

Analiza podataka košarkaša Joel Embiid-a u R programskom jeziku

Adi Džombić

December 11, 2024

Ovaj rad analizira performanse košarkaša Joel Embiid-a korištenjem stečenog znanja iz statistike i vjerovatnoće. Koristeći statističke alate i modele u R-u, istražujemo faktore koji utječu na ishode igre, kao što su učinkovitost igrača kao i ostale faktore koje utiču na njegovu igru. Putem metoda kao što su regresijska analiza i prediktivno modeliranje, cilj nam je otkriti obrasce koji mogu pružiti korisne uvide za poboljšanje izvedbe i optimizaciju strategija u košarci.

1. Uvod

Data Science. Integracija znanosti o podacima u sport revolucionirala je način na koji se igre analiziraju, razvijaju strategije i donose odluke. Košarka, sa svojom brzom akcijom i bogatim skupovima podataka, nudi savršenu arenu za primjenu tehnika znanosti o podacima. Ovo istraživanje usmjereno je na analizu podataka o košarci kako bi se otkrili obrasci i uvidi koji mogu poboljšati razumijevanje igre i informirati strateške odluke.

Cilj našeg istraživanja. Primarni cilj ovog rada je istraživanje različitih aspekata košarkašove igre kroz analizu podataka. Ključna područja fokusa uključuju procjenu učinka igrača, razumijevanje dinamike tima i identifikaciju faktora koji doprinose uspješnim ishodima igre. Koristeći statističke i računske tehnike, ova studija ima za cilj da izvuče značajne informacije iz neobrađenih skupova košarkaških podataka, pretvarajući ih u uvide koji se mogu primijeniti.

2. Biblioteke i podaci koje koristimo

Podaci koje obrađujemo u ovom radu su prikupljeni od strane nekoliko košarkaških portala, te spakovani u jednostavnu R biblioteku pod nazivom "**nbastatR**"¹. Biblioteka se sastoji od mnogobrojnih košarkaša, međutim zbog kratke dužine ovog kursa, obradit ćemo samo jednog igrača pod nazivom Joel Embiid², usporediti njegovu igru sa par igrača njegovog nivoa i probat pronaći razloge šta ga čini jednim od najboljih igrača u NBA današnjice.

Pored te biblioteke, iskoristit ću snagu "**Tidyverse**"³ biblioteke. **Tidyverse** je kolekcija R biblioteka dizajniranih za obradu podataka, a razvili su je i održavali Hadley Wickham i RStudio tim. Izgrađen je oko strukture podataka nazvane "tibble", koja je slična podatkovnom okviru, ali s dodatnom funkcionalnošću. Tidyverse slijedi skup pravila nazvanih "tidy data" koji promoviraju strukturiran, lako razumljiv i učinkovit proces manipulacije podacima. Matematičari i analitičari podataka naširoko koriste ove principe i biblioteke u "tidyverse", čineći ga moćnim i praktičnim alatom za analizu podataka.

¹<https://www.rdocumentation.org/packages/nbastatR/>

²<https://en.wikipedia.org/wiki/Embiid/>

³<https://www.tidyverse.org/>

2.1. Pristup podacima o igraču

Za početak instalirajmo prethodno navedene biblioteke za obradu podataka:

```
1 install.packages("nbastatR")
2 install.packages("tidyverse")
3 install.packages("tvthemes")
4
5 library(nbastatR)
6 library(tidyverse)
7 library(tvthemes)
```

Biblioteka "nbastatR" sadrži hiljade igrača... Potrebno je izdvojiti podatke našeg ciljanog igrača i njegove vrijednosti smjestiti u varijablu. To ćemo uraditi na sljedeći način:

```
1 Sys.setenv(VROOM_CONNECTION_SIZE=500072)
2
3 jojo_game_log <- suppressWarnings(
4   game_logs(
5     seasons = c(2015:2023),
6     league = "NBA",
7     result_types = "player",
8     season_types = "Regular Season")) %>%
9   filter(namePlayer == 'Joel Embiid')
```

Prva komanda proširuje default prostor varijable s obzirom da se radi o ogromnoj količini podataka. Veoma korisna stvar kod biblioteke "nbastatR" koju vrijedi napomenuti je da samo uz par linija koda, uspjeli smo "izvući" podatke iz njegovih 367 utakmica koje je odradio u periodu 2015-2023 godine.

2.2. Prosjek poena po utakmici po sezoni



Korištenjem sljedećeg koda dobili smo prethodni graf koji prikazuje prosjek poena po utakmici po sezoni našeg igrača. Ako pogledamo graf, možemo primijetiti da je prosjek poena drastično pao u sezoni 2019-2020. Nakon malo istraživanja, dolazimo do dva glavna razloga zbog ovakvog perfomansa u toj sezoni, a to su: ozljeda koljena i leđa, te pandemija COVID-19...

Ozljeda koljena je uzrokovala da propusti čak 22 utakmice⁴ tokom te sezone, ali sezone poslije te nastavlja da se konstantno poboljšava.

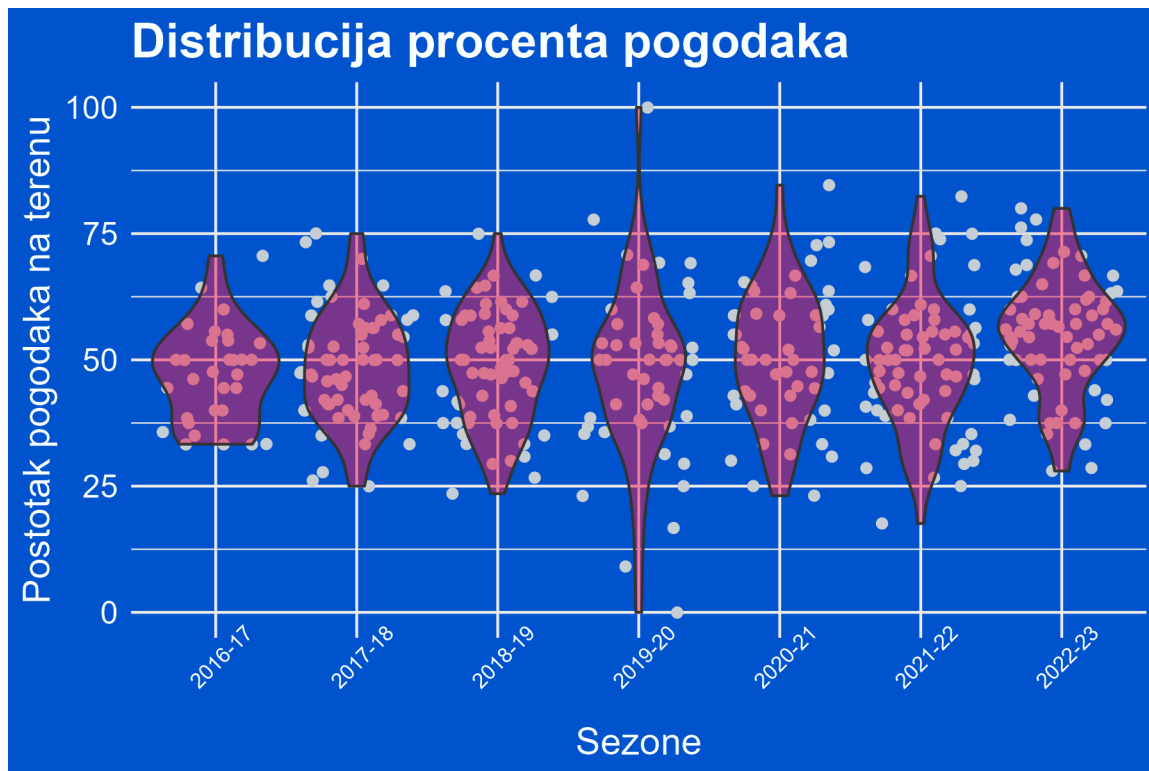
```
1 joel_points_per_season <- jojo_game_log %>%
2   group_by(slugSeason) %>%
3   summarise(avg_points = mean(as.numeric(pts), na.rm = TRUE)) %>%
4   arrange(slugSeason)
5
6 ggplot(joel_points_per_season, aes(x = slugSeason, y = avg_points, group = 1))
7   +
8   geom_line(color = "white", size = 1.5) +
9   geom_point(color = "red", size = 4) +
10  labs(
11    title = "Joel Embiid Prosjek poena po utakmici po sezoni",
12    x = "Sezone",
13    y = "Prosjek poena po utakmici"
14  ) +
15  theme_brooklyn99() +
16  theme(
17    axis.text.x = element_text(angle = 45, hjust = 1, size = 8),
18    plot.title = element_text(size = 10, face = "bold"),
19    axis.title = element_text(size = 10)
20  )
```

2.3. Distribucija postotka šuteva

Metrika koja nas takođe zanima je Joel Embiid-ov postotak šuta. Distribucije u nastavku pokazuju da on konstantno šutira preko 50%. Kako mu je karijera napredovala, počeo je stalno pucati preko 60%. Gledajući prosjek sezone, on je iznad 53,6%, što je rekord karijere. Bez iznenađena, najmanji procenat šuteva je napravio u sezoni 2019-2020 zbog povrede.

```
1 jojo_game_log %>%
2   ggplot(aes(x=slugSeason,y=round(pctFG*100,1))) +
3   geom_jitter(color='#C4CED4')+
4   geom_violin(alpha=.5,fill='#ED174C')+
5   theme_brooklyn99()+
6   theme(legend.position = 'none')+
7   labs(x = "Sezone", y = "Postotak pogodaka na terenu") +
8   ggtitle('Distribucija procenta pogodaka')
```

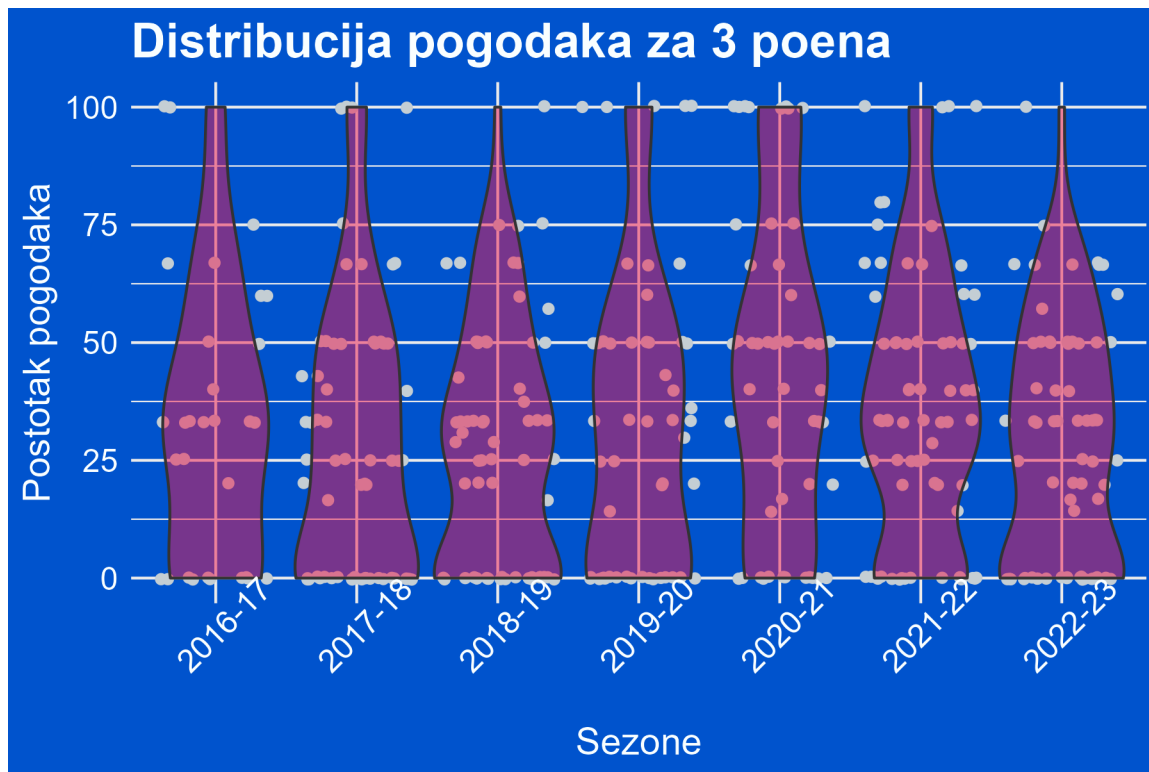
⁴<https://bit.ly/4iozfSp>



2.4. Distribucija pogodaka za "tricu"

Danas u košarci sve više igrača šutira trice. Visoki košarkaši nisu iznimka. Zapravo, ove sezone Joel Embiid je pokušao prosječno 3,25 šuteva za tricu po utakmici, malo niže nego u proteklih nekoliko sezona. Gledajući njegov postotak šuteva za tri poena, mnogo je manje konstantan.

```
1 jojo_game_log %>%
2   ggplot(aes(x=slugSeason,y=round(pctFG3*100,1))) +
3   geom_jitter(color = '#C4CED4')+
4   geom_violin(alpha=.5,fill='#ED174C')+
5   theme_brooklyn99()+
6   theme(legend.position = 'none', axis.text.x = element_text(angle = 45))+
7   labs(x = "Sezone", y = "Postotak pogodaka") +
8   ggtitle('Distribucija pogodaka za 3 poena')
```

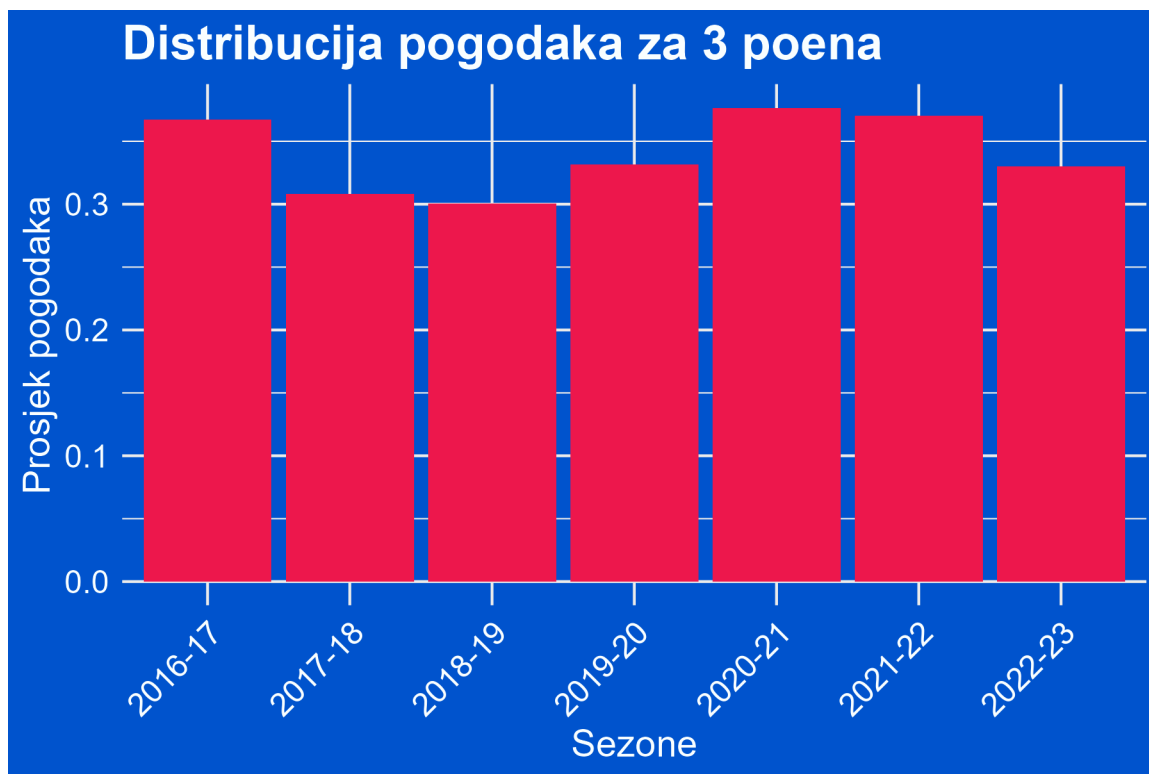


Iako iz utakmice u utakmicu prebacuje malo preko 3,25 šuteva za tricu, gledajući statistiku sezone, ostao je igrač koji pogađa 30% šuteva za tri poena. Još jedan prikaz ovih rezultata možemo posmatrat sljedećim grafom:

```

1 jojo_game_log %>%
2   group_by(slugSeason) %>%
3   summarise(games_played = n(),
4             avg_minutes_played = mean(minutes),
5             fgm = sum(fg3m),
6             fga = sum(fg3a)) %>%
7   mutate(fg_pct = fgm / fga) %>%
8   ggplot(aes(x = slugSeason, y = fg_pct)) +
9     geom_bar(stat = 'identity', fill = '#ED174C') +
10    theme_brooklyn99() +
11    theme(
12      legend.position = 'none',
13      axis.text.x = element_text(angle = 45, hjust = 1)
14    ) +
15    labs(x = "Sezone", y = "Prosjeak pogodaka") +
16    ggtitle("Distribucija pogodaka za 3 poena")

```



3. Usporedba Joel Embiid-a sa ostalim All-Star igračima

U ovom poglavlju poredit ćemo Joel Embiid-a sa ostalim All-Star igračima u 3 kategorije, a to su: prosjek pogodaka po utakmici u sezoni 2022-2023, prosjek bacanja na koš i prosjek pobjeda svakog igrača.

Prvo, moram izvući novi skup podataka pojedinih All-Stars igrača i Joel za sezonu 2022-2023. Srećom, biblioteka "nbastatsR" nam olakšava izvlačenje skupova podataka o igračima, timovima i sezonama.

3.1. Izvlačenje podataka

Sa narednim kodom izvlačim statistiku 22-23 igrača i filtriram samo igrače iz All-Star tima i Joela Embiid-a:

```
1 all_stars = c('Giannis Antetokounmpo', 'Kevin Durant',  
2             'Jayson Tatum', 'Donovan Mitchell',  
3             'Kyrie Irving', 'LeBron James', 'Nikola Jokic',  
4             'Zion Williamson', 'Stephen Curry',  
5             'Luka Doncic', 'Joel Embiid')  
6  
7 allStar_log <- suppressWarnings(  
8     game_logs(seasons = c(2023),  
9     league = "NBA",
```

```

10         result_types = "player",
11         season_types = "Regular Season"))%>%
12         filter(namePlayer %in% all_stars)

```

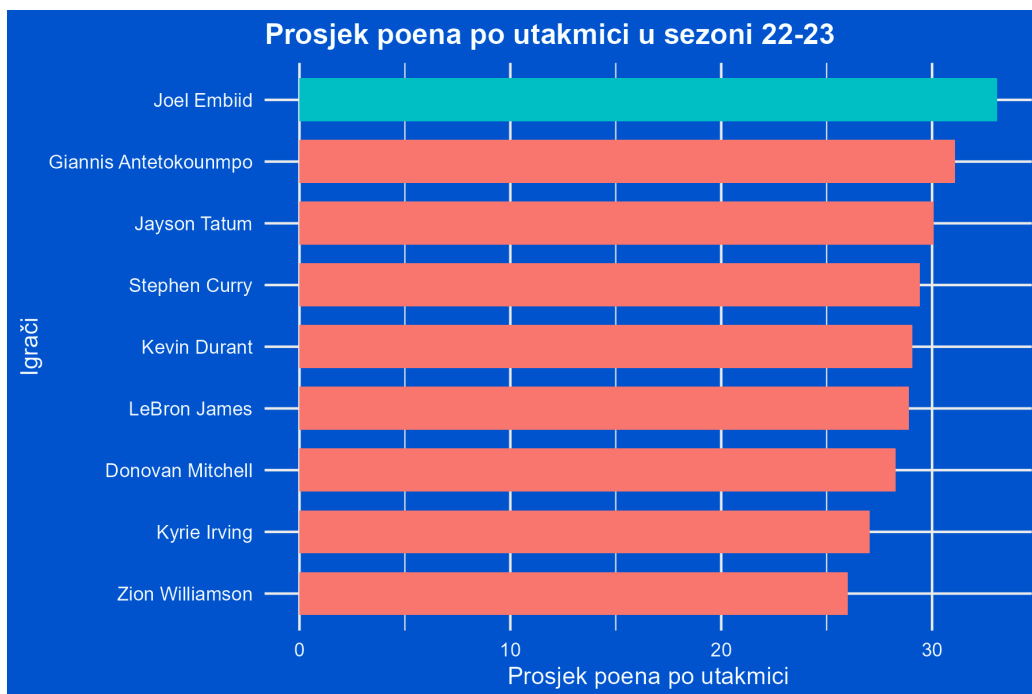
3.2. Prosjek pogodaka po utakmici u sezoni 2022-2023

Prva statistika koju sam pogledao bio je prosjek poena po utakmici. Joel zabija više poena po utakmici nego bilo koji od odabranih All Starsa. Naivno bi mogli zaključiti da je on jednostavno bolji od svih ostalih, međutim, može se desiti i da on "baca" na koš dosta više od svih ostalih košarkaša... Iz tog razloga je potrebno pogledati koji je njegov prosjek bacanja na koš u sljedećem poglavlju.

```

1 all_stars = c('Giannis Antetokounmpo', 'Kevin Durant',
2               'Jayson Tatum', 'Donovan Mitchell',
3               'Kyrie Irving', 'LeBron James', 'Nikola Jokic',
4               'Zion Williamson', 'Stephen Curry',
5               'Luka Doncic', 'Joel Embiid')
6
7 allStar_log <- suppressWarnings(
8     game_logs(seasons = c(2023),
9                   league = "NBA",
10                  result_types = "player",
11                  season_types = "Regular Season"))%>%
12     filter(namePlayer %in% all_stars)

```



3.3. Prosjek bacanja na koš

Iako Joel nije na vrhu popisa, on ima bolji postotak od barem polovice All-Star igrača. Ipak, istok izgleda malo konkurentniji od zapada kada je riječ o ovoj statistici. Prosjek bacanja na koš dobijamo sljedećim kodom:

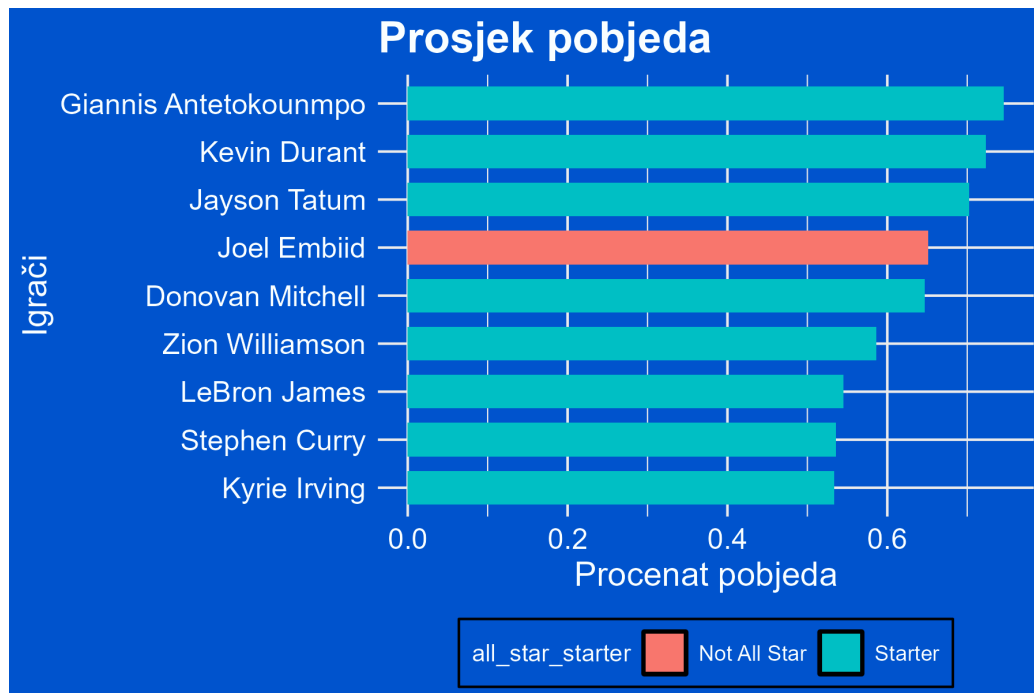
```
1 allStar_log %>%
2   group_by(namePlayer) %>%
3   summarise(games_played = n(),
4             avg_minutes_played = mean(minutes),
5             mean_points = mean(pts),
6             fgm = sum(fgm),
7             fga = sum(fga)) %>%
8   mutate(fg_pct = fgm / fga,
9          all_star_starter = ifelse(namePlayer == 'Joel Embiid', 'Joel Embiid',
10                                   'Starter')) %>%
11   ggplot(aes(x=reorder(namePlayer, fg_pct), y=fg_pct, fill=all_star_starter)) +
12   coord_flip() +
13   geom_bar(stat='identity') +
14   theme_brooklyn99() +
15   labs(x = "Igraci", y = "Postotak bacanja") +
16   ggtitle('Prosjek bacanja na kos')
```



3.4. Prosjek pobjeda u 2022-2023 sezoni

Konačno, koliki je postotak pobjeda ovih igrača. Ove sezone Sixersi su trenutno treći i prilično im dobro ide. Joel je bio ozlijeđen ranije tokom sezone, ali je ove sezone igrao onoliko minuta koliko je igrao u prethodnim sezonama. Njegov prosjek pobjeda je veći od 60%, a samo nekoliko od svih All-Star igrača će vjerovatno imati veći prosjek od njega. Sljedećim kodom ćemo dobiti graf prosjeka pobjeda nekoliko igrača:

```
1 allStar_log %>%
2   group_by(namePlayer) %>%
3   summarise(games_played = n(),
4             avg_minutes_played = mean(minutes),
5             mean_points = mean(pts),
6             fgm = sum(fgm),
7             fga = sum(fga),
8             win = sum(iffelse(isWin == TRUE , 1,0))) %>%
9   mutate(fg_pct = fgm / fga,
10          win_pct = win / games_played,
11          all_star_starter = iffelse(namePlayer == 'Joel Embiid','Not All Star'
12                                     , 'Starter')) %>%
13   ggplot(aes(x=reorder(namePlayer,win_pct),
14                     y=win_pct, fill=all_star_starter)) +
15   coord_flip()+
16   geom_bar(stat='identity', width = 0.7)+
17   theme_brooklyn99() +
18   labs(x = "Igraci", y = "Procenat pobjeda") +
19   ggtitle('Prosjek pobjeda')
```



3.5. Nastavak...