

3.4. The Marketing Plan

Advertising datası, rsm su 7 soruyu cevaplıyoruz!

1. Reklam bütçesiyle satışlar arasında bir ilişki var mı?

Bu soruya cevap verebilmek için;

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon$$

TV, radio ve newspaper üzerinden sales'in multiple regression modelini uyguluyoruz ve $H_0: \beta_{TV} = \beta_{radio} = \beta_{newspaper} = 0$ hipotezini test ediyoruz.

F-statistiği null hipotezi reddedip reddetmeyeceğimizi belirlemek için kullanılabileceğimizi öğrenmiştik. Bu da değere kullanarak hesapladığımız değerler aşağıdaki gibiydi.

Quantity	Value
Residual standard error	1.69
R^2	0.897
F-statistic	570

F-statistiği 570 bulmuştuk. Bununla ilişkili p-value değeri oldukça küçük olduğu için null hipotezi reddetmiştik ve bu da reklamlar ile satışlar arasında bir ilişki olduğunu göstermişti.

2. İlişki ne kadar güçlü?

Null hipotezini reddettikten sonra modelin verilere ne ölçüde uyduğunu değerlendirmek için RSE (residual standard error) ve R^2 kullanıyoruz.

• RSE açıklanmayan $(Y = \beta_0 + \beta_1 X + \epsilon)$ e'nin standart sapmasının bir tahminiydi. Diğer bir deyişle; Y'nin gerçek regresyon çizgisinden sapacağı ortalama miktardı.

$$RSE = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Advertising datası için $RSE = 1.69$ çıkmıştı. Y'nin ortalama değeri 14,022 olduğuna göre kabaca %12'lik bir hata vermekte model.

• İkinci olarak R^2 , tahmin ediciler kullanılarak açıklanabilen Y 'deki değişkenlik modelin R^2 varyansını verirdi.

$$R^2 = \frac{TSS - RSS}{TSS} = \frac{\sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = \frac{\text{Modelin açıkladığı varyans}}{\text{Toplam varyans}}$$

Braim modelimizde $R^2 = 0.897$ çıkmıştı. Yani tahmin ediciler, y 'deki değişkenliğin neredeyse %90'ını açıklıyor.

3. Hangi medya satışlarla ilişkilidir?

Aşağıdaki formüller kullanarak her tahmin edicinin t -istatistiğini hesaplıyoruz

$$\sigma^2 = \text{Var}(e)$$

$$SE(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \quad SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$t = \frac{\hat{\beta}_i - 0}{SE(\hat{\beta}_i)}$$

Bu formüller simple linear regression için hesaplamıştık. Advertising dataseti için β_2 de eklenmesi gerekecek ve formüller değişecek. Her bir tahmin edicinin t -istatistiği ile ilişkili p -value'larını aşağıdaki şekilde tabloya yazmıştık.

	Coefficient	Std. error	t-Statistic	p-value
Intercept	2.939	0.3119	9.42	< 0.0001
TV	0.046	0.0014	32.81	< 0.0001
radio	0.189	0.0086	21.89	< 0.0001
newspaper	-0.001	0.0059	-0.18	0.8599

Her bir tahmin edicinin t -istatistiği ile ilişkili p -value'larına baktığımızda TV ve radyo için p -value'nun düşük olduğunu ama newspaper için düşük olmadığını görüyoruz. Bu, yalnızca TV ve radyonun satışlarla ilgili olduğunu göstermektedir.

4. Her bir ortam ve satış arasındaki ilişki ne kadar büyük?

β_j 'nin standart hatasının β_j için güven aralıkları oluşturmak için kullanılabileceğini görmüştük.

β_j %95 güven aralığında $\beta_j \pm 2 SE(\beta_j)$ aralığında yer alırdı. Advertising dataseti için her bir medya türünün bir multiple regression modelindeki katsayıları için %95 güven aralığını hesaplamak için yukarıdaki tabloyu kullanabiliriz. Güven aralıkları TV için (0.042, 0.049), radyo için (0.172, 0.206) ve newspaper için (-0.013, 0.011)

TV ve radyo için güven aralıkları dar ve 0'dan uzaktır ve bu medya türünün satışlarla ilgili olduğuna dair kanıt sağlar. Ancak newspaper aralığının dar olması, TV ve radyo değerleri gittikçe arttıkça değişkenin istatistiksel olarak anlamlı olmadığını göstermektedir.

✓ Collinearity'nin çok geniş standart hatalara yol açabileceğini görmekte. Newspaper ile güven analizi'nin bu kadar geniş olmasının nedeni collinearity olabilir mi? VIF ile kontrol olup olmadığına karar veriyorduk.

$$VIF = \frac{1}{1 - R^2_{X_j|X_{-j}}}$$

$R^2_{X_j|X_{-j}}$, X_j 'nin diğer tüm tahmin ediciler üzerindeki regresyonundan elde edilen değerdir. 1'e yakın olması X_j 'yi diğer X_j ile açıklayabildiğini, dolayısıyla korelasyon olduğu anlamına gelir. R^2 'nin 1'e yakın olması ise VIF değerini büyütecektir böyle bir durumda Genel olarak 5 veya 10'dan aşan VIF değeri bize varoluşlu bir collinearity haber veriyordu.

Ek not:

VIF 5 ise;

$$1 - R^2 = \frac{1}{5} \Rightarrow R^2 = \frac{4}{5} = 0,8$$

VIF 10 ise;

$$1 - R^2 = \frac{1}{10} \Rightarrow R^2 = \frac{9}{10} = 0,9$$

Advertising verisinde TV, radyo ve newspaper için VIF puanları sırasıyla 1.005, 1.115 ve 1.115'tir. Dolayısıyla collinearity olduğuna dair herhangi bir kanıt yoktur.

✓ Her bir ortamın satışlarla ilişkisini ayrı ayrı değerlendirmek için 3 ayrı simple linear regression gerçekleştirebiliyorduk. TV, radyo ve newspaper için ayrı ayrı yapmış olduğumuz simple linear regressionların istatistiksel değerlerini şöyle göstermiştik.

Simple regression of sales on TV

	Coefficients	Std. error	t-statistic	p-value
Intercept	7.0325	0.4578	15.36	< 0.0001
TV	0.0675	0.0027	12.67	< 0.0001

Simple regression of sales on radio

	Coefficients	Std. error	t-statistic	p-value
Intercept	9.312	0.563	16.54	< 0.0001
radio	0.203	0.020	9.92	< 0.0001

Simple regression of sales on newspaper

	Coefficients	Std. error	t-statistic	p-value
Intercept	12.351	0.621	19.88	< 0.0001
newspaper	0.055	0.017	3.30	0.0015

Bu tablolarda p-value'lara baktığımızda televizyon ile satış arasında ve radyo ile satış arasında son derece güçlü bir ilişki varken güneş ve satış arasında hafif bir ilişki olduğunu görüyoruz.

5. Gelecekteki satışları ne kadar doğru tahmin edebiliriz?

Yanıtı $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$ formülüne kullanarak tahmin edebiliriz. Bu tahminin ilgili accuracy $y = f(x) + e$ bireysel bir yanıtı veya $f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ ortalama yanıtı isteyip istemediğimize bağlıdır.

✗ $y = f(x) + e$ için bir tahmin aralığı kullanırız. Belirli bir şehir için satışları kapsayan tahmin belirsizliğini ölçmek için tahmin aralığı kullanıyorduk ve bu aralık güven aralığına kıyasla çok daha geniş bir aralık oluyordu. Çünkü güven aralığından farklı olarak irreducible hata (e)'yi de hesaba katıyorduk.

✗ $f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ için güven aralığı kullanırız. Çok sayıda şehirdeki ortalama satışlardaki belirsizliği ölçmek için güven aralığı kullanıyorduk. Bu güven aralığını, aralıkların %95'inin bize gerçek $f(x)$ değerini vermediği şeklinde yorumlamıştık.

6. İlişki lineer mi?

Non-linearity'yi belirlemek için residual plots kullanıyorduk. Eğer ilişki lineerse residual plotlar bir pattern göstermiyordu.

TV, radyo ve sales olmak üzere 3 eksenli olan grafikte residualları gösterdiğimizde (plasma şeklinde), pozitif residualların (yüzeğin üzerinde gösterilen) TV ve radyo kutucuklarının epi olarak bölündüğü 17 derecelik kutu boyunca uzanma eğilimindeyken, negatif residuallar kutucukların daha dengeli olduğu bu açıdan uzak durma eğilimindeydiler. Non-linearity ilişkisi burada gösterilebiliyordu.

Bunun yerine fitted values (\hat{y}_i) ve residuals eksenleriyle residual plot grafiği çizilip residualların non-linear bir pattern içermelerinden de yemin non-lineerliğini gösterilebiliyordu.

Yine non-linearity'yi katmak için polynomial regressionla (x^2) veya $\log x$, \sqrt{x} gibi yöntemlerle predictorı değiştirilerek model non-lineerlik katılabiliyordu ve yine lineer regresyon olarak kalıyordu.

$$mpg = \beta_0 + \beta_1 \times horsepower + \beta_2 \times horsepower^2 + e$$

7. Reklam medya araçları arasında bir sinerji var mı?

Standart lineer regresyon modelleri tahmin ediciler ve yanıt arasında additive bir ilişki olduğunu varsayıyor. Fakat bu varsayım her veri için geçerli olmayabiliyor.

Interaction term ekleyip bu deęişkenin kat sayısının p-value'una baki-
yorduk.

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times (\text{radio} \times \text{TV}) + \epsilon$$

	Coefficient	Std. error	t-statistic	p-value
Intercept	6.7502	0.248	27.23	< 0.0001
TV	0.0191	0.002	12.70	< 0.0001
radio	0.0289	0.009	3.24	0.0014
TV x radio	0.0011	0.000	20.73	< 0.0001

Bu terimin p-value'ı bu kadar küçükse böyle bir ilişki olabileceęi anlamına geliyor. Advertising data'sı için bu terimi eklediğimizde R²'nin % 30'undan neredeyse % 37'ye çıktığını gördük.