

# Multivariable Linear Regression

## NCAA Men's Basketball Power Scores

We are given a dataset of the game results of each game in 2022-2023

Season	WLoc	WTeamName	LTeamName	ScoreDiff	WScore	LScore
	<dbl>	<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

$i$  = index of  $i^{th}$  game

$H(i)$  = home team index  
in game  $i$

$A(i)$  = away team

$y_i$  = score diff

Suppose each team  $j$  has a latent (unobserved) power score  $\beta_j$  and we model the outcome of a game (score diff) by

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

$\varepsilon_i$  is noise  $E\varepsilon_i = 0$  independent

$\beta_0$  is home court advantage

$$y_i = X_i \cdot \beta + \varepsilon_i$$

$$y_1 = \beta_0 + \beta_{\text{Defam}} - \beta_{\text{Logosla}} + \varepsilon_1$$

$$Y_2 = \beta_0 + \beta_{\text{Duke}} - \beta_{\text{Jack}} + \varepsilon_2$$

In matrix-vector form;

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ \text{red dot} \end{bmatrix} = \begin{bmatrix} 1 & -1 & 0 & 0 & \dots \end{bmatrix} \cdot \begin{bmatrix} \beta_0 \\ \beta_{\text{Dept}} \\ \beta_{\text{Logis}} \\ \beta_{\text{Source}} \\ \beta_{\text{Prod}} \\ \vdots \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \text{blue wavy line} \end{bmatrix}$$

*interpret cols*

*Dept* *Logis* *Source* *Prod* ...

$\bar{y} = X\bar{\beta} + \bar{\varepsilon}$

Overtime (data)  $\hookrightarrow$  scheduling matrix  $\hookrightarrow$  data

	Season	WTeamName	LTeamName	WScore	LScore	WLoc	ScoreDiff
	<dbl>	<chr>	<dbl>	<dbl>	<chr>	<dbl>	<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

How do we estimate the coefficients (the power scores)  $\beta$  from the data  $(X, y)$ ?

In Simple Linear regression, we estimated  $(\beta_0, \beta_1)$  by minimizing RSS and we do the same thing here

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\underbrace{\beta_{\text{H(i)}} + \beta_{\text{A(i)}} + \beta_0}_{\text{add like to use the } X \text{ matrix}}))^2$$

add like to use the  $X$  matrix

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \vec{x}_i \cdot \vec{\beta})^2$$

ith Row of  $X$

$$(X\beta)_i = x_i \beta$$

$$x_i^T = [1 \quad \underbrace{\quad \quad \quad}_{\text{intercept}} \quad \underbrace{-1}_{H(i)} \quad \underbrace{-1}_{A(i)}]$$

Set gradient equal to 0 and solve

$$\nabla_{\beta} \text{RSS}(\beta) = 0 \quad \text{and solve}$$

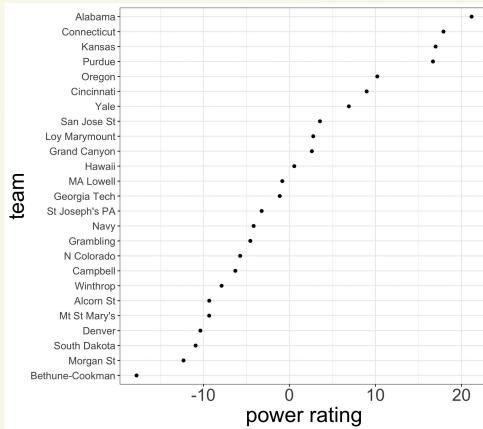
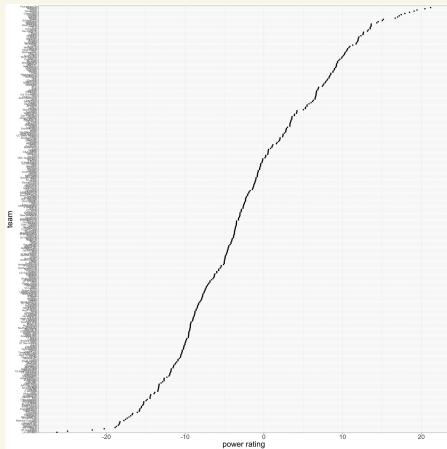
$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$X_{n \times k}$

$$\left\{ \begin{array}{l} \vec{y} = X \vec{\beta} + \vec{\varepsilon} \\ X^T y = X^T X \beta \\ (X^T X)^{-1} X^T y = \beta \end{array} \right.$$

So for our NCAA basketball power rating model  $y = X \cdot \beta + \varepsilon$  and how we can estimate  $\beta$ .

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



```
> tibble(teams=colnames(X), power_ratings=uname(power_ratings)) %>%
+   drop_na() %>%
+   arrange(power_ratings) %>%
+   head(5)
# A tibble: 5 × 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=uname(power_ratings)) %>%
+   drop_na() %>%
+   arrange(-power_ratings) %>%
+   head(5)
# A tibble: 5 × 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
```

Intercept  $\hat{\beta}_0 = 2$

Being at home increases the predicted score diff by 2 pts above being on a neutral court.

# Expected Points in American Football

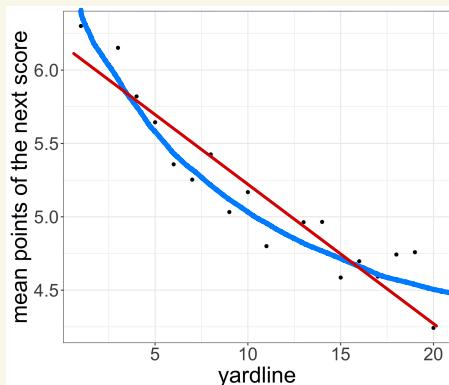
We are given a dataset of NFL plays

$\left\{ \begin{array}{l} i = \text{index } i^{\text{th}} \text{ play} \\ ydl_i = \text{yardline} \\ y_i = \text{net points of the next score} \\ \text{in the half} \in \{7, 3, 2, 0, -2, -3, -7\} \end{array} \right.$

Define expected point (EP) by a model

$$E[y_i] = f(ydl_i)$$

Plot



Might not  
be linear...

empirical mean

Begin with just the red zone for simplicity.

How can we use linear regression to capture a relationship that is not linear?

## data transformation

Linear model  $y_i = \beta_0 + \beta_1 \cdot y_{\text{dl}} + \epsilon_i$

Quadratic model  $y_i = \beta_0 + \beta_1 \cdot y_{\text{dl}} + \beta_2 \cdot y_{\text{dl}}^2 + \epsilon_i$

$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$

use multivariable regression with new features that we created transformed from the original features

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

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \end{bmatrix} = \begin{bmatrix} 1 & y_{\text{dl}} & y_{\text{dl}}^2 \\ 1 & y_{\text{dl}2} & y_{\text{dl}2}^2 \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \epsilon$$

$$y = X \cdot \beta + \epsilon$$

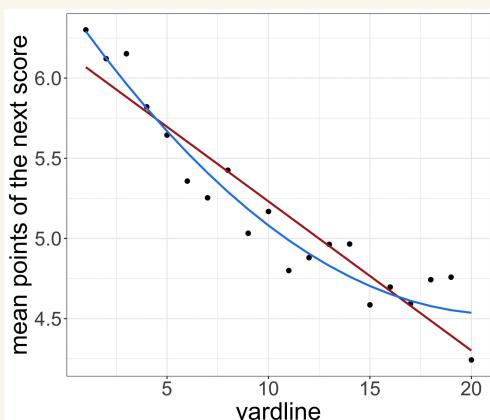
```
> m_ep_linear = lm(data=D3r, pts_next_score ~ yardline_100)
> m_ep_linear

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
```



# NFL Draft Expected Value Curves

We are given a dataset of NFL draft picks,

$i = i^{\text{th}}$  draft pick in our dataset  
 $X_i = \text{draft pick number } \in \{1, 2, \dots, 256\}$   
 $y_i = \text{player who was drafted's first contract performance "Value"}$

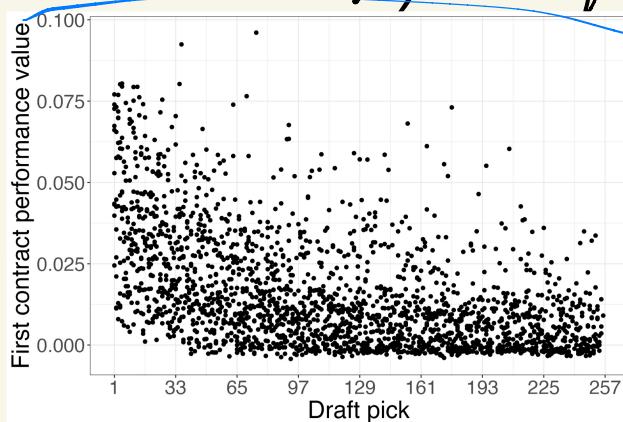
player_id	player_name	year	t	draft_pos	firstContractPerformanceValue
40688	A.J. Bouye	2013	1	N/A	1.659893e-02
42410	A.J. Cann	2015	1	67	3.541377e-02
35558	A.J. Edds	2011	1	119	-7.510774e-04
37077	A.J. Green	2011	1	4	8.018761e-02
30819	A.J. Hawk	2006	1	5	4.689117e-02
35863	A.J. Jefferson	2010	1	N/A	4.419914e-03
38560	A.J. Jenkins	2012	1	30	5.734494e-03
40096	A.J. Klein	2013	1	148	1.305431e-02
30972	A.J. Nicholson	2006	1	157	-2.287106e-03

Model

$$\mathbb{E}[y_i] = f(x_i)$$

ignore position  
ignore everything else

Before specifying a specific model, Plot



\$  
salary cap

the draft is very noisy.

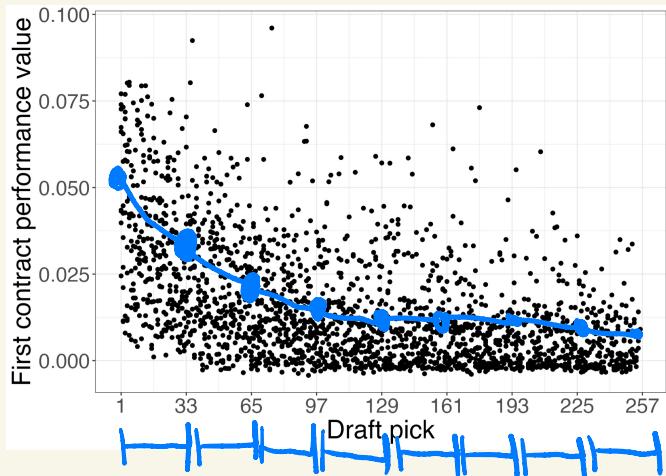
the expected value curve is nonlinear.

Convex: the dropoff in value b/t picks  $t$  and  $t+1$  decreases as  $t$  increases

## Spline:

to model a general nonlinear shape, using a Spline often works.

To fit a spline, you fit a separate polynomial (usually a cubic) to different subsections of the data



a Spline also forces the curve to be Smooth, even at the Knots:

at each knot, we mandate that the curve has

{ the same left y value and right y value  
the same left derivative and right derivative  
the same left 2<sup>nd</sup> derivative and right 2<sup>nd</sup> derivative

Suppose we fit a cubic spline with one knot at  $x = k$   
(e.g.  $k = 129$  in middle of draft)

we model  $y_i = f(x_i | \beta) + \varepsilon_i$   
where  $f$  is a spline and  $\beta$   
are the spline parameters,

$$f(x|\beta) = \begin{cases} \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 & \text{if } x \leq K \\ \beta_4 + \beta_5 x + \beta_6 x^2 + \beta_7 x^3 & \text{if } x > K \end{cases}$$

Such that

$$\left\{ \begin{array}{l} \lim_{x \rightarrow K^-} f(x|\beta) = \lim_{x \rightarrow K^+} f(x|\beta) \\ \lim_{x \rightarrow K^-} f'(x|\beta) = \lim_{x \rightarrow K^+} f'(x|\beta) \\ \lim_{x \rightarrow K^-} f''(x|\beta) = \lim_{x \rightarrow K^+} f''(x|\beta) \end{array} \right.$$

enforce

$$\left\{ \begin{array}{l} \beta_0 + \beta_1 K + \beta_2 K^2 + \beta_3 K^3 = \beta_4 + \beta_5 K + \beta_6 K^2 + \beta_7 K^3 \\ \beta_1 + 2\beta_2 K + 3\beta_3 K^2 = \beta_5 + 2\beta_6 K + 3\beta_7 K^2 \\ 2\beta_2 + 6\beta_3 K = 2\beta_6 + 6\beta_7 K \end{array} \right.$$

$$\left\{ \beta_7 = (2\beta_2 + 6\beta_3 k - 2\beta_6) / (6k) \right.$$

$\beta_7$  is dead

↓  
each knot in  
a cubic spline  
"kills" 3 parameters

with 1 knot, we need  
to estimate 5 parameters

$$f(x|\beta) = \begin{cases} \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 & \text{if } x \leq k \\ \beta_4 + \beta_5 x + \beta_6 x^2 + \beta_7 x^3 & \text{if } x > k \end{cases}$$



$$\gamma_1 \tilde{x}_{1i} + \gamma_2 \tilde{x}_{2i} + \gamma_3 \tilde{x}_{3i} + \gamma_4 \tilde{x}_{4i} + \gamma_5 \tilde{x}_{5i}$$

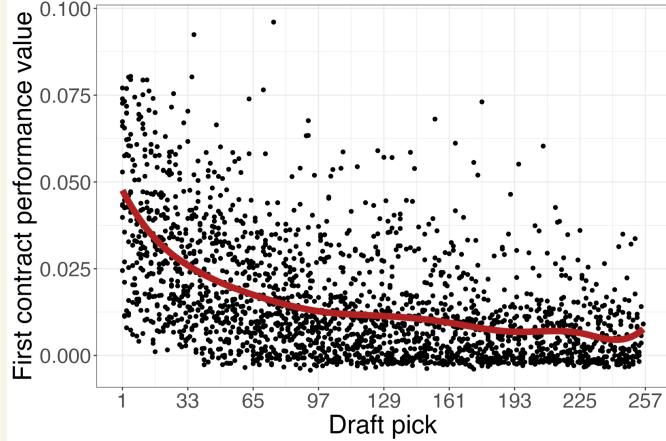
↳ the Spline basis

```

draft_model1 = lm(
  firstContractPerformanceValue ~ splines::bs(draft_pos, degree=3, knots=seq(33, 32*8, by=32)),
  data=df_draft
)
draft_model1

```

$$y \sim \text{Spline}(X, \text{deg}=3, \text{knots}=\binom{33}{65, \dots})$$



```

draft_model2 = lm(
  firstContractPerformanceValue ~ splines::bs(draft_pos, degree=3, df=5),
  data=df_draft
)
draft_model2

```

degrees of freedom

