

## Regression, Part 2: Multivariable Linear Regression

Q Devise power scores for college basketball teams which take into account strength of schedule and the score differential of each game and home court advantage.

variables

- $i$  is the index of the  $i^{\text{th}}$  game
- team indices  $\{1, \dots, N\}$
- $Y_i$  is the score differential of game  $i$   
(Home score minus away score)
- $\beta_{H(i)}$  is the (unknown) power score of the Home team  $H(i)$  of game  $i$
- $\beta_{A(i)}$  is same for away team  $A(i)$

Model

$$Y_i = \beta_0 + \beta_{H(i)} - \beta_{A(i)} + \varepsilon_i$$

mean-zero noise  $\varepsilon_i$   
so,  $E\varepsilon_i = 0$

# 2023 NCAA Mens Basketball DataFrame

Season	WLoc	WTeamName	LTeamName	ScoreDiff	WScore	LScore
<dbl>	<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>
2023	H	DePaul	Loyola MD	6	72	66
2023	H	Duke	Jacksonville	27	71	44
2023	A	Evansville	Miami OH	-4	78	74
2023	A	FL Gulf Coast	USC	-13	74	61
2023	H	Florida	Stony Brook	36	81	45
2023	H	Florida Intl	Houston Chr	11	77	66

$$Y_1 = \beta_0 + \beta_{\text{DePaul}} - \beta_{\text{Loyola}} + \epsilon_1$$

$$Y_2 = \beta_0 + \beta_{\text{Duke}} - \beta_{\text{Jacksonville}} + \epsilon_2$$

$$Y_3 = \beta_0 + \beta_{\text{Miami}} - \beta_{\text{Evansville}} + \epsilon_3$$

⋮

In Matrix-Vector Form:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \end{bmatrix} = \begin{bmatrix} 1 & 1 & -1 & 0 & 0 & 0 & 0 & \dots \\ 1 & 0 & 0 & 1 & -1 & 0 & 0 & \dots \\ 1 & 0 & 0 & 0 & 0 & 1 & -1 & \dots \\ \vdots & & & & & & & \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_{\text{DePaul}} \\ \beta_{\text{Loyola}} \\ \beta_{\text{Duke}} \\ \beta_{\text{Jacksonville}} \\ \vdots \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \end{bmatrix}$$

# Vectorized Model

$$\vec{y} = X \vec{\beta} + \vec{\epsilon}$$

$\vec{y} = (y_1, \dots, y_n)$  vector of observed score differentials

$\vec{\beta} = (\beta_0, \beta_{\text{Akron}}, \beta_{\text{loyola}}, \dots)$  vector of unknown power ratings, to be estimated

$\vec{\epsilon} = (\epsilon_1, \dots, \epsilon_n)$  vector of mean zero noise  $E\epsilon_i = 0$

$X =$  Scheduling matrix so,

$j=1, X_{i1} = 1$  (intercept term)

$j>1, X_{ij} = X[\text{row } i, \text{column } j]$

$= \begin{cases} 1 & \text{if home team in} \\ & \text{game } i \text{ is team } j-1 \\ -1 & \text{if away team in} \\ & \text{game } i \text{ is team } j-1 \\ 0 & \text{else} \end{cases}$

```
> df_ncaamb2[1:5,]
# A tibble: 5 x 7
  Season WTeamName LTeamName WScore LScore WLoc ScoreDiff
  <dbl> <chr> <chr> <dbl> <dbl> <chr> <dbl>
1 2023 Abilene Chr Jackson St 65 56 H 9
2 2023 Akron S Dakota St 81 80 H 1
3 2023 Alabama Longwood 75 54 H 21
4 2023 Arizona Nicholls St 117 75 H 42
5 2023 Arizona St Tarleton St 62 59 H 3
> X[1:5, c(1:5, 131)]
(Intercept) Abilene Chr Air Force Akron Alabama Jackson St
[1,] 1 1 0 0 0 -1
[2,] 1 0 0 1 0 0
[3,] 1 0 0 0 1 0
[4,] 1 0 0 0 0 0
[5,] 1 0 0 0 0 0
```

Remove the vector superscripts for simplicity:  $y = X\beta + \epsilon$

How do we estimate the coefficients (e.g., the power ratings)  $\beta$  from observed data  $(X, y)$ ?

Recall that in simple linear regression, we estimated  $(\beta_0, \beta_1)$  by minimizing the Residual Sum of Squares.

Similarly, in multivariable linear regression we minimize the RSS,

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta} \operatorname{RSS}(\beta) \\ &= \operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - x_i^T \beta)^2\end{aligned}$$

where  $x_i$  is the  $i^{\text{th}}$  row of  $X$

and  $x_i^T \beta = x_i \cdot \beta = x_{i1} \beta_0 + x_{i2} \beta_1 + \dots + x_{i(k+1)} \beta_k$   
 $= \sum_{j=0}^k x_{i,j+1} \beta_j$  is the dot product

$$= \operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta)$$

in matrix form

$$= \operatorname{argmin}_{\beta} y^T y - 2\beta^T X^T y + \beta^T X^T X \beta$$

Multivariable Calculus: set the gradient equal to 0.  
The gradient is the analog of the derivative.

Gradient  $\nabla_{\beta} f(\beta) = \nabla_{\beta} f(\beta_0, \dots, \beta_k)$   
 $= \left( \frac{\partial f}{\partial \beta_0}, \dots, \frac{\partial f}{\partial \beta_k} \right)$   
is the vector of partial derivatives.

$$\begin{aligned} 0 &= \nabla_{\beta} \text{RSS}(\beta) \\ &= \nabla_{\beta} (y^T y - 2\beta^T X^T y + \beta^T X^T X \beta) \\ &= -2X^T y + 2(X^T X)\beta \end{aligned}$$

$$\Rightarrow X^T X \beta = X^T y$$

$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

This is the matrix form of linear multivariable regression!

HW suppose  $K=1$  so  $X = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix}$

and  $\beta = (\beta_0, \beta_1)$ ; this is simple linear regression.

Show that  $\hat{\beta} = (X^T X)^{-1} X^T y$  matches our previous simple linear regression formula for  $\beta_0, \beta_1$ .

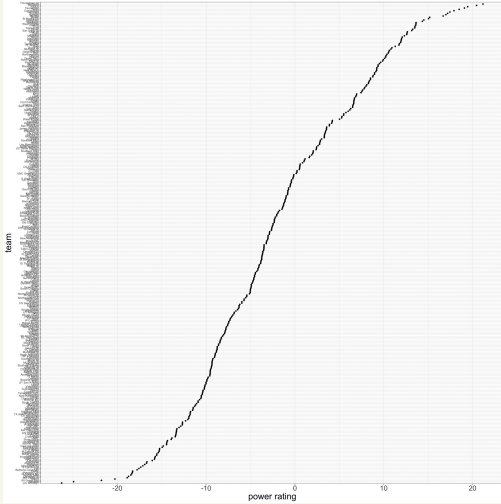
So, for our NCAA Basketball power ratings model  $y = X\beta + \epsilon$ , we now know how to estimate  $\hat{\beta}$ .

Let's run the computation and see what it says!

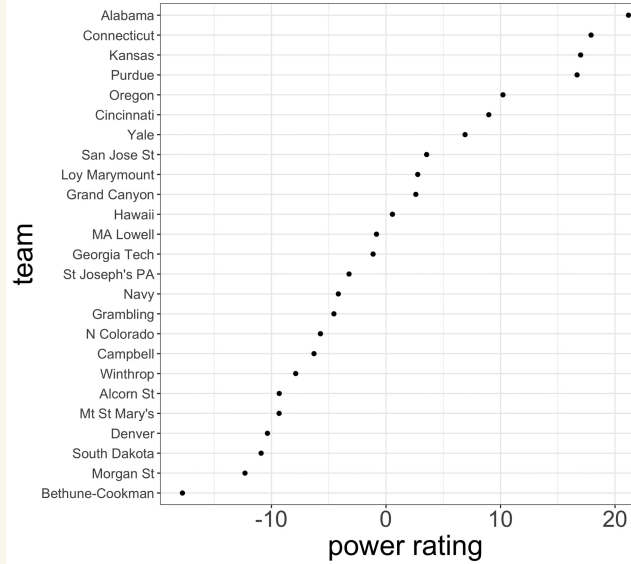
```
## get power ratings using multivariable linear regression  
power_ratings_model = lm(df_ncaamb2$ScoreDiff ~ X + 0)  
power_ratings = power_ratings_model$coefficients
```

Intercept  $\hat{\beta}_0 = 2 \rightarrow$  Home Court Advantage!

Too many teams to see,



Some power ratings:



```
> tibble(teams=colnames(X), power_ratings=unnamed(power_ratings)) %>%
+ drop_na() %>%
+ arrange(power_ratings) %>%
+ head(5)
# A tibble: 5 x 2
  teams      power_ratings
  <chr>      <dbl>
1 LIU Brooklyn -26.3
2 Hartford     -24.9
3 WI Green Bay -21.8
4 IUPUI        -20.3
5 MS Valley St -18.9
> tibble(teams=colnames(X), power_ratings=unnamed(power_ratings)) %>%
+ drop_na() %>%
+ arrange(-power_ratings) %>%
+ head(5)
# A tibble: 5 x 2
  teams      power_ratings
  <chr>      <dbl>
1 Alabama    21.2
2 Houston    20.5
3 UCLA       19.4
4 Tennessee  19.1
5 Texas      18.5
```

\* Biggest problem with this power Rating method:  
only works well if there is a "path" from  
each team to each other team.

Isolated groups of teams means we can't  
compare teams across groups. This is  
the problem with comparing international  
soccer teams in different leagues using  
power scores.

Q Use the power scores to simulate  
the NCAA tournament.

Model  $Y_i = \beta_0 + \beta_{H(i)} - \beta_{A(i)} + \epsilon_i$

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Since  $E\epsilon_i = 0$  still, we can use the same linear  
regression process to estimate  $\beta$ .

But now, we need to also estimate  $\sigma^2$ ,  
the (unknown) variance of the noise.



\* Assume we have  $\hat{\sigma}^2$ , an estimate of  $\sigma^2$ .

```
> sigma(power_ratings_model)
[1] 10.93223
```

\* Then we can simulate the result (score differential) of a basketball game by sampling  $\epsilon_i \sim N(0, \sigma^2)$

and then computing  $y_i = \hat{\beta}_0 + \hat{\beta}_{H(i)} - \hat{\beta}_{A(i)} + \epsilon_i$ .

\* We can then simulate the March Madness tournament by simulating each of the 65 games.

\* How to estimate  $\sigma^2$ :

$$\text{Error } \epsilon_i = y_i - x_i^T \beta$$

$$\text{Residual } \hat{\epsilon}_i = y_i - x_i^T \hat{\beta}$$

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{\epsilon}_i^2$$

$k = \# \text{ columns of } X$   
excluding intercept

Thm  $E[\hat{\beta}^2] = \sigma^2$  (unbiased)

HW Prove it. OR ask me to do it.

HW OR, Prove it in the case of simple linear regression  $K=1$

HW Create power slopes and then simulate the March Madness tournament.

Q Predict 400 meter dash time from a database of previous Races, which includes Runner names.

Variables  $i =$  index of  $i^{\text{th}}$  Runner-Race in the dataset

$Y_i =$  race time of  $i^{\text{th}}$  runner-race

$X_{i0} = 1$

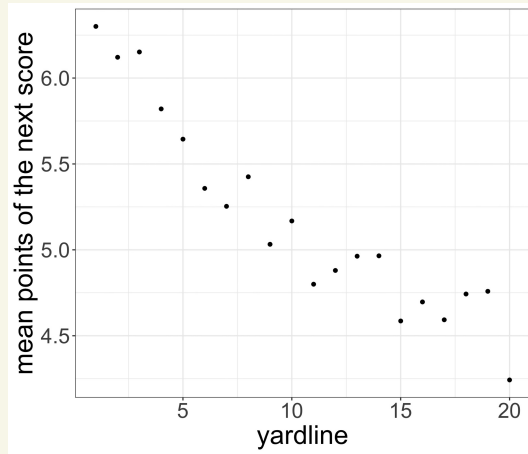
$X_{ik} = 1$  if player  $k$  is the Runner in the  $i^{\text{th}}$  Row, else 0

Model  $Y_i = \beta_0 + \sum_{j=1}^K X_{ij} \beta_j + \epsilon_i$ ,  $E\epsilon_i = 0$

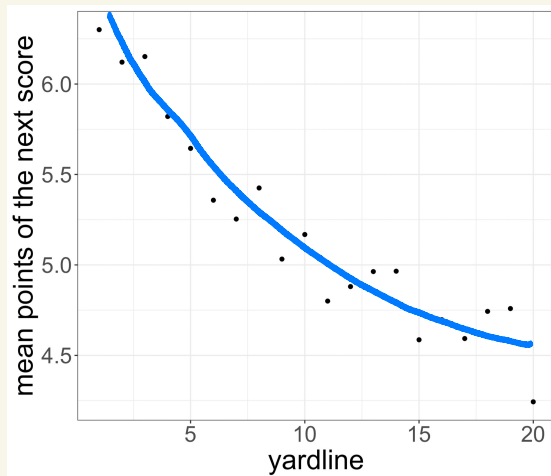
We now know how to estimate using linear regression,  $\hat{\beta} = (X^T X)^{-1} X^T y$ .

Q Estimate the expected points of the next score in a half in American football as a function of yardline, in the Redzone.

Generally, it is smart to begin with plotting



The Relationship looks quadratic, not linear



How can we use linear Regression to capture a nonlinear Relationship??

Data transformations!

Variables  $i$  = index of the  $i^{\text{th}}$  play in our dataset

$Y_i$  = points of the next score in the half after play  $i$  relative to the team with possession  
(a number in  $\{7, -7, 3, -3, 2, -2, 0\}$ )

$X_i$  = yardline of play  $i$   
( $X_i$  yards from opponent's endzone)

Linear Model  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

where  $\mathbb{E}\epsilon_i = 0$  (mean zero)

Quadratic Model  $Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2$

where  $\mathbb{E}\epsilon_i = 0$

The 2<sup>nd</sup> model has an additional parameter,  $\beta_2$ .  
Next time: how to estimate a regression model  
with  $k \geq 3$  parameters.

```
> m_ep_linear = lm(data=D3r, pts_next_score ~ yardline_100)
> m_ep_linear

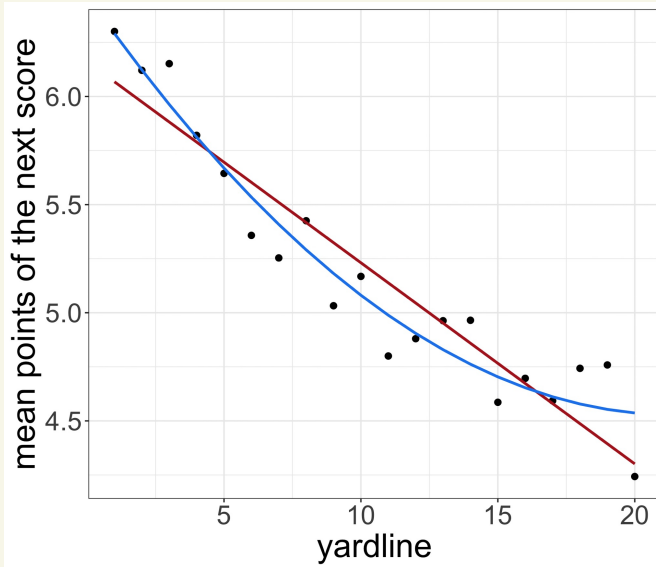
Call:
lm(formula = pts_next_score ~ yardline_100, data = D3r)

Coefficients:
(Intercept)  yardline_100
 6.16098      -0.09299

> ## quadratic model
> m_ep_quad = lm(data=D3r, pts_next_score ~ yardline_100 + I(yardline_100^2))
> m_ep_quad

Call:
lm(formula = pts_next_score ~ yardline_100 + I(yardline_100^2),
    data = D3r)

Coefficients:
(Intercept)      yardline_100  I(yardline_100^2)
 6.467712      -0.180798      0.004212
```



Quadratic model looks better.

HW Predict an NFL player's 2<sup>nd</sup> contract value (a proxy for his on field value) as a function of his draft position.

This relationship should be highly nonlinear!! Compare a linear model to a nonlinear model (say, a quartic polynomial), using linear regression to fit both. What fits better?

If we had more time: splines

Model any continuous shape

