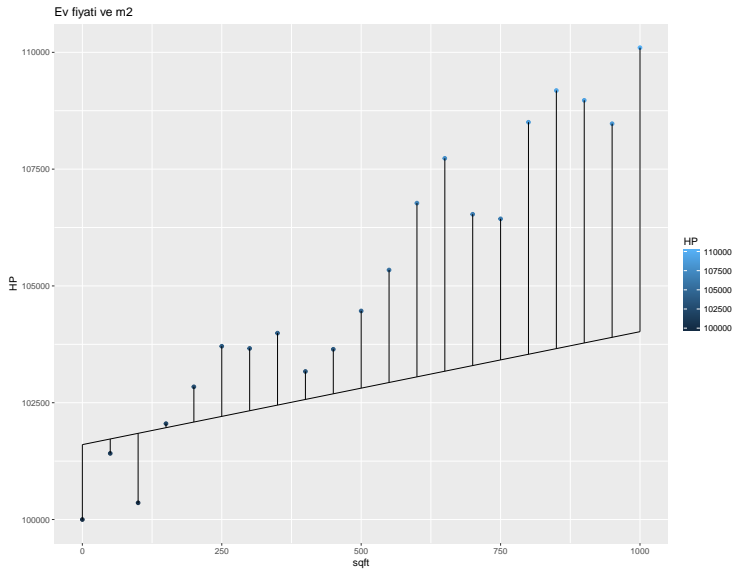




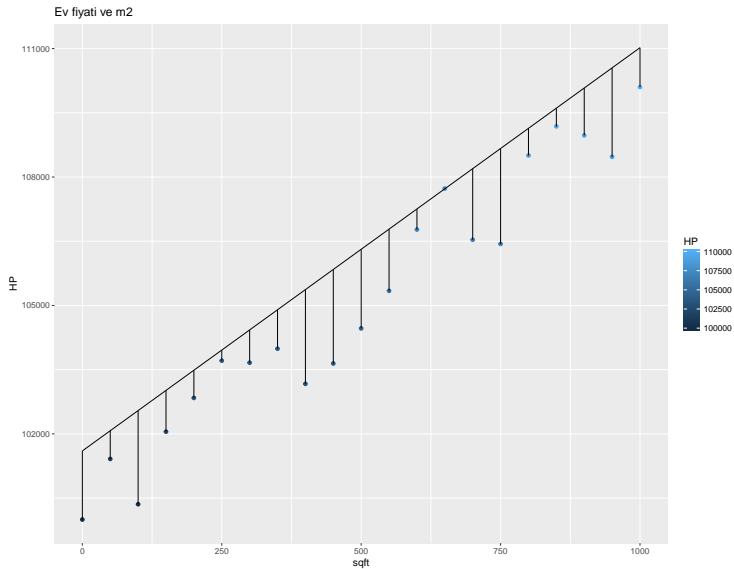
# Hangi Model



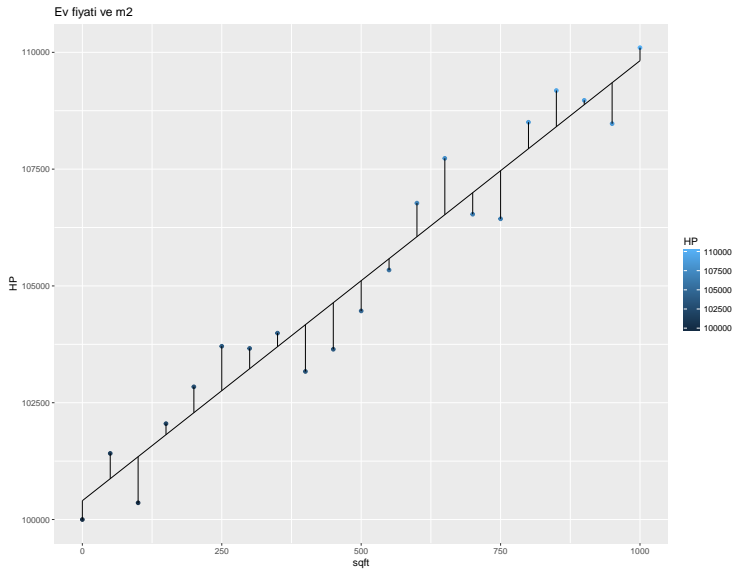
# Hangi Model



# Hangi Model



# Hangi Model



# En İyi Modeli Oluşturmak

- ▶ En iyi model hata terimlerinin karelerin toplamının (RSS) hesaplanması ile oluşturulur:

$$RSS = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2 \quad (4)$$

- ▶ RSS'i minimize eden parametre kümesi ( $\beta_i$ ) matematiksel olarak en iyi model(best fit)dir.
- ▶ Ancak örneklem kümesi ile elde edilen modelin popülasyon için de geçerli olabilmesi için belli varsayımları sağlaması gerekir.

# Çoklu Regresyon (Multiple Regression)

- ▶ Çoklu regresyon aşağıdaki formülle gösterilir:

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \cdots + \beta_n X_{nt} + \epsilon_t \quad (5)$$

- ▶ Ayrıca:

$$Y_t = \beta_0 + \sum_{i=1}^N \beta_i X_{it} + \epsilon_t \quad (6)$$

- ▶ ya da matrix formunda:

$$Y = X\beta + \epsilon \quad (7)$$

olarak gösterilir.

# Regresyonun Varsayımları

- ▶ Normality:  $\epsilon_i$ 'lerin dağılımı  $N(0, \sigma^2)$  olmalı,
- ▶ Homoskedasticity: Her bir  $x_i$  için  $\sigma^2$  sabit olmalı,
- ▶ No serial autocorrelation: Her bir  $\epsilon_i$ 'nin bir diğeri ile arasında korelasyon olmamalı.
- ▶ No multicollinearity: Çoklu regresyon için  $X_{ij}$  ve  $X_{ik}$  arası korelasyon olmamalı.

# Modelin Açıklayıcılığı

- ▶ Bir modelin kuvvetini ölçen başlıca istatistikler:
  - ▶  $R^2$  ve Adjusted  $R^2$
  - ▶ t-test, F-test
  - ▶ AIC, BIC



# Modelin Açıklayıcılığı, $R^2$

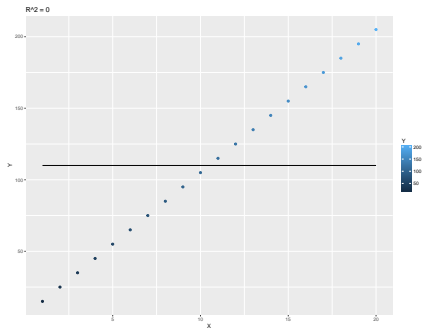


Figure:  $R^2 = 0$

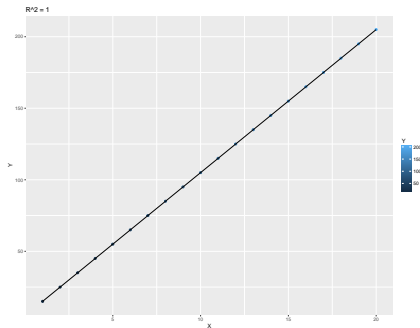
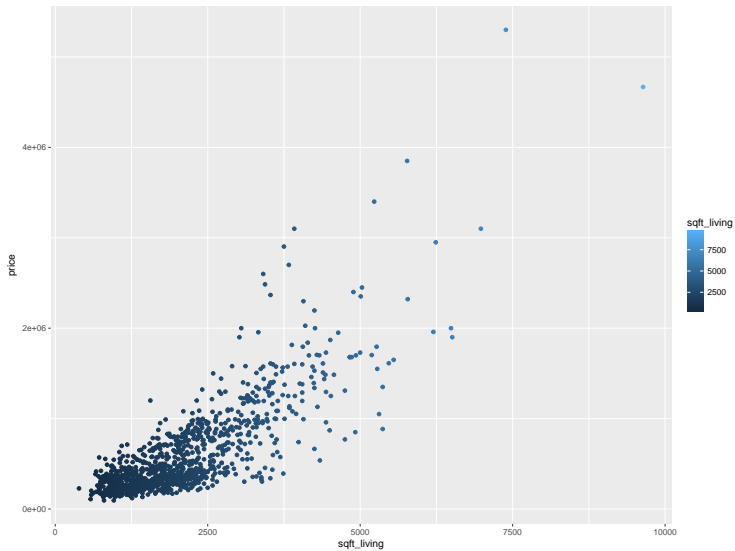


Figure:  $R^2 = 1$

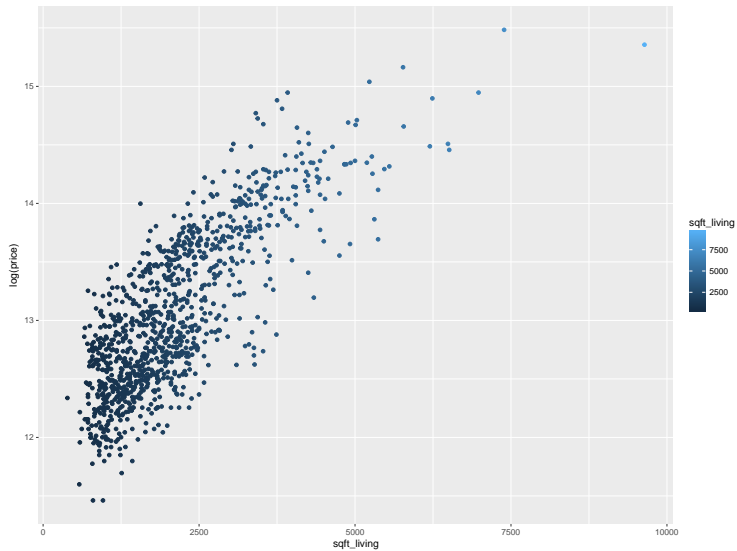
# Lineer Olmayan İlişkiler

- ▶ Bazı durumlarda iki değer arasında lineer dışı ilişki söz konusu olur.
- ▶ Bu regresyon varsayımlarının sağlanamamasına neden olacaktır.
- ▶ Analizden önce iki değer arasında öncelikle scatter plot yapılması bu açıdan önemlidir.
- ▶ Bu gibi durumlarda değerlerden biri ya da ikisinin log değeri alınabilir. Log değerler veriler arasındaki açıklıkları düşürerek normal dağılıma yakınsamalarını sağlayabilir, öte yandan heteroskedasticity sorununu da çözebilir.
- ▶ Sorun çözülmediği takdirde polinom modellere başvurulabilir.

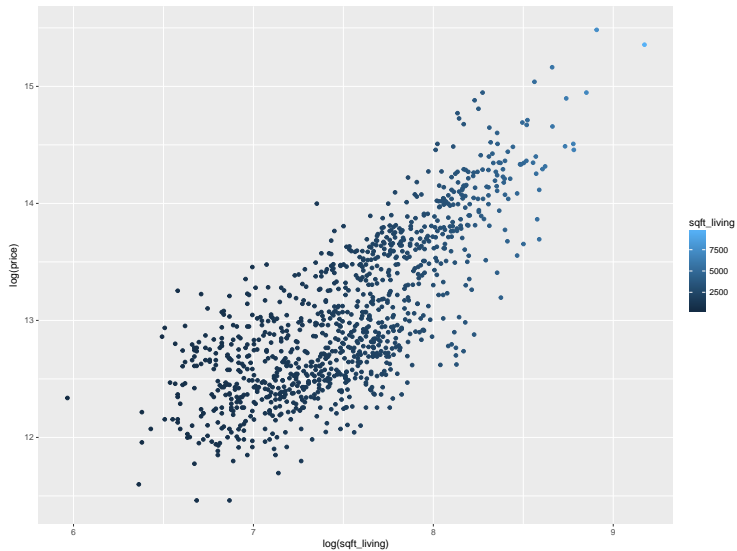
hp vs sqft, Washington



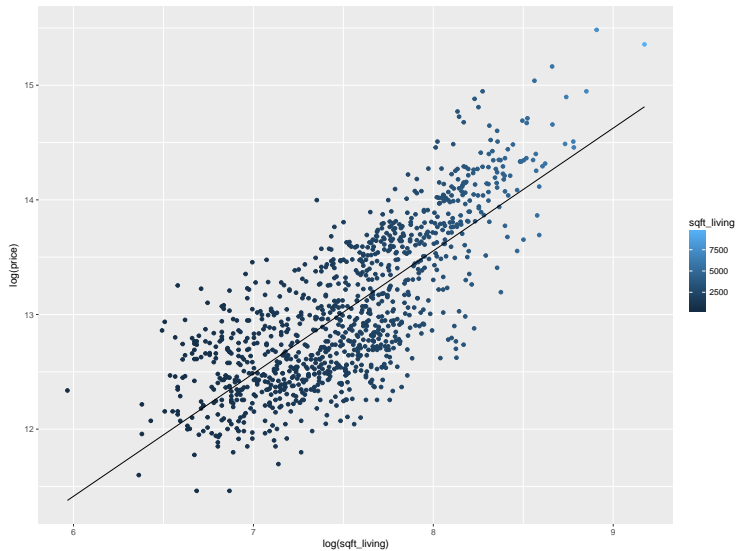
log(hp) vs sqft, Washington



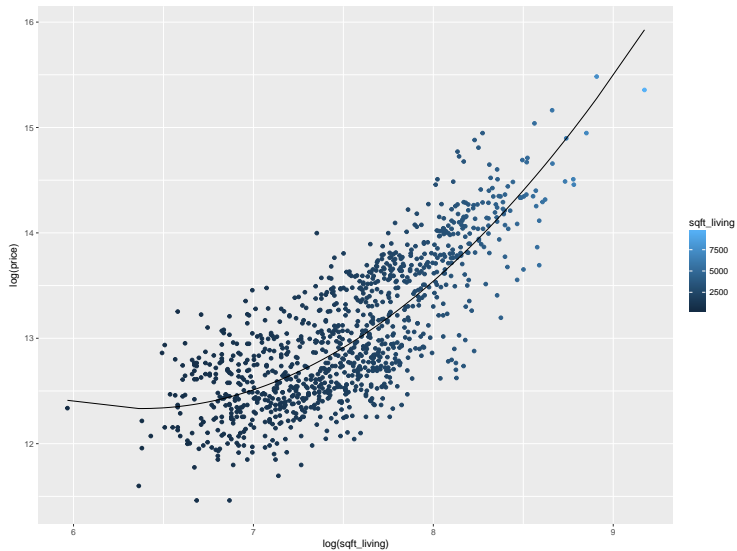
log(hp) vs log(sqft), Washington



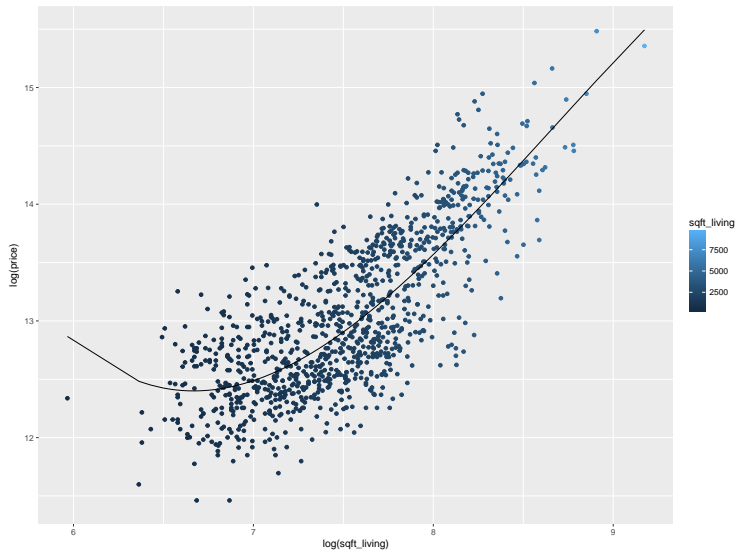
log(hp) vs log(sqft), Washington



log(hp) vs log(sqft), degree=2, Washington



log(hp) vs log(sqft), degree=3, Washington

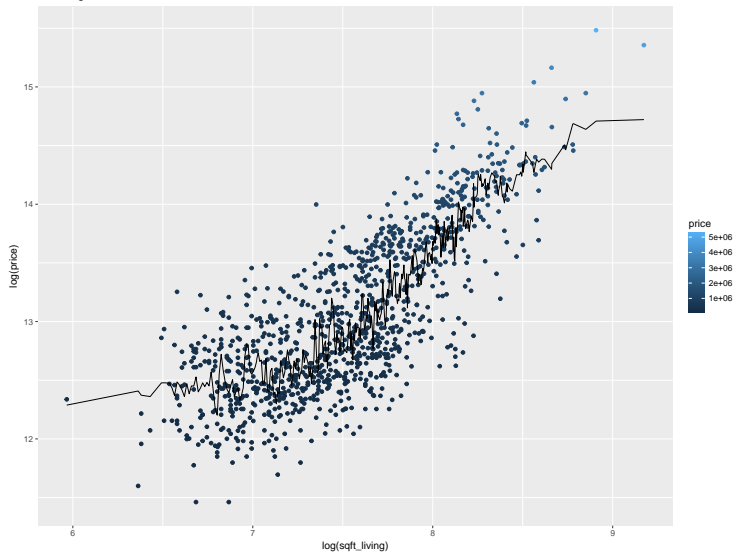




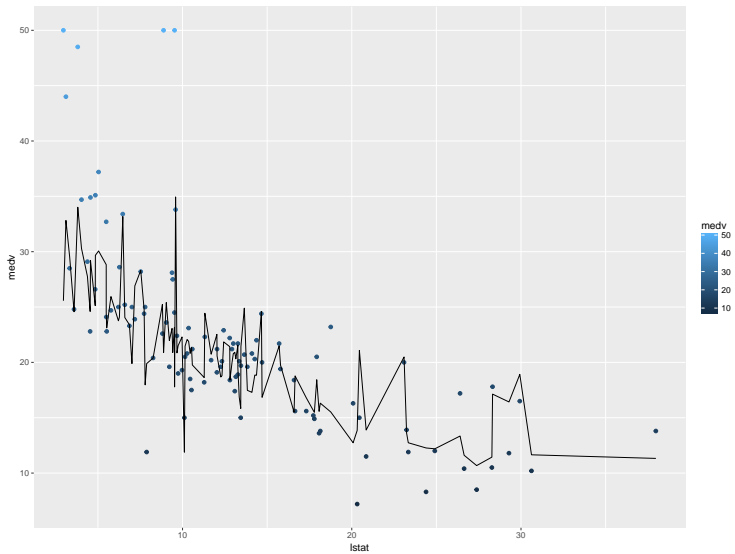
# K-NN Regressions

- ▶ Lineer olmayan regresyonlardan biri de k-NN regresyonlardır.
- ▶ Arkasında yatan motivasyon, yeni bir değer için tahmin istendiğinde, veri setinde ona en yakın değeri vermektir.
- ▶ Başarısı lineer ve polinom modele göre yüksektir.
- ▶ Ancak çoklu modeller ile ilgili çok fazla çalışma bulunmamaktadır.

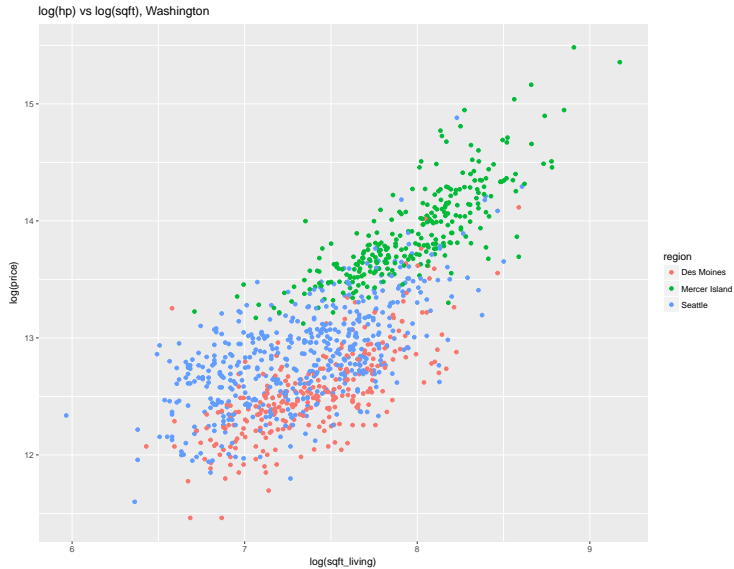
Washington, k = 10



Boston Data, Multivariate KNN, k = 50



# Verinin Yapısı



# Lineer Model için R kodları

<code>lm()</code>	Lineer model
<code>summary()</code>	Modelin özet istatistikleri
<code>coefficients()</code>	Hesaplanan katsayılar
<code>confint()</code>	Hesaplanan katsayılar için güven aralıkları
<code>fitted()</code>	Her bir gözlem için modelin hesapladığı değerler
<code>residuals()</code>	Hesaplanan değer ile gerçekleşen arasındaki farklar
<code>AIC()</code>	AIC değeri
<code>BIC()</code>	BIC değeri
<code>plot()</code>	Modelin başarısı hakkında grafikler
<code>predict()</code>	Modeli kullanarak yeni değerler için tahmin

# Kaynakça I



Kabacoff, Robert I,  
*R in Action*,  
2010.



Dalpiaz, David  
<https://daviddalpiaz.github.io/r4sl/k-nearest-neighbors.html>,  
2017.