



Universidade de Brasília
Departamento de Estatística

Modelagem de dados para análise da influência de testes físicos obtidos no
Combine no desempenho dos jogadores de Skill Positions na NFL

Gabriel Vêras Monteiro

Projeto apresentado para o Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários para obtenção do grau de Bacharel em Estatística.

Brasília
2023

Gabriel V́eras Monteiro

**Modelagem de dados para análie da influência de testes físicos obtidos no
Combine no desempenho dos jogadores de Skill Positions na NFL**

Orientador: Prof. Donald Matthew Pianto

Projeto apresentado para o Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários para obtenção do grau de Bacharel em Estatística.

**Brasília
2023**

Agradecimentos

Aos meus pais, que sempre me apoiaram e deram condições de fazer esse trabalho.

À Carlyne, que é uma pessoa tão especial que conheci durante a graduação, e que me acompanhou e ajudou a passar por todas as etapas dentro e fora da UnB.

Ao meu orientador Donald, que aceitou minha proposta de tema e tornou isso possível.

A todos os colegas da estatística que compartilharam bons momentos de aprendizado.

Ao *NFLScoutBR*, que tornou todo o trabalho possível, pela ajuda a coletar os dados, e todos os ensinamentos e discussões sobre o tema.

Resumo

Este trabalho almeja investigar, de acordo com dados da empresa *Pro Football Focus*, que apresenta análises de desempenho na *National Football League*, e dados dos testes físicos do Combine, uma relação entre o desempenho de um jogador de posições de habilidade (*Quarterbacks*, *Wide receivers* e *Running backs*) na liga e suas medidas físicas obtidas. Para tal, foram desenvolvidos modelos que, tendo em vista a amostra geral dessas posições, e também cada uma individualmente, buscavam se ajustar da melhor maneira possível para obter resultados mais profundos sobre a seleção de jogadores para um time na liga. Assim, a partir do *Software* de programação R, foram ajustados então modelos de regressão linear, que indicam variáveis distintas nos ajustes para cada posição em relação ao modelo geral, o que é de suma importância, pois traz uma riqueza maior para o estudo de cada uma em específico. Após do estudo dos modelos, constatou-se baixa explicação da variável resposta pelos testes do *Combine*, para todas as amostras.

Palavras-Chave: NFL, Combine, Modelagem, Regressão.

Abstract

This work aims to investigate, according to data from the company Pro Football Focus, which provides performance analysis in the National Football League, and data from the physical tests of the Combine, a relationship between the performance of a skill position player (Quarterbacks, Wide receivers and Running backs) in the league and their obtained physical measurements. To achieve this, regression models were developed, taking into account the overall sample of these positions, as well as each one individually, to obtain more in-depth results. At the end of the study, using the R programming software, linear regression models were then adjusted, indicating distinct variables in the adjustments for each position in relation to the general model. This is of utmost importance as it provides greater accuracy to the study of each specific position. After studying the models, it was found that the explanation of the response variable was low according to the Combine tests for all samples.

Keywords: NFL, Combine, Modeling, Regression.

Lista de Tabelas

1	Banco de dados Combine	15
2	Banco de dados da <i>Pro Football Focus</i>	15
3	Tamanhos amostrais dos bancos de dados finais para modelagem	16
4	Tabela ANOVA	22
5	Tabela de medidas de posição e resumo	28
6	Tabela com p-valores para o teste de correlação de Spearman entre a variável resposta e as independentes	45
7	Análise do modelo inicial de regressão linear múltipla	48
8	Seleção stepwise com critério de informação AIC	48
9	Análise do modelo final de regressão linear múltipla das <i>Skill positions</i> . .	49
10	Tabela ANOVA do modelo final para a amostra geral	49
11	Valores do modelo final geral de regressão linear múltipla	50
12	P-valores dos testes de pressupostos para modelo de regressão linear múltiplo final	51
13	Valores de VIF de cada variável explicativa	52
14	Análise do modelo inicial de regressão linear múltipla para os <i>Quarterbacks</i>	53
15	Análise do modelo inicial de regressão linear múltipla para os <i>Running backs</i>	54
16	Análise do modelo inicial de regressão linear múltipla para os <i>Wide receivers</i>	54
17	Seleção stepwise com critério de informação AIC para os <i>Quarterbacks</i> . .	55
18	Análise do modelo final de regressão linear múltipla para os <i>Quarterbacks</i> .	55
19	Tabela ANOVA do modelo final para os <i>Quarterbacks</i>	56
20	Seleção stepwise com critério de informação AIC para os <i>Running backs</i> . .	57
21	Análise do modelo final de regressão linear múltipla para os <i>Running backs</i>	57
22	Tabela ANOVA do modelo final para os <i>Running backs</i>	57
23	Seleção stepwise com critério de informação AIC para os <i>Wide receivers</i> . .	58
24	Análise do modelo final de regressão linear múltipla para os <i>Wide receivers</i>	58
25	Tabela ANOVA do modelo final para os <i>Wide receivers</i>	59
26	Valores dos modelos finais de regressão linear múltipla	59

27	P-valores dos testes de pressupostos para modelo de regressão linear múltiplo final	61
28	Valores de VIF de cada variável explicativa para os <i>Quarterbacks</i>	61
29	Valores de VIF de cada variável explicativa para os <i>Running backs</i>	61
30	Tabela de comparação de variáveis presentes nos modelos de regressão finais	64
31	Medidas de diagnóstico para as previsões dos modelos finais	64
32	AIC para os modelos de regressão	65

Lista de Figuras

1	Explicação da variável Grade pela Pro Football Focus	16
2	Histogramas e Boxplots para a variável desempenho	29
3	Histogramas e Boxplots para a variável desempenho por posição	29
4	Colunas de frequência das categorias de notas de desempenho	30
5	Colunas de frequência das categorias de notas de desempenho por posição .	30
6	Histogramas e Boxplots para a variável altura	31
7	Histogramas e Boxplots para a variável altura por posição	32
8	Histogramas e Boxplots para a variável peso	33
9	Histogramas e Boxplots para a variável peso por posição	33
10	Histogramas e Boxplots para a variável comprimento do braço	34
11	Histogramas e Boxplots para a variável comprimento do braço por posição	34
12	Histogramas e Boxplots para a variável tamanho da mão	35
13	Histogramas e Boxplots para a variável tamanho da mão por posição . . .	35
14	Histogramas e Boxplots para a variável 3-cone drill	36
15	Histogramas e Boxplots para a variável 3-cone drill por posição	36
16	Histogramas e Boxplots para a variável Shuttle	37
17	Histogramas e Boxplots para a variável Shuttle por posição	37
18	Histogramas e Boxplots para a variável altura do salto vertical	38
19	Histogramas e Boxplots para a variável altura do salto vertical por posição	38
20	Histogramas e Boxplots para a variável distância do salto horizontal	39
21	Histogramas e Boxplots para a variável distância do salto horizontal por posição	39
22	Histogramas e Boxplots para a variável tempo no tiro de 40 jardas	40
23	Histogramas e Boxplots para a variável tempo no tiro de 40 jardas por posição	40
24	Histogramas e Boxplots para a variável número de repetições no supino, com 102kg	41
25	Histogramas e Boxplots para a variável número de repetições no supino, com 102kg por posição	42

26	Dispersão da varável resposta pelas variáveis explicativas para as <i>Skill positions</i>	43
27	Dispersão da varável resposta pelas variáveis explicativas para os <i>Quarterbacks</i>	43
28	Dispersão da varável resposta pelas variáveis explicativas para os <i>Wide receivers</i>	44
29	Dispersão da varável resposta pelas variáveis explicativas para os <i>Running backs</i>	44
30	Mapa de calor com p-valores do teste de correlação de Spearman entre as variáveis independentes para Skill positions	45
31	Mapa de calor com correlação entre as variáveis independentes para os <i>Quarterbacks</i>	46
32	Mapa de calor com correlação entre as variáveis independentes para os <i>Wide receivers</i>	46
33	Mapa de calor com correlação entre as variáveis independentes para os <i>Running backs</i>	47
34	Pressupostos modelo geral	51
35	Análise gráfica de pontos influentes do modelo geral	52
36	Normalidade	60
37	Homocedasticidade	60
38	Independência	60
39	Análise gráfica de pontos influentes dos <i>Quarterbacks</i>	62
40	Análise gráfica de pontos influentes dos <i>Running backs</i>	63
41	Análise gráfica de pontos influentes dos <i>Wide receivers</i>	63

Sumário

1 Introdução	11
1.1 Objetivos	13
1.1.1 Objetivos Específicos	13
1.2 Metodologia.	13
1.3 Conjuntos de dados	14
1.3.1 Combine	14
1.3.2 PFF Grade	15
2 Referencial Teórico	17
2.1 Análise Exploratória	17
2.2 Testes de hipótese	18
2.3 Função relacional entre 2 variáveis	19
2.4 Modelo de Regressão Linear Simples	20
2.5 Modelo de Regressão Linear Múltiplo	20
3 Revisão bibliográfica	26
4 Resultados	27
4.1 Análise Descritiva dos dados	27
4.1.1 Medidas de posição e resumo	27
4.1.2 Análise univariada da variável resposta	29
4.1.3 Análise univariada das variáveis explicativas	31
4.1.4 Análise bivariada	43
4.1.5 Testes de correlação entre as variáveis	45
4.2 Modelo de regressão linear múltipla	47
4.2.1 Modelo das <i>Skill positions</i>	47
4.2.2 Modelos por posição	53
4.2.3 Diagnóstico e comparação dos modelos	63
5 Conclusão	66

1 Introdução

A NFL (*National Football League*) é a liga de futebol americano dos Estados Unidos, e, assim como as outras ligas de esporte estadunidenses (Basquete, Hóquei e Baseball), a NFL possui, para a entrada de jogadores, um processo chamado de *Draft* (KOZ; FRASER-THOMAS; BAKER, 2012), tal qual leva cerca de 250 jogadores à liga anualmente. Com esse propósito, todos os times da NFL recebem escolhas que correspondem à ordem inversa do seu desempenho na classificação no ano anterior ao *Draft* (dessa maneira, o pior time da liga em uma temporada fica com a melhor escolha no próximo *Draft*, ao passo que o campeão fica com a pior), com a finalidade de os piores times da liga terem acesso aos melhores talentos (BERRI; SIMMONS, 2011). Esse sistema de seleção se repete sete vezes, ao longo de sete rodadas, com cada uma das 32 franquias da NFL fazendo suas escolhas, tendo, uma média de 259 jogadores escolhidos por ano. Muitas vezes, times trocam jogadores já consolidados por escolhas que se transformarão em novas promessas da liga. Assim, os times da liga valorizam muito mais as escolhas mais altas (selecionar antes os jogadores), em relação às mais baixas, por escolherem os jogadores de maior destaque durante sua vida esportiva.

Dessa maneira, é de suma importância compreender que, para construção de seu time no esporte, objetivando o sucesso, é crucial a predição de potencial desses jogadores. Para tal, o processo que envolve as escolhas dos jogadores se trata, majoritariamente, da análise do desempenho deles, que praticaram o futebol americano universitário. A fim de reduzir as escolhas desperdiçadas em jogadores que não corresponderão às expectativas, são utilizados diversos métodos, a cargo do time que irá fazer a seleção. Portanto, existem diversos critérios que os times dão maior enfoque para a seleção dos jogadores, podendo ser: Selecionar um jogador com base no potencial futuro; escolher com base nos atributos físicos; pode-se também priorizar um jogador que já está consolidado no futebol universitário; ou ainda, optar por aspectos psicológicos e de liderança demonstrados nas entrevistas conduzidas com os atletas.

Para tal, antes do *Draft*, é conduzido um evento chamado *Combine*, que leva, anualmente, em média 327 jogadores oriundos das universidades para fazer uma bateria de testes físicos, medições e entrevistas com donos de times, fazendo sua propaganda, mostrando seus talentos, para assim, serem selecionados e atuarem profissionalmente no esporte que sempre sonharam (COOK et al., 2020). Já para os donos e os técnicos dos times, resta a grande tarefa de acompanhar todas as etapas de perto, a fim de não desperdiçarem sua escolha em um jogador que pode não contribuir para o sucesso de sua equipe. Para o desenvolvimento do estudo, os dados coletados são dos testes dos jogadores durante o evento, que fazem pesagem, medição de altura, tamanho da mão e dos braços. Além de testes como contagem de repetições no supino com 102kg, tempo de tiro de 40

jardas (aproximadamente 36,6 metros), dois tipos de circuitos de agilidade, para os quais medem-se os respectivos tempos de conclusão, medição da distância de pulo horizontal e vertical, e também exercícios de cada posição promovidos para análise de mecânicas do jogador (KUZMITS; ADAMS, 2008).

Dado este prognóstico, existem na liga diversos estigmas que são levados em conta quando times vão selecionar seus futuros jogadores. Dentre esses, pode-se citar, por exemplo, que muitos analistas, ao estudarem um jogador da posição *Quarterback*, ou seja, quem lança a bola, consideram o tamanho da mão e a altura deste jogador como fatores fundamentais. Ou também, para jogadores de defesa, muitas vezes o comprimento de seus braços é visto como um fator importante para seu sucesso como profissional. Além disso, para selecionar recebedores, muitos levam em conta suas respectivas alturas (TERAMOTO; CROSS; WILLICK, 2016), velocidades e saltos verticais, dados os comportamentos dessas variáveis ao analisar atletas de elite nessa posição (HEDLUND, 2018).

Além disso, existem posições que são consideradas *Skill positions*, nas quais se considera que há uma maior importância das habilidades dos jogadores com a bola que em suas qualidades físicas propriamente, tais posições são, *Quarterback* que passam a bola, *Running Back* que correm, e *Wide Receiver*, que são os recebedores (KUZMITS; ADAMS, 2008; BEAULIEU-JONES et al., 2017). Dadas essas informações, é válido questionar se tais posições realmente são posições de *skill*, ou se os fatores físicos são de fato determinantes no desempenho deles na liga profissional. Caso exista, quais métricas seriam então mais relacionadas a esse desempenho? Além disso, objetiva-se, um encorajamento dos times da NFL a levarem em consideração e avaliarem as diferentes medidas obtidas no combine para selecionar os jogadores dessas posições, de tal forma preterindo até algumas características de habilidade que podem depois ser desenvolvidas ou supridas, para assim obter um sistema que possa facilitar a predição da possibilidade de sucesso dos jogadores na liga (KUZMITS; ADAMS, 2008).

Assim, a NFL, com seu complexo sistema de recrutamento e seleção de talentos, oferece uma visão única sobre como as habilidades atléticas e técnicas se entrelaçam na construção de equipes vencedoras, especialmente para as *Skill Positions*. Por conseguinte, a fim de entender melhor esse fenômeno, é fundamental explorar os dados e as métricas associadas a essas posições e como elas se traduzem em sucesso na liga, e a partir da análise detalhada desses aspectos, pode-se obter informações valiosas para aprimorar a estratégia de recrutamento e seleção na NFL.

1.1 Objetivos

Como principal objetivo do trabalho, será levada em conta a análise de se há, a partir dos dados dos anos em estudo, uma tendência observada, demonstrando uma relação causal, na qual o desempenho de um jogador de *Skill position* na liga é uma consequência de seus atributos físicos.

1.1.1 Objetivos Específicos

- Realizar uma análise exploratória dos dados de ambas as fontes, para resumir suas características;
- Ajustar modelos de regressão linear a fim de estudar a relação entre os dados;
- Verificar se após a modelagem, os pressupostos foram atendidos;
- Realizar testes para validar e analisar os modelos propostos.

1.2 Metodologia

Para esse estudo ser desenvolvido, foram levados em consideração os dados de jogadores que entraram na NFL entre 2007 e 2022, juntamente às suas medições, características físicas e medidas de atletismo.

Dessa maneira, foram coletados dois tipos de bancos de dados: o primeiro, referente aos testes físicos de cada jogador, desde a edição do *Combine* de 2007 até 2018, com todos os atletas participantes no evento, e também, dados que contém as notas dos jogadores de 2007 a 2022 pela PFF (*Pro Football Focus*). Assim, foram cruzadas essas fontes a fim de obter informação sobre os jogadores da NFL em dois âmbitos, antes do *Draft*, ou seja, durante o combine, e após o *Draft*, que corresponde ao desempenho dos jogadores realmente jogando profissionalmente na liga. Portanto, será feita a análise de dados utilizando técnicas estatísticas, sendo elas a organização e análise crítica dos dados, análise descritiva, uso de regressão para predição, análise de seus pressupostos e validação do modelo (NETER et al., 1996).

Primeiramente, após obtenção dos bancos, como citado anteriormente, estes foram analisados, excluindo-se as variáveis que não eram de interesse para o estudo, e depois, uma mescla, utilizando três variáveis chave para garantir que não haveria perda ou confronto de informações, sendo elas: Nome do jogador, sua posição e o ano em que foi selecionado no *Draft*.

Além disso, foi feita uma divisão dos jogadores no grupo posicional já comentado,

com base no banco de dados dos testes do *Combine*, a partir da separação de jogadores de *Skill positions*, que é uma nomenclatura bem comum no âmbito do esporte, cuja definição é de posições em que as habilidades dos jogadores são levadas mais em consideração que suas características físicas (BEAULIEU-JONES et al., 2017), de tal maneira, sendo mais relevante como esse jogador passa, dribla ou recebe a bola. Tal separação que está também presente no trabalho de Kuzmits e Adams (2008). Dessa forma, o grupo será, portanto, de jogadores que ocupam as posições *Wide Receiver*, *Running Back* e *Quarterback*.

Após isto, será feita a análise descritiva dos dados, a fim de descrever e resumir características do banco, além de detecção de valores discrepantes. Para melhor visualização, organização e identificação de tendências, facilitando portanto no entendimento dos dados em estudo (MORETTIN; BUSSAB, 2017).

Dessa forma, ajusta-se um modelo de regressão linear múltipla abarcando todos os jogadores de *Skill position* da amostra, com o objetivo de traçar uma relação entre os dados, a fim de examinar a possibilidade de estabelecer predição da variável resposta, ou seja, o desempenho dos jogadores, em relação às variáveis explicativas, os atributos físicos destes jogadores, para os jogadores de posições de habilidade. A técnica em questão torna necessária outro passo na análise, que é a análise dos resíduos do modelo, para checagem de seus pressupostos, e por último usar a amostra de teste para validação da capacidade preditiva do modelo ajustado (MONTGOMERY; PECK; VINING, 2021).

Além disso, ajusta-se as três posições de habilidade individualmente, a fim de aumentar o ajuste e capacidade preditiva do modelo para cada uma delas, e estudar os diferentes comportamentos da variável de desempenho em relação a cada um dos testes físicos.

1.3 Conjuntos de dados

1.3.1 Combine

Uma vez obtidos os dados do *Combine* (SUNDAY, 2019), se tratando de uma amostra com os atletas convidados para o programa de 2000 a 2018, tais quais possuíam 17 variáveis distintas, sendo elas: Nome do participante, Posição, Altura (em polegadas), Peso (em libras), Comprimento dos Braços (em polegadas), Tamanho das mãos (em polegadas), Tempo em segundos da corrida de 40 jardas, Altura do salto vertical (em polegadas), Número de repetições com 102kg no supino, Distância do salto horizontal (em polegadas), Tempo em segundos de percurso de agilidade entre 3 cones em formato de "L", com 4,3 metros entre cada um deles, Tempo em segundos de percurso com 16,3 metros, Ano do *Combine* que o jogador atendeu, ID do participante, Time que selecionou

o jogador no *Draft*, Rodada de seleção, e Número da escolha em que o jogador foi selecionado em seu draft (ROBBINS, 2010). As variáveis selecionadas para o estudo estão com o seguinte nome no banco de dados original:

Tabela 1: Banco de dados Combine

Nº	Variável	Nome	Descrição
1	Y	Grade	Nota de desempenho
2	X_1	Ht	Altura
3	X_2	Wt	Massa
4	X_3	Hand	Comprimento da mão
5	X_4	Arm	Comprimento do braço
6	X_5	Forty	Tempo no tiro de 40 jardas
7	X_6	Vert	Altura do salto vertical
8	X_7	Broad	Distância do salto horizontal
9	X_8	Cone	Tempo no percurso de 3 cones
10	X_9	Shuttle	Tempo no percurso em L
11	-	Bench	Repetições no supino
12	-	Name	Nome do jogador
13	-	Position	Posição do jogador (no <i>College</i>)
14	-	Year	Ano do <i>Combine</i>

1.3.2 PFF Grade

Já os dados de notas, com fonte em PFF (2023), possuem cerca de 50 variáveis distintas, mas dentre elas, foram utilizadas para o estudo, as seguintes: Nome do participante, Ano cujas notas foram obtidas, de tal forma que alguns os jogadores têm mais de uma observação neste banco, e as notas gerais do jogador.

Tabela 2: Banco de dados da *Pro Football Focus*

Nº	Nome	Descrição
1	Name	Nome do jogador
2	Year	Ano da temporada
3	Offensive Grade	Nota de desempenho no ataque

De acordo com a PFF (2023), as notas de desempenho têm como método uma escala que está no intervalo de $(-2, 2)$, e que varia em 0.5. Dessa forma, o 0 é uma jogada neutra, média, ou esperada, e a cada grau de melhora da jogada por parte do jogador especificamente, acrescenta 0.5 à sua nota. Já para cada grau de piora, ou seja, uma

jogada negativa, decresce-se 0.5 de sua nota, assim como se segue a imagem explicativa do site:

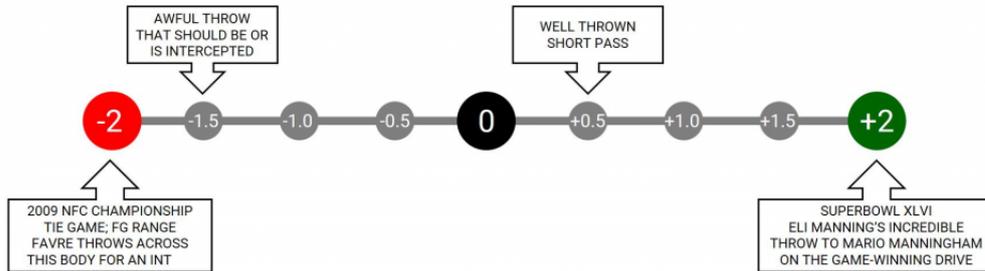


Figura 1: Explicação da variável Grade pela Pro Football Focus

De tal maneira, como explicado, para cada jogada de um jogador, ele recebe uma pontuação, com o 0 sendo uma jogada média, pela imagem, tem-se a explicação de que, uma jogada que some 0.5 à nota de um jogador (no exemplo, da posição *Quarterback*), pode se tratar de um bom passe curto, considerado uma jogada mais fácil. Além disso, para uma jogada que rende 2 pontos à nota do jogador, se trata de um passe magnífico, que tem um grau de dificuldade extremo e tem grande valor sobre a vitória do seu time. Em contraponto, uma decisão muito ruim, de alta influência na derrota do time, exemplificando uma jogada que vale -2 pontos para a nota. Soma-se portanto as jogadas ao decorrer de cada partida, calculando qual a nota que um jogador teve ao final, em uma escala de 0 a 100.

Para a construção do modelo, foi retirada a variável *Bench*, de número de repetições de supino, pois para o grupo posicional estudado, apresentava muitos valores inexistentes, por ser um teste que não é recorrente nas posições em estudo, em especial para os *Quarterbacks*, os quais na amostra, nenhum havia feito este teste, com isso, obtendo um banco de dados contendo 320 jogadores. Além disso, foi feita a divisão do banco de dados para uma amostra de treino, contendo 80% das observações, e o de validação contendo os outros 20% a fim de obter a informação de se o modelo final é preciso ou não. Os tamanhos amostrais de cada subgrupo são os que se seguem:

Tabela 3: Tamanhos amostrais dos bancos de dados finais para modelagem

Sub-grupo	Amostra	
	Treino	Teste
Geral	256	64
<i>Quarterbacks</i>	40	9
<i>Running backs</i>	87	21
<i>Wide receivers</i>	131	32

2 Referencial Teórico

2.1 Análise Exploratória

A análise descritiva de dados fornece algumas medidas de posição e variabilidade, como a média e variância, por exemplo. Além disso, tomando como base Tukey et al. (1977), esta análise utiliza principalmente técnicas gráficas, em oposição a apenas resumos numéricos. Dessa maneira, faz-se a análise de boxplots, histogramas e gráficos de dispersão na análise exploratória. Portanto, sumários não só devem ser obtidos, mas uma análise exploratória de dados não deve se limitar a calcular apenas essas medidas. (MORETTIN; BUSSAB, 2017).

Dessa maneira, tendo em vista os livros em questão, é possível obter uma maneira mais didática e moderna de se prosseguir com a análise descritiva dos dados. Esta é então, classificada em três capítulos, que se tratam de:

Resumo de dados

A distribuições de frequências objetiva conhecer o comportamento da variável, analisando a ocorrência de suas observações. A frequência absoluta representa o número bruto de observações em cada classe da variável, e a frequência relativa, que é a proporção de cada classe em relação ao total de observações, ou seja, $f_i = n_i/n$. Esse resumo pode ser utilizado tanto em variáveis discretas, quanto contínuas, contudo, apenas a partir da categorização das contínuas.

Análise gráfica para este resumo pode ser feita por gráficos de pizza ou de colunas, apresentando ambas as frequências relativa e absoluta para cada classe.

Medidas resumo

As medidas podem ser a média aritmética, que pode ser expressa por

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Além disso, pode-se expressar a variância por

$$Var(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Analogamente, o desvio-padrão é a raiz quadrada da variância, sendo então

$$dp = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2} = \sqrt{Var(X)}.$$

A mediana, é o quantil 50%, ou seja, que corresponde ao valor central da série, quando seus dados estão em ordem crescente. Caso o número de observações seja par, a mediana corresponde à média aritmética das duas observações centrais. Pode-se também obter os quantis dos dados. Para uma análise descritiva, o comum é utilizar o 1º e o 3º quartis, que, analogamente à mediana, respectivamente indicam o 1º quarto e o 3º quarto da amostra (25% e 75%).

Por último, é de praxe calcular a amplitude de classe, que corresponde à diferença do maior valor, para o menor valor, expresso portanto por

$$\text{Amplitude} = \text{valor máximo} - \text{valor mínimo}.$$

Análise Bidimensional

Para analisar a associação entre duas variáveis quantitativas, pode-se não só utilizar gráfico de dispersão, para visualização da distribuição de uma variável em relação à outra, mas também, acrescentar a utilização do coeficiente de correlação, que é expresso por

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{dp(X)dp(Y)},$$

sendo que

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}.$$

2.2 Testes de hipótese

Teste de correlação de Spearman

O teste é uma medida não paramétrica, ou seja, testa para dados não normais a correlação, sob hipóteses

$$\begin{cases} H_0 : \text{A correlação entre as variáveis é estatisticamente igual a zero} \\ H_1 : \text{Há correlação estatística entre as variáveis} \end{cases}$$

Teste de Shapiro-Wilk

Teste de hipótese cujo objetivo é verificar a normalidade da variável em estudo, dessa forma, trazendo as hipóteses

$$\begin{cases} H_0 : \text{A amostra segue a distribuição Normal} \\ H_1 : \text{A amostra não segue a distribuição Normal} \end{cases}$$

Teste de Breusch-Pagan

Teste para a homocedasticidade de um modelo de regressão linear, com hipóteses

$$\begin{cases} H_0 : \text{Os resíduos são homocedásticos} \\ H_1 : \text{Os resíduos não são homocedásticos} \end{cases}$$

Teste de Durbin-Watson

Teste de hipótese para checar independência dos resíduos, com as hipóteses sendo

$$\begin{cases} H_0 : \text{Os resíduos são independentes} \\ H_1 : \text{Os resíduos não são independentes} \end{cases}$$

2.3 Função relacional entre 2 variáveis

Escrita pela fórmula matemática a seguir,

$$Y = f(X),$$

na qual X é uma variável independente e Y uma variável dependente. Dessa maneira, pode-se exemplificar para um problema simples, como quando o salário de um chefe de uma empresa (Y), depender do salário de seu subordinado (X), cuja relação é de que o salário do chefe sempre será o dobro do salário de seu subordinado, a função relacional poderá ser escrita por:

$$Y = 2X.$$

2.4 Modelo de Regressão Linear Simples

A partir do conceito abordado sobre a função relacional, um modelo de regressão linear simples é aquele em que há apenas uma variável dependente e uma independente, analogamente ao exemplo anterior.

Declaração Formal do Modelo:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

em que

- Y_i é a variável resposta na i -ésima observação;
- β_0 e β_1 são parâmetros;
- ϵ_i é um erro aleatório com média $E(\epsilon_i) = 0$. ϵ_i e ϵ_j são não correlacionados, dessa maneira, sua covariância $\sigma(\epsilon_i, \epsilon_j) = 0, \forall i, j; i \neq j$.

2.5 Modelo de Regressão Linear Múltiplo

Tipicamente, pode-se buscar, contudo, mais de uma variável independente, pois, sabe-se que há mais de uma variável importante que influencia na resposta. Dessa forma, de maneira similar, pode-se obter um modelo de primeira ordem com duas ou mais variáveis explicativas (NETER et al., 1996).

Então, considerando o caso em que há $p - 1$ variáveis independentes X_1, \dots, X_{p-1} ; é feita uma generalização para um modelo de regressão linear, que assume termos de erro normais, assumindo então a função

$$\hat{Y}_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_{p-1} X_{ip-1} + \epsilon_i,$$

em que

- Y_i é a variável resposta na i -ésima tentativa;
- $\beta_0, \beta_1, \dots, \beta_{p-1}$ são parâmetros;
- X_{i1}, \dots, X_{ip-1} são constantes conhecidas;
- ϵ_i seguem $N(0, \sigma^2)$ independentes;
- $i \in 1, \dots, n$.

Estimativa dos coeficientes de regressão

O método mais utilizado para tal é o método dos mínimos quadrados (MONTGOMERY; PECK; VINING, 2021), e pode ser calculado por

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

sendo que

- $\hat{\beta}$ é a estimativa do coeficiente β ;
- X é a matriz que contém os valores das variáveis independentes;
- y é o vetor da variável resposta.

Análise de variância (ANOVA)

A análise de variância é uma técnica estatística fundamental utilizada na regressão múltipla, sendo utilizada para avaliar a significância global do modelo de regressão e a importância dos coeficientes individuais (NETER et al., 1996).

Contém as seguintes fontes de variação:

- **Regressão:** Representa a variação explicada pelo modelo de regressão. Se trata então da influência conjunta de todas as variáveis independentes na dependente.
- **Erro:** Representa a variação que não é explicada pelo modelo de regressão. Ela inclui todas as outras fontes de variação na variável dependente que não são explicadas pelas variáveis independentes do modelo.

É comum realizar testes de hipóteses individuais para cada coeficiente de regressão β_j para determinar se cada variável independente é significativa no modelo. Os testes de hipóteses para cada β_j são conduzidos sob:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \dots = \beta_n = 0 \\ H_1 : \exists \beta_j \neq 0, \forall j = 1, \dots, n \end{cases}$$

A tabela ANOVA fornece estatísticas de teste, como o valor-p associado ao teste F global e os valores-p associados aos testes t individuais para os coeficientes de regressão. Com base nos valores p, podemos determinar se o modelo de regressão é significativo e quais variáveis independentes individuais são significativas para a previsão da variável dependente. A tabela tem o formato que se segue:

Fonte de Variação	Graus de Liberdade	Soma de Quadrados	Quadrado Médio	Estatística F
Regressão	k-1	$\sum(\hat{Y}_i - \bar{Y})^2$	$\frac{SQReg}{k-1}$	$\frac{QMReg}{QMRes}$
Resíduos	n-k	$\sum(Y_i - \hat{Y}_i)^2$	$\frac{SQRes}{n-k}$	
Total	n-1	$\sum(Y_i - \bar{Y})^2$		

Tabela 4: Tabela ANOVA

Sendo que:

- $SQReg$: Soma de quadrados da regressão;
- $SQRes$: Soma de quadrados dos resíduos;
- $QMReg$: Quadrado médio da regressão;
- $QMRes$: Quadrado médio dos resíduos.

Fator de inflação de variância (VIF)

A multicolinearidade é um problema sério que atrapalha a utilidade de um modelo de regressão é quando não há independência entre as variáveis explicativas (MONTGOMERY; PECK; VINING, 2021).

O VIF é um método para detecção da Multicolinearidade, que mede o quanto de variância dos coeficientes estimados da regressão são inflados, quando comparado a um modelo em que os coeficientes não apresentam correlação linear. O VIF é definido por

$$VIF_j = \frac{1}{1 - R_j^2},$$

tal que:

- VIF_j é o VIF para o j -ésimo coeficiente de regressão;
- R_j^2 é o coeficiente de determinação múltiplo.

A sua interpretação é de que $VIF_j \geq 10$ é um problema sério de multicolinearidade para a regressão (MONTGOMERY; PECK; VINING, 2021).

Pontos influentes

Ao identificar *outliers* que influenciam os valores de X e de Y, devemos considerá-los influentes se eles causarem diferenças significativas na função de regressão. Dessa maneira, com a finalidade de identificar tais observações, pode-se prosseguir com a análise das técnicas que se seguem (NETER et al., 1996):

DFFITS

$$(DFFITS)_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{QMRes_{(i)}h_{ii}}}.$$

Sua interpretação é que um valor é considerado influente quando:

$$\begin{cases} |DFFITS_i| \geq 1, \text{ Para amostras pequenas e médias} \\ |DFFITS_i| \geq 2\sqrt{\frac{p}{n}}, \text{ Para amostras grandes} \end{cases}$$

Distância de Cook

$$D_i = \frac{\sum(\hat{Y}_j - \hat{Y}_{j(i)})^2}{p \cdot QMRes},$$

em que p é o número de coeficientes de regressão. Sua interpretação é que um valor é considerado influente quando:

$$D_i \geq 0,8$$

DFBETAS

$$DFBETAS_{(i)} = \frac{\hat{\beta} - \hat{\beta}_{k(i)}}{\sqrt{QMRes_{(i)}c_{kk}}}.$$

Sua interpretação é que um valor é considerado influente quando:

$$\begin{cases} |DFBETAS_{j(i)}| \geq 1, \text{ Para amostras pequenas e médias} \\ |DFBETAS_{j(i)}| \geq \frac{2}{\sqrt{n}}, \text{ Para amostras grandes} \end{cases}$$

Métodos de seleção de variáveis

A fim de calcular quais variáveis devem ser incluídas ou excluídas de um modelo, são usadas as seguintes técnicas estatísticas (MONTGOMERY; PECK; VINING, 2021):

Critério de informação de Akaike (AIC_p)

O critério é usado para fazer a seleção de variáveis e comparação de modelos, de forma a adicionar penalidade para a adição de variáveis, da forma que se segue (AKAIKE, 1974):

$$AIC_p = n \ln(SSE_p) - n \ln(n) + 2p.$$

Dessa forma, selecionando o modelo que apresentar o menor AIC.

Backward

O modelo inicial é feito com todas as variáveis explicativas, e a partir dele, é testada a permanência de cada uma no modelo, retirando as variáveis e analisando o impacto da retirada dela no modelo, visando a redução dos critérios de seleção.

Forward

Oposto ao método *Backward*, inicia-se com o modelo contendo apenas o intercepto, e cada passo consiste na adição de uma variável independente que reduz o valor dos critérios.

Stepwise

É a junção dos dois métodos anteriores, de forma que a cada passo, são adicionadas e retiradas variáveis do modelo, até sobrar um possível modelo final, que obteve menor valor de AIC ou outro critério escolhido.

Medidas de diagnóstico do modelo

PRESS

O PRESS é uma medida do quão bem um modelo vai performar prevendo dados novos (MONTGOMERY; PECK; VINING, 2021), ele pode ser expresso por:

$$PRESS = \sum (\hat{y}_i - y_i)^2 = \sum e_i^2,$$

em que:

- y_i é o valor real do banco de teste para a variável resposta na i -ésima observação;

- \hat{y}_i é a predição da variável resposta na i -ésima observação baseada no modelo de regressão ajustado;
- e_i são os resíduos do modelo

Coefficiente de determinação R^2

O coeficiente de determinação indica o quanto da variância dos dados é explicada pelo modelo. A partir do PRESS, é possível calcular o R^2 de um modelo (MONTGOMERY; PECK; VINING, 2021), e assim, chegar no seguinte resultado:

$$R^2 = 1 - \frac{\text{PRESS}}{\text{SQTot}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2},$$

em que:

- R^2 é o percentual de variação explicado por esse modelo;
- SQTot é a soma dos quadrados total da regressão;
- \bar{y} é a média das observações.

Erro quadrático médio

O erro quadrático médio mede a performance de um modelo. Para modelos em que os resíduos seguem distribuição normal, ele é mais apropriado que o erro médio absoluto (CHAI; DRAXLER, 2014), assim, ele tem a forma a seguir:

$$EQM = \sqrt{\frac{1}{n} \sum e_i^2}.$$

3 Revisão bibliográfica

Kuzmits e Adams (2008), Utilizando análise de correlação, não foi encontrada relação estatística entre os testes e a performance de jogadores, com excessão a testes de velocidade para *Running Backs*.

Berri e Simmons (2011) concluem baixa correlação entre as avaliações dos times sob os *Quarterbacks selecionados* e os seus desempenhos na liga.

Cook et al. (2020) modelaram, a partir de uma regressão linear múltipla, e concluíram que o *Combine* explicou apenas 2.6% da variância da resposta (performance média dos jogadores na NFL na temporada de calouro).

Teramoto, Cross e Willick (2016) estudaram se as medidas do Combine da NFL podem prever o desempenho futuro de *running backs* e *wide receivers* na NFL. Concluiu-se a partir de análise de regressão, correlação e de componentes principais, que o *Combine* tem valor na predição de performance nessas posições na NFL.

Hedlund (2018) prosseguiram um estudo focado nos jogadores de elite da liga, e concluíram q as características físicas são fundamentais no sucesso na NFL.

4 Resultados

4.1 Análise Descritiva dos dados

Após a manipulação e agrupamento dos dados, prossegue-se com a análise exploratória, dando ênfase à comparação entre os jogadores das *Skill positions*, para os quais o estudo é voltado (MORETTIN; BUSSAB, 2017), e as posições individualmente.

4.1.1 Medidas de posição e resumo

De início, é feita a análise das medidas de posição das variáveis dos testes físicos e da nota de desempenho dos jogadores, os números são os que se seguem:

Tabela 5: Tabela de medidas de posição e resumo

	1º Quartil	3º Quartil	Máximo	Média	Mediana	Mínimo	Abstenção	Desvio-Padrão
Ht Skill	178.4	188	201	183.4	182.9	167.3	0	5.72
QB	188.3	194.6	201	191.3	190.8	179.4	0	4.59
WR	180.3	188	198.4	184.2	184.5	170.2	0	5.67
HB	176.8	182.4	189.6	179.5	179.4	167.3	0	4.17
Wt Skill	89.81	100.24	120.2	95.05	95.25	70.76	0	6.52
QB	98.88	105.23	120.2	102.57	102.51	89.36	0	4.61
WR	87.09	95.71	108.86	91.43	91.17	70.76	0	6.08
HB	93.22	101.6	112.04	96.98	97.52	78.02	0	5.96
Arm Skill	78.74	82.55	91.44	80.55	80.64	71.12	88	3.11
QB	80.64	84.15	88.9	82.26	82.25	76.2	13	2.65
WR	79	83.09	91.44	81.14	81.28	71.12	40	3.15
HB	77.17	80.72	87	79.2	79.22	71.12	35	2.62
Cone Skill	6.81	7.09	7.56	6.95	6.95	6.42	193	0.19
QB	6.9	7.16	7.52	7.04	7.06	6.66	18	0.18
WR	6.76	7.03	7.4	6.89	6.9	6.42	88	0.17
HB	6.86	7.15	7.56	7	6.97	6.5	87	0.20
Hand Skill	22.86	24.76	27.94	23.8	23.83	20.65	75	1.26
QB	23.83	25.4	27.64	24.63	24.61	22.56	1	1.23
WR	23.19	24.68	27.94	23.82	23.83	20.65	40	1.30
HB	22.86	24.13	27.64	23.45	23.5	20.65	34	1.19
Shuttle	4.15	4.34	5.01	4.25	4.24	3.81	162	0.14
QB	4.2	4.41	4.53	4.3	4.31	4	16	0.13
WR	4.13	4.32	5.01	4.22	4.21	3.81	72	0.13
HB	4.18	4.37	4.67	4.27	4.25	3.93	74	0.15
Vert Skill	82.55	92.71	114.3	87.96	87.63	64.77	72	7.49
QB	74.93	85.09	97.79	80.13	78.74	64.77	14	7.40
WR	85.09	93.98	114.3	89.98	90.17	66.04	35	7.57
HB	82.55	92.71	109.22	88	87.63	71.12	23	7.39
Bench Skill	14	20	35	17.16	17	6	200	5.26
QB	-	-	-	-	-	-	77	-
WR	12	18	35	15.02	15	6	79	4.37
HB	16	22	32	19.54	19	8	44	4.97
Broad Skill	292.1	312.4	353.1	303	304.8	259.1	80	14.15
QB	274.3	294.6	320	285.5	284.5	259.1	14	15.23
WR	299.7	315	353.1	308.1	307.3	266.7	38	13.98
HB	292.1	309.9	342.9	302.5	302.3	264.2	28	14.09
Forty Skill	4.45	4.61	5.14	4.55	4.53	4.22	29	0.09
QB	4.68	4.89	5.14	4.8	4.81	4.47	4	0.15
WR	4.42	4.55	4.78	4.49	4.5	4.22	12	0.09
HB	4.46	4.6	4.84	4.53	4.54	4.26	13	0.09
Grade Skill	57.6	68.5	89	62.92	63.1	27	0	7.90
QB	51.35	69.75	89	59.85	60.4	27	0	14.12
WR	58.2	68.55	88.4	63.61	62.6	45.6	0	7.87
HB	57.7	68	87.1	63.14	63.8	39.9	0	7.96

4.1.2 Análise univariada da variável resposta

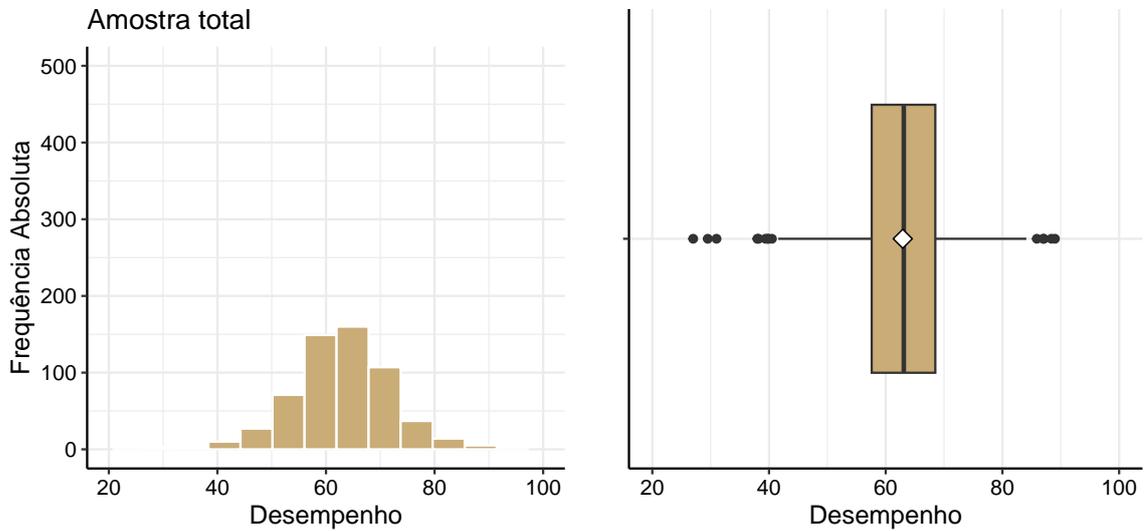


Figura 2: Histogramas e Boxplots para a variável desempenho

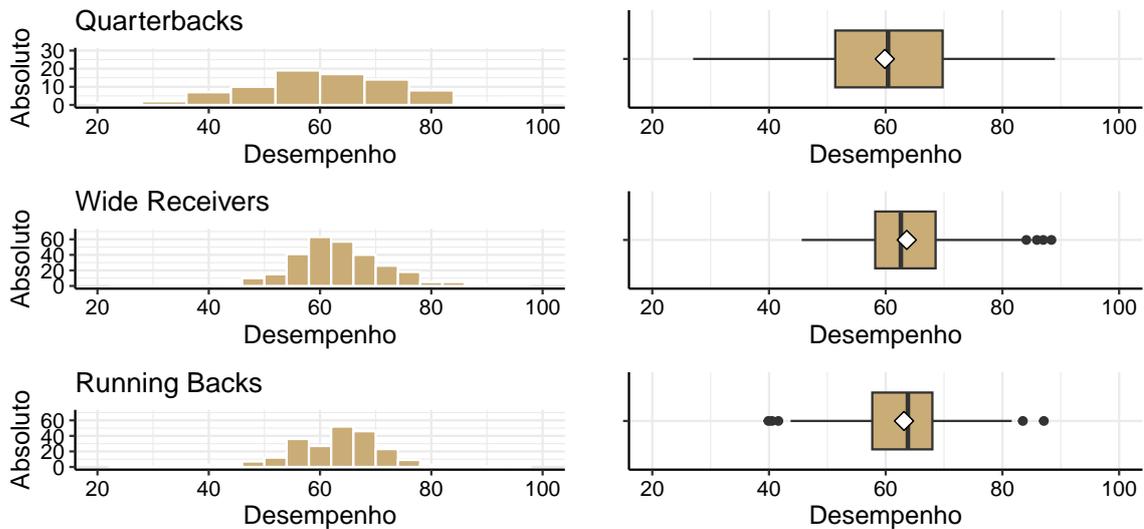


Figura 3: Histogramas e Boxplots para a variável desempenho por posição

Pode-se observar a maior proximidade das notas gerais em relação aos recebedores e corredores, o que pode apenas ser reflexo de suas amostras serem consideravelmente maiores que para os *Quarterbacks*, mas, apesar disso, ainda pode indicar a dificuldade que é de se jogar na posição mais importante do jogo.

Agora, pode-se visualizar as frequências da categorização da variável de notas dos jogadores:

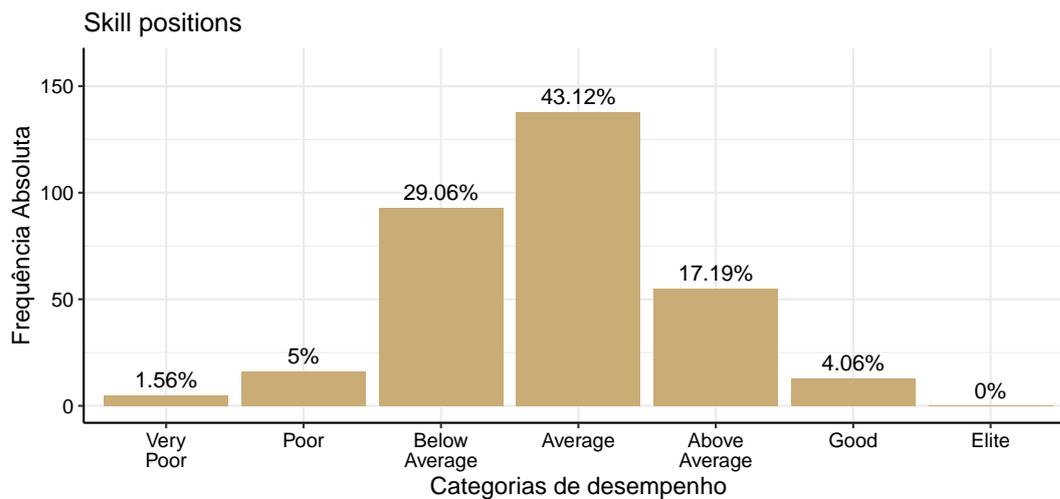


Figura 4: Colunas de frequência das categorias de notas de desempenho

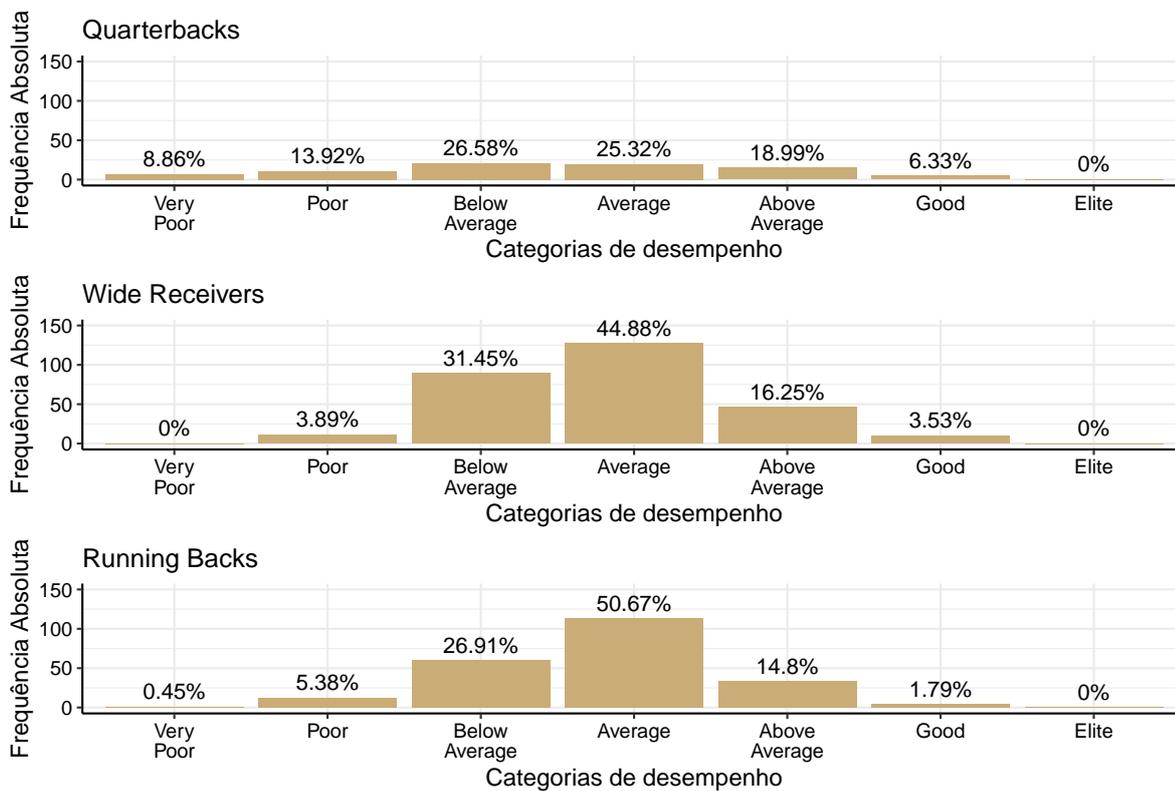


Figura 5: Colunas de frequência das categorias de notas de desempenho por posição

Assim, pode-se perceber como a *NFL* é uma liga que apresenta um grande grau de dificuldade, visto que, para a amostra em estudo, dentre os 320 jogadores, na média de nota de sua carreira inteira, existe mais que o dobro de jogadores na categoria *Below Average* que na *Above Average*.

Além disso, pode-se ver que ao estudar a variável de desempenho graficamente, não se nota muita diferença novamente entre as *Skill positions* e as duas posições de maior

amostra, contudo, para os *Quarterbacks*, pode-se perceber, retornando aos histogramas e comparando às frequências, que há uma grande variabilidade nas notas, pois na amostra em estudo os jogadores estão bem dispersos em relação à média, com mais jogadores na classificação boa de desempenho, mas com poucos na média e maior frequência na classe *Poor*, se comparado às outras posições.

4.1.3 Análise univariada das variáveis explicativas

Agora, é realizada a análise gráfica para as variáveis explicativas, prosseguindo com a análise dos histogramas e boxplots para cada amostra em estudo:

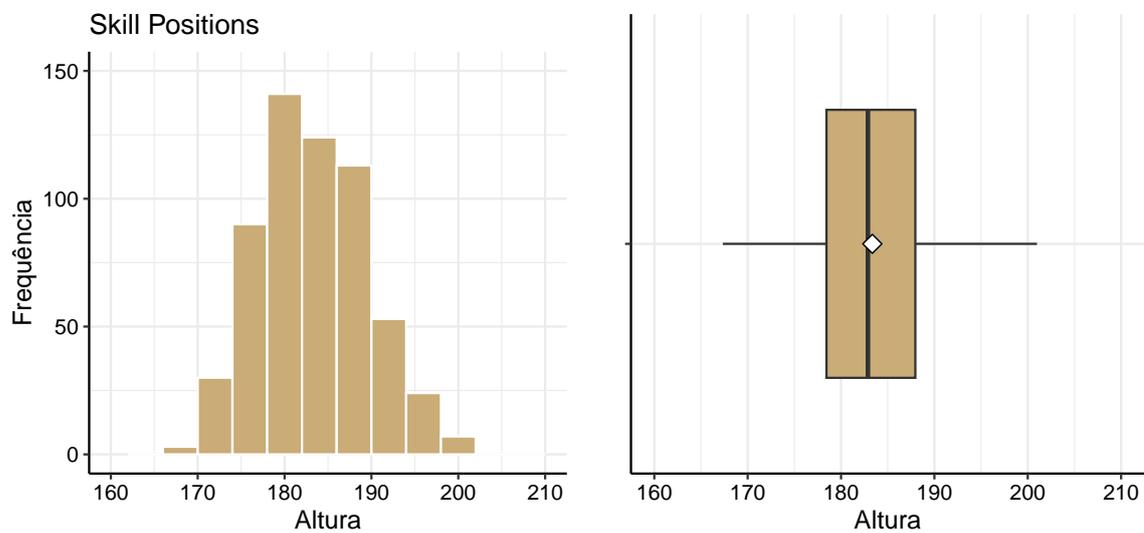


Figura 6: Histogramas e Boxplots para a variável altura

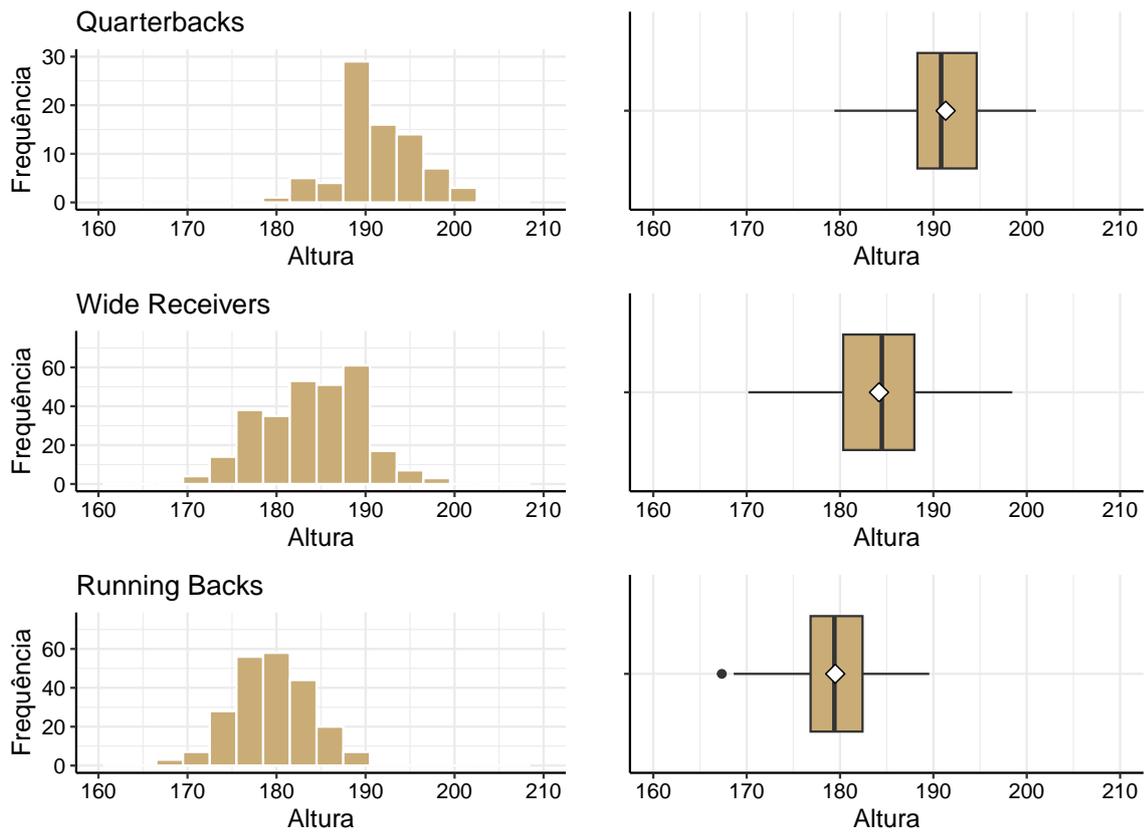


Figura 7: Histogramas e Boxplots para a variável altura por posição

A partir dos gráficos, pode-se notar que a tendência para os jogadores das *Skill positions* é de estarem mais concentrados entre 1,80 e 1,90 de altura, com sua mediana próxima dos 1,85.

É visto que a maior parte dos *Running backs* estão mais distribuídos em torno de 180cm, de modo que os *Wide receivers* em 185cm e os *Quarterbacks* estão em maioria além dos 190.

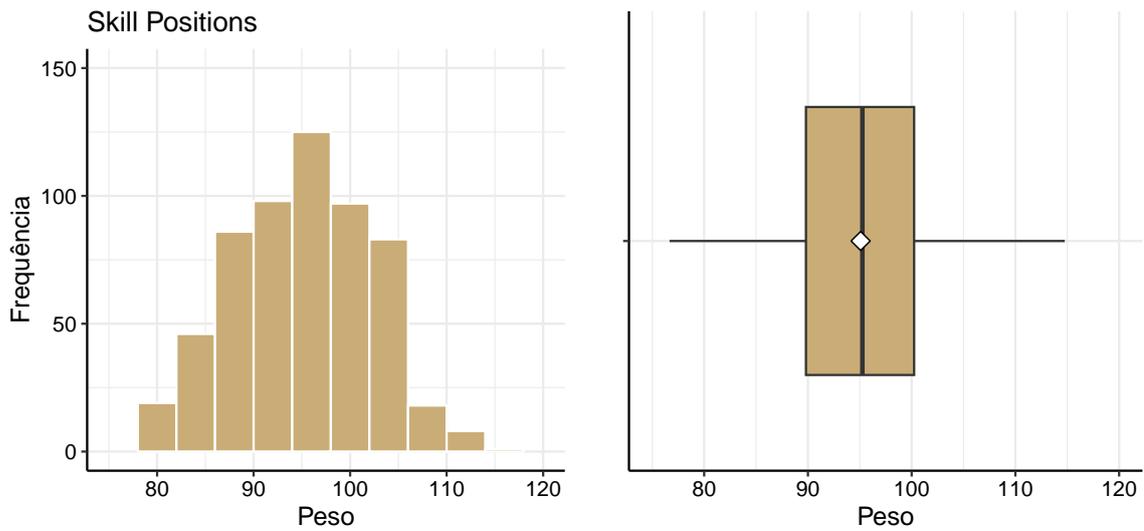


Figura 8: Histogramas e Boxplots para a variável peso

A partir da amostra, percebe-se que os jogadores de habilidade estão bem distribuídos em volta da média para a variável peso, que é de 95kg.

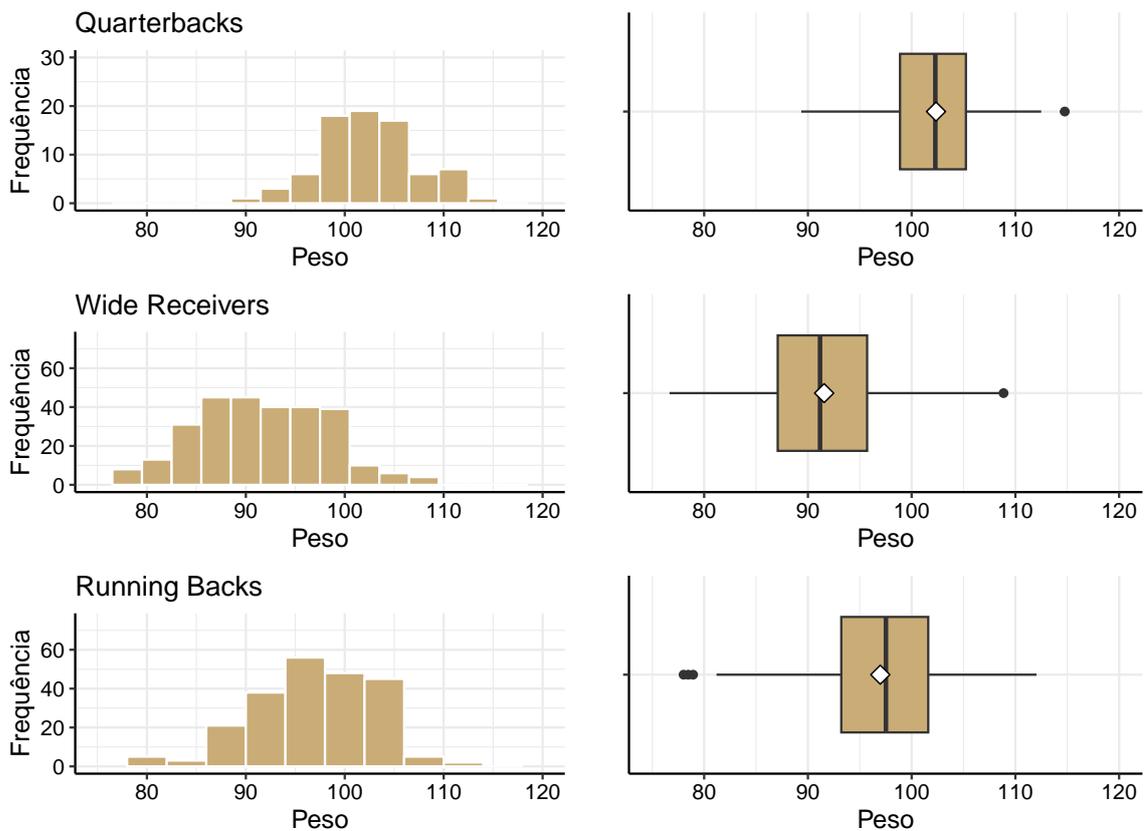


Figura 9: Histogramas e Boxplots para a variável peso por posição

Pode-se notar que os corredores são, no geral, jogadores mais pesados que os recebedores, e os passadores são os que testaram com maior peso, no geral.

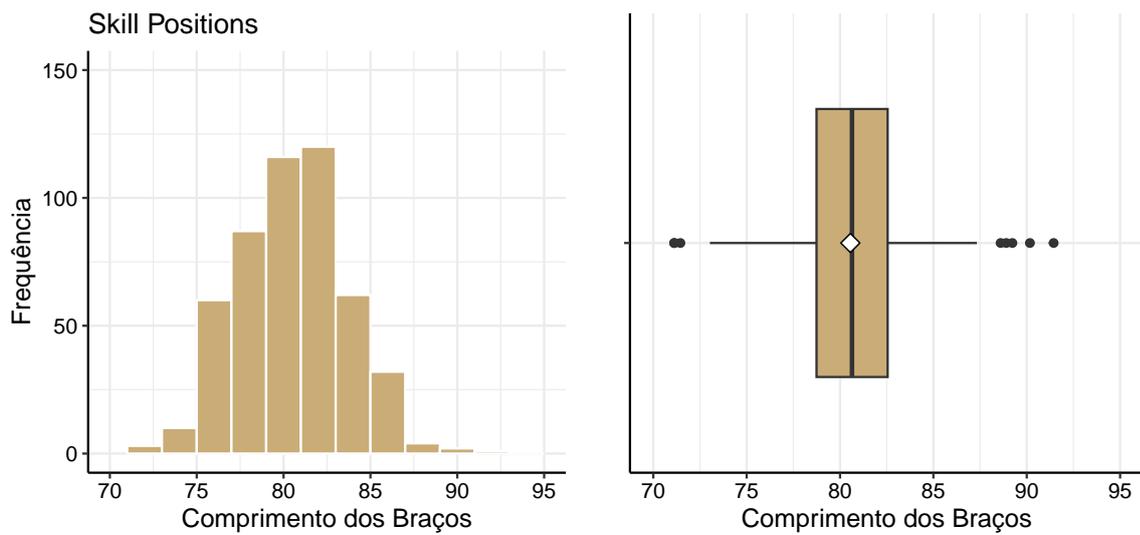


Figura 10: Histogramas e Boxplots para a variável comprimento do braço

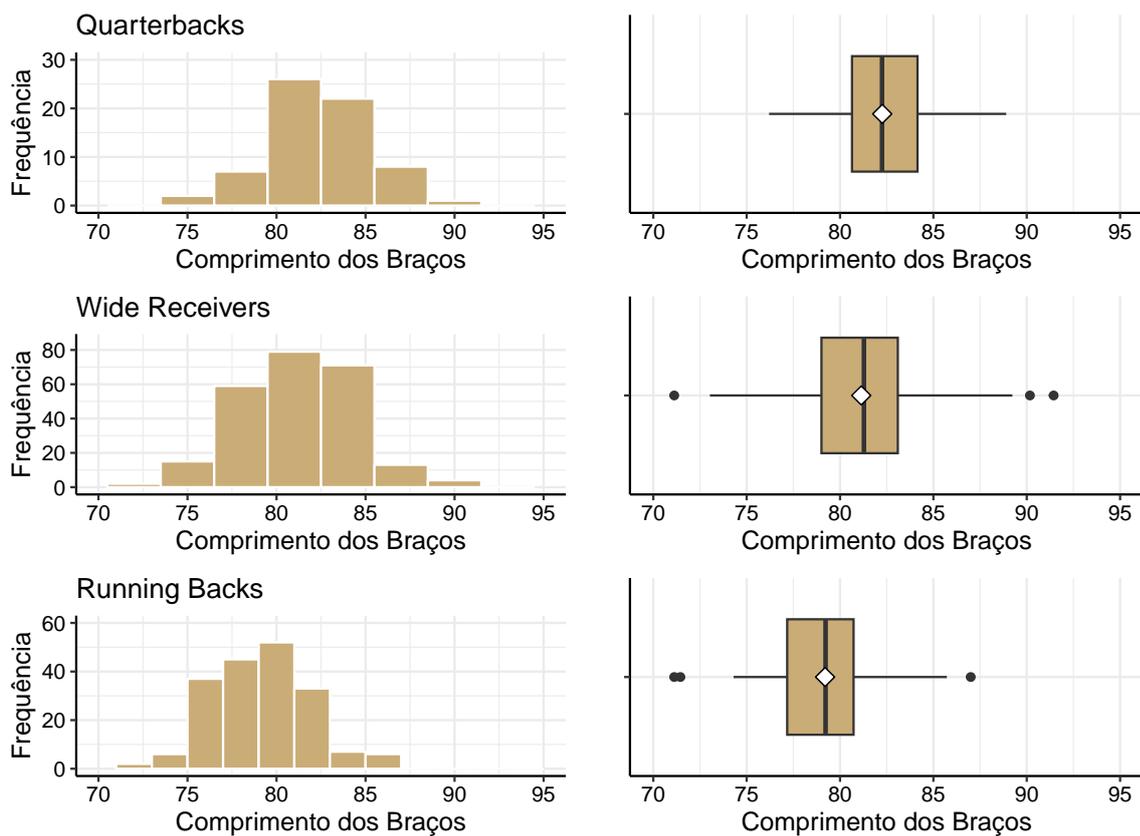


Figura 11: Histogramas e Boxplots para a variável comprimento do braço por posição

É perceptível que os recebedores possuem tamanho de braço, geralmente, maior que dos corredores, e um pouco abaixo até da amostra dos passadores.

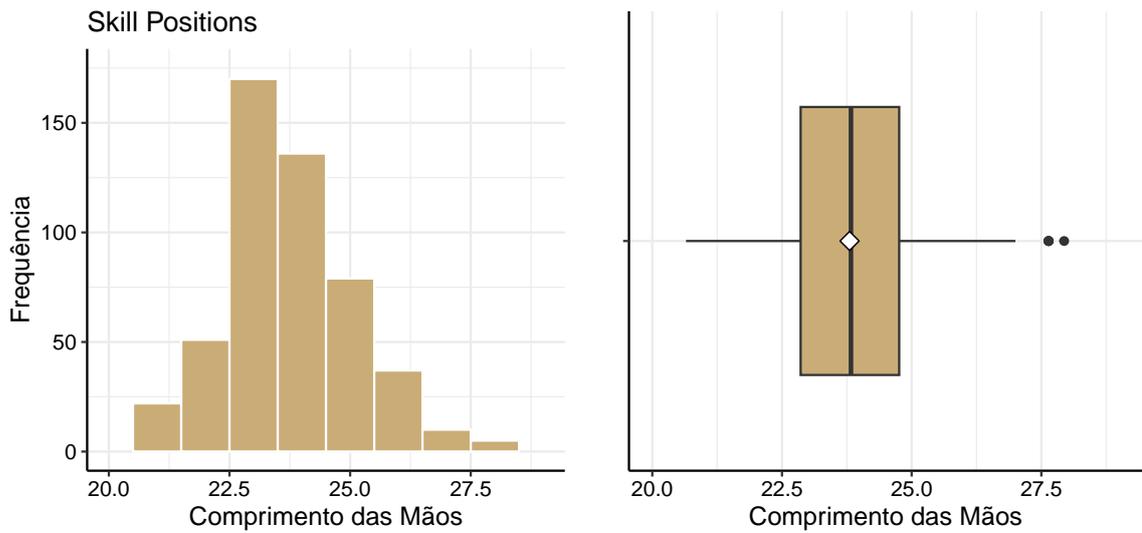


Figura 12: Histogramas e Boxplots para a variável tamanho da mão

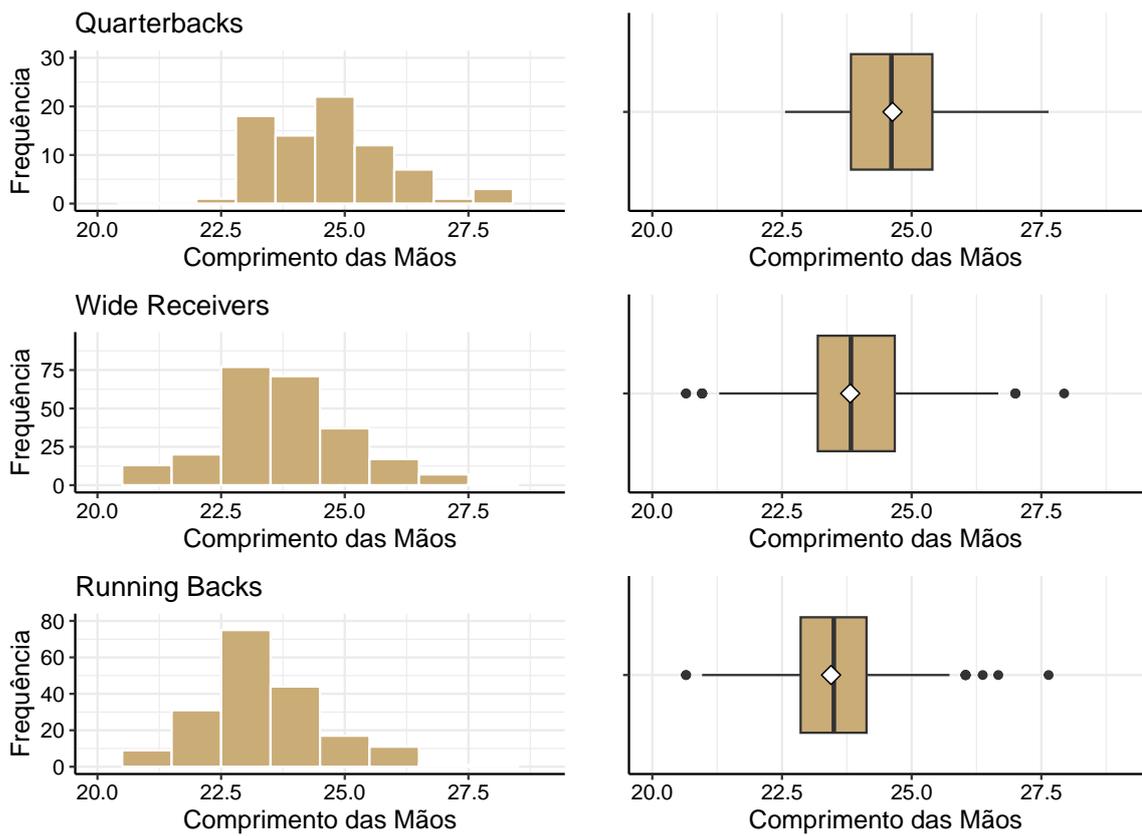


Figura 13: Histogramas e Boxplots para a variável tamanho da mão por posição

Aqui também é notável o comportamento semelhante, e vale observar a partir das variáveis altura, peso, tamanho das mãos e dos braços, que podem ser variáveis com grande correlação, o que vale observar com cuidado no momento de modelar os dados.

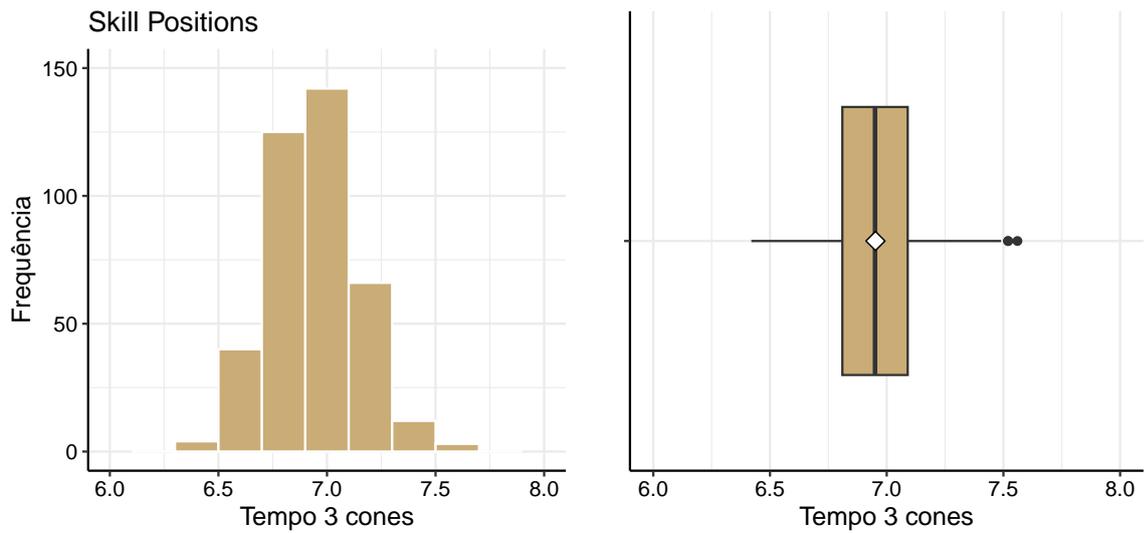


Figura 14: Histogramas e Boxplots para a variável 3-cone drill

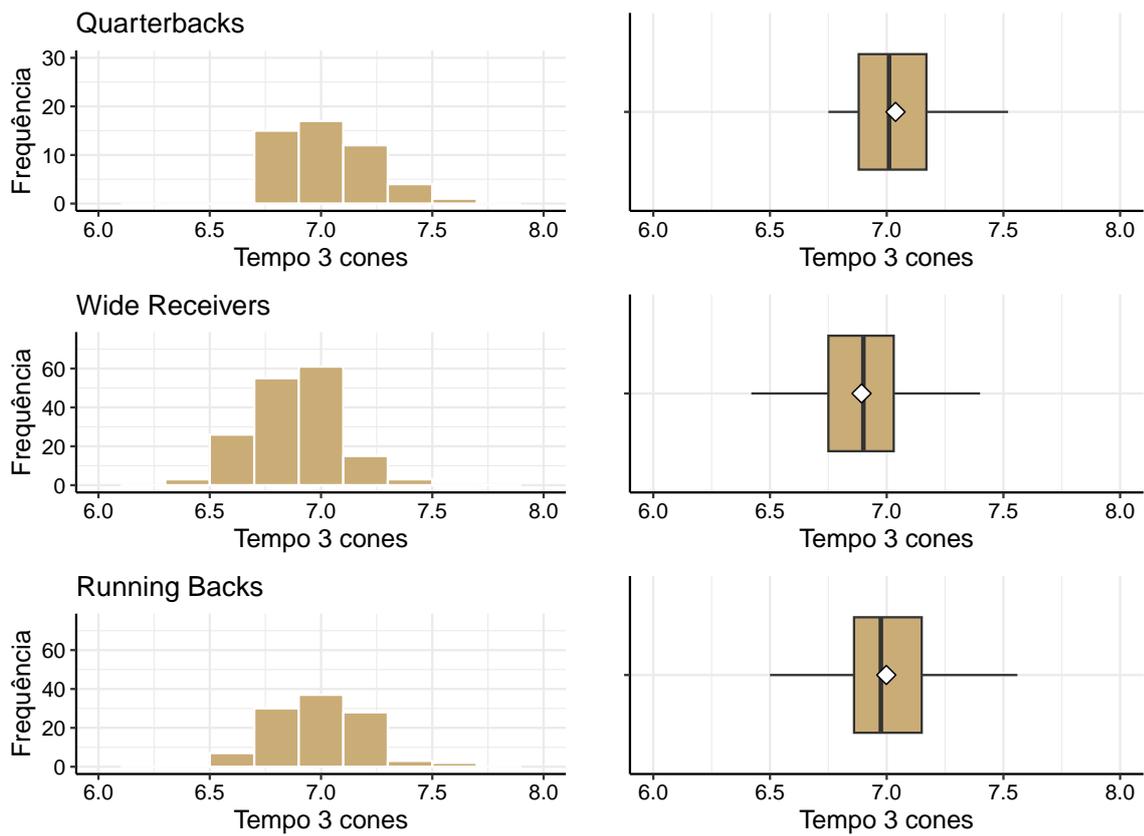


Figura 15: Histogramas e Boxplots para a variável 3-cone drill por posição

Aqui é interessante perceber como recebedores são em média mais ágeis para os 3 cones, o que pode ser reflexo de serem a posição de menor massa do estudo.

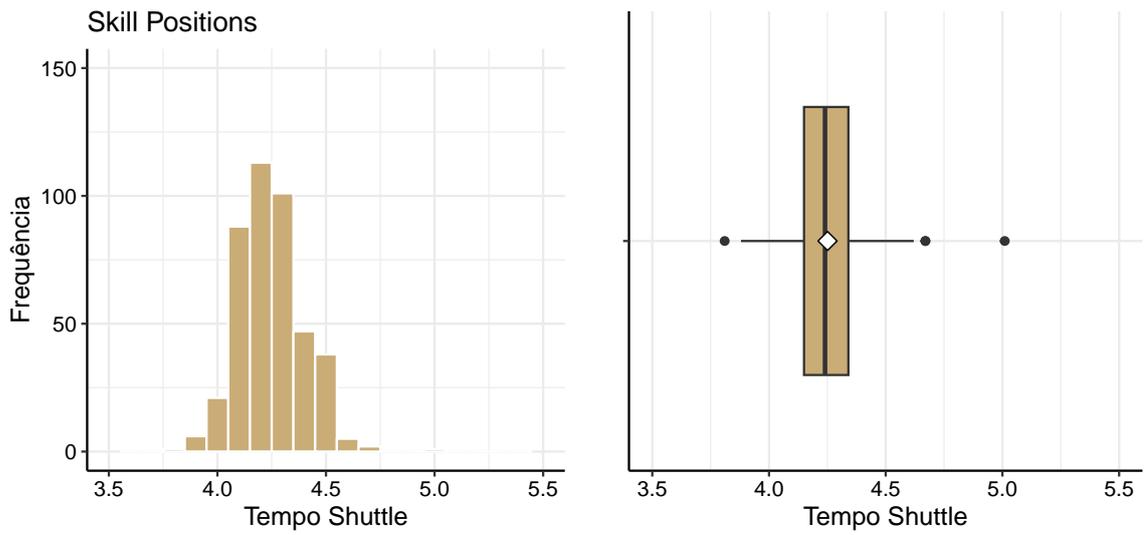


Figura 16: Histogramas e Boxplots para a variável Shuttle

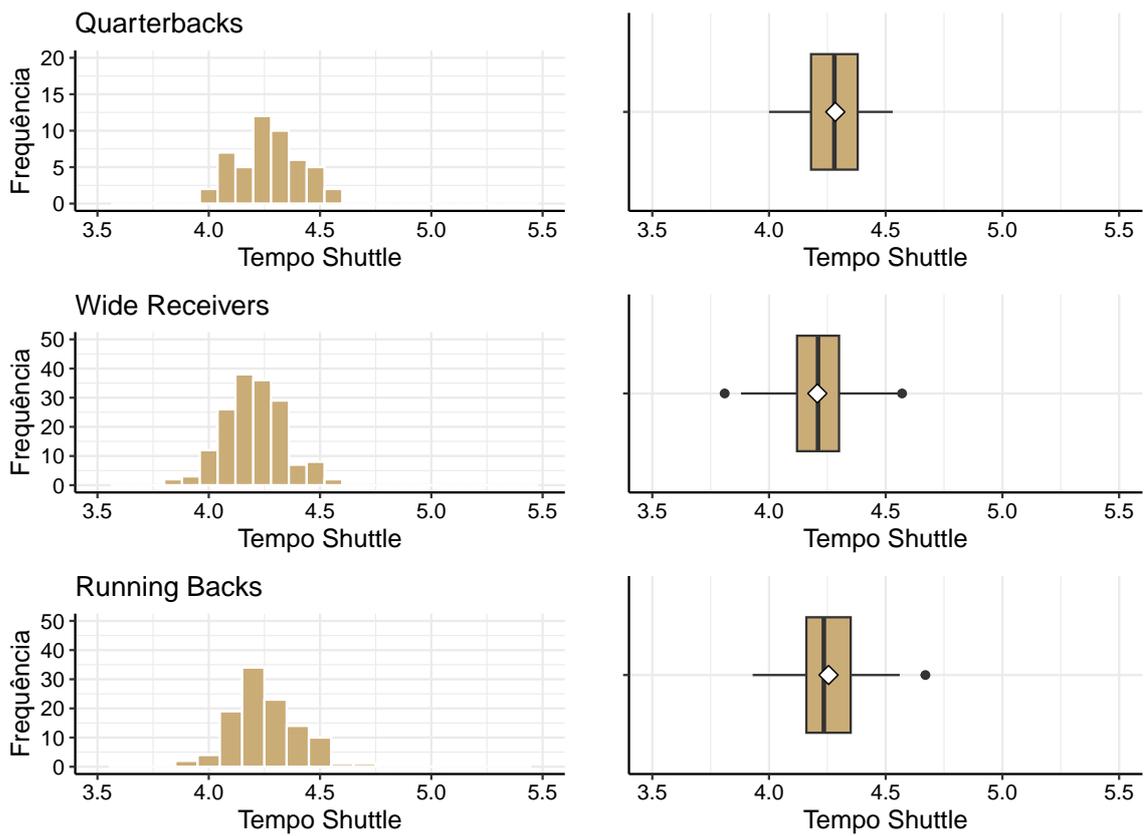


Figura 17: Histogramas e Boxplots para a variável Shuttle por posição

Para a variável do teste de agilidade *Shuttle*, nota-se um comportamento semelhante ao teste dos cones, contudo, talvez com uma proximidade maior entre as posições

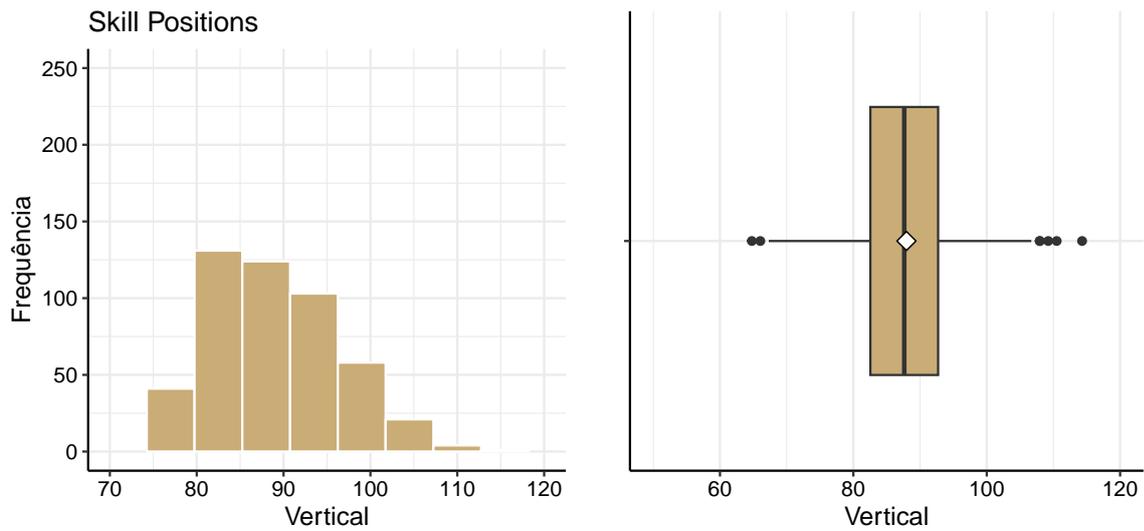


Figura 18: Histogramas e Boxplots para a variável altura do salto vertical

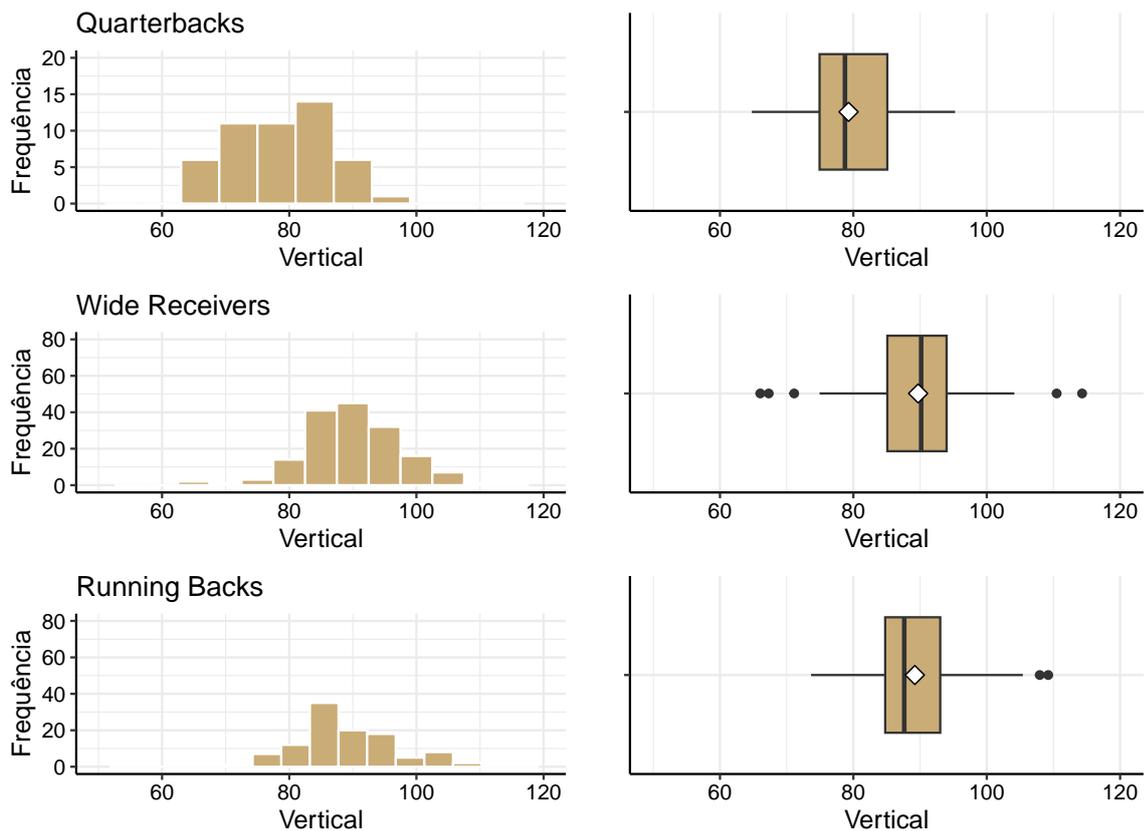


Figura 19: Histogramas e Boxplots para a variável altura do salto vertical por posição

Pode-se perceber pelo histograma que os recebedores e corredores obtiveram, em média, resultados maiores que os passadores, o que é de se esperar pela diferença de atletismo que é necessária para cada uma das posições.

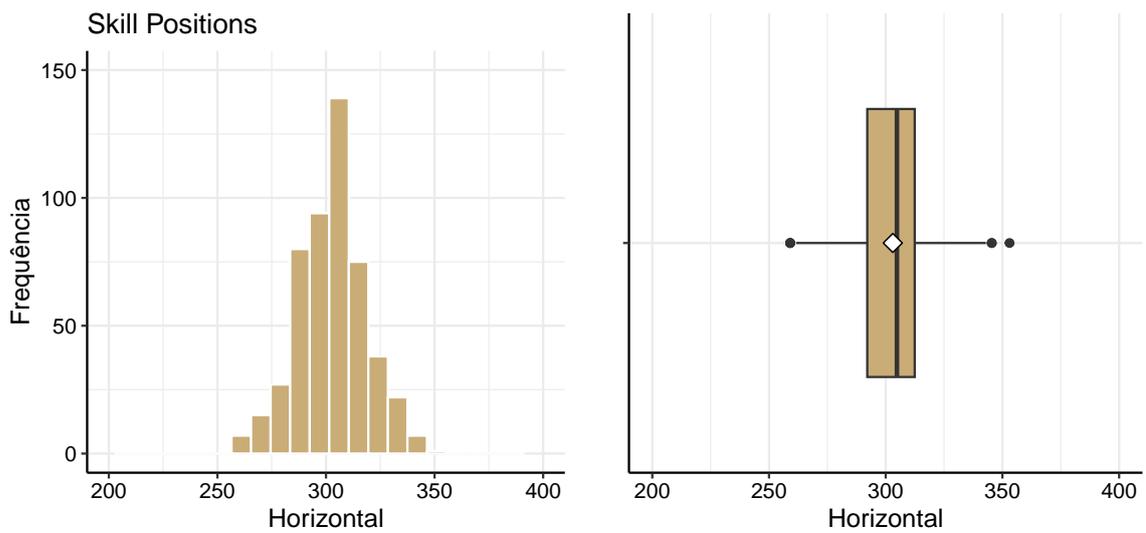


Figura 20: Histogramas e Boxplots para a variável distância do salto horizontal

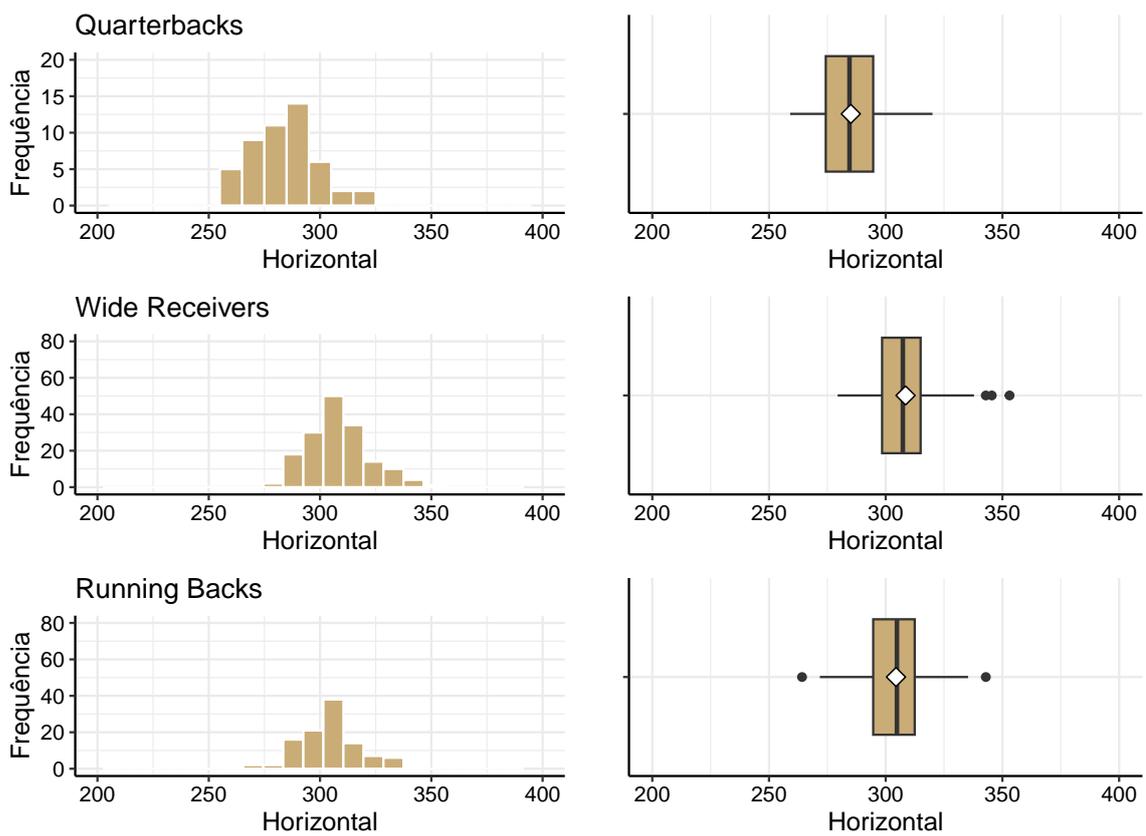


Figura 21: Histogramas e Boxplots para a variável distância do salto horizontal por posição

Para o salto horizontal há um comportamento semelhante ao vertical para a amostra de jogadores de *Skill positions*.

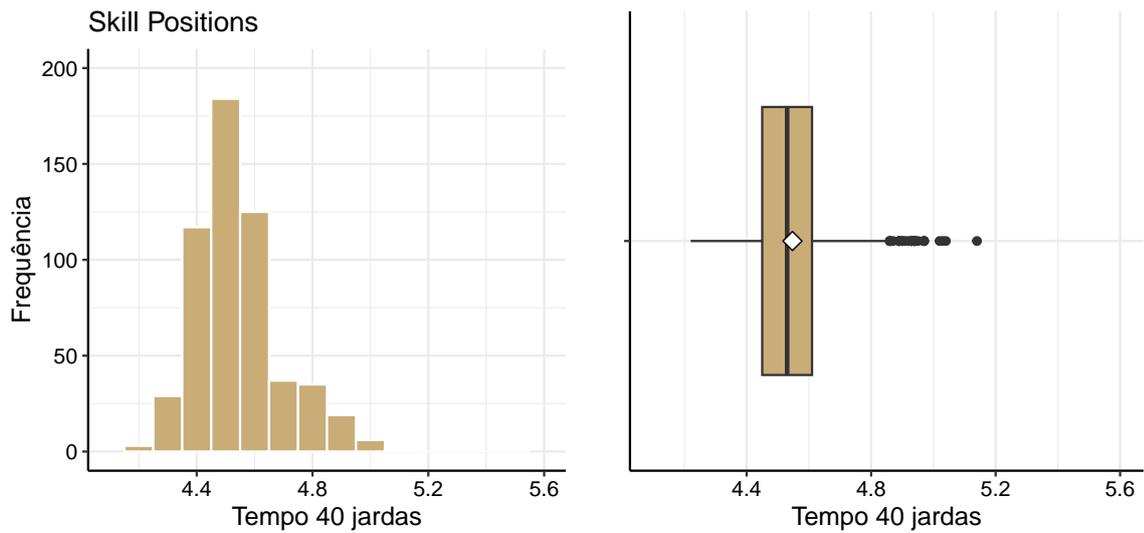


Figura 22: Histogramas e Boxplots para a variável tempo no tiro de 40 jardas

Pelo histograma, é notável que a maior concentração de atletas das *Skill positions* está perto dos 4,4 a 4,6 segundos

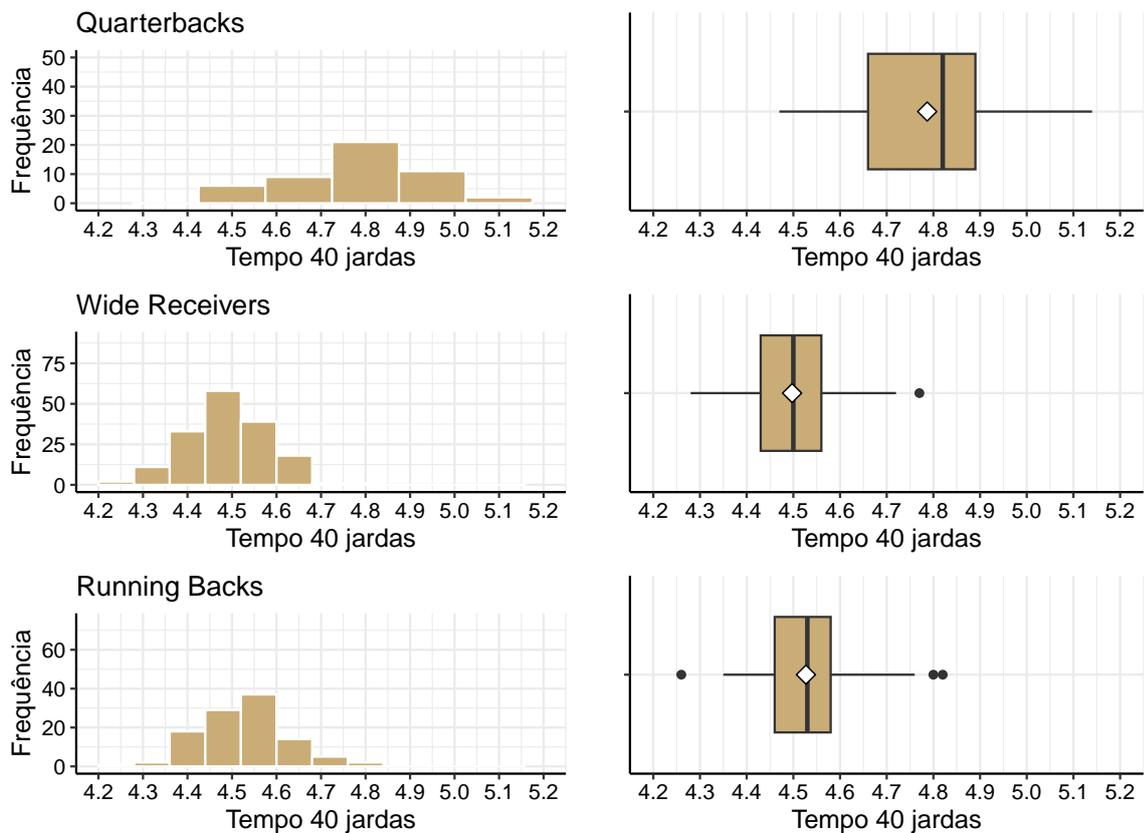


Figura 23: Histogramas e Boxplots para a variável tempo no tiro de 40 jardas por posição

No provável teste mais importante e renomado do *Combine*, percebe-se pela Tabela de medidas de posição uma média de conclusão das 40 jardas consideravelmente

inferior nos recebedores e corredores em relação aos passadores e também indicando a grande proximidade dos números, com uma amplitude de menos de 1s para estes jogadores no teste.

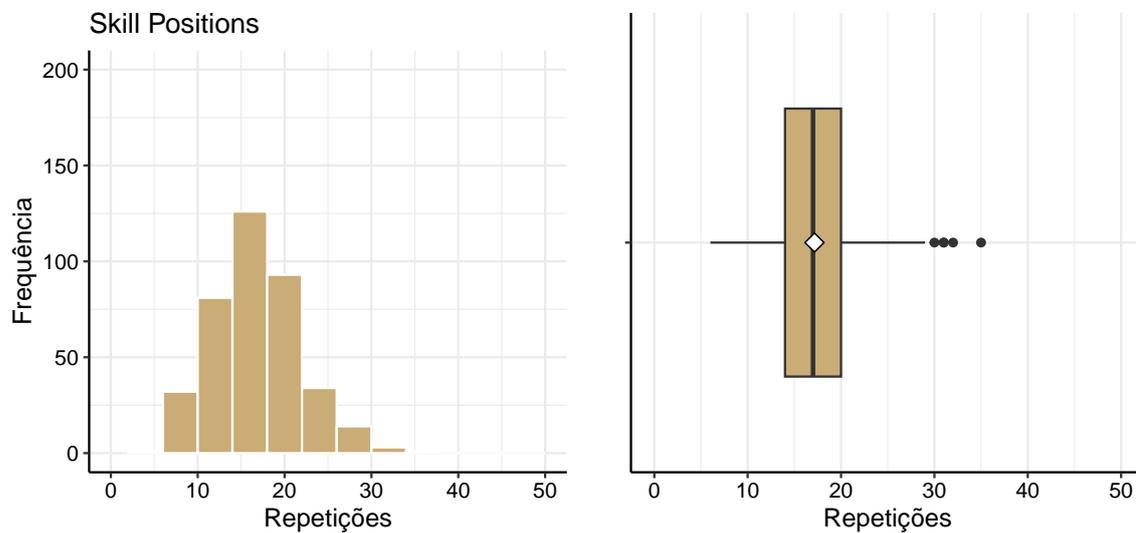


Figura 24: Histogramas e Boxplots para a variável número de repetições no supino, com 102kg

Pelo boxplot, percebe-se que o terceiro quartil dos jogadores de habilidade são 20 repetições, este teste não sendo o foco principal dos atletas dessas posições.

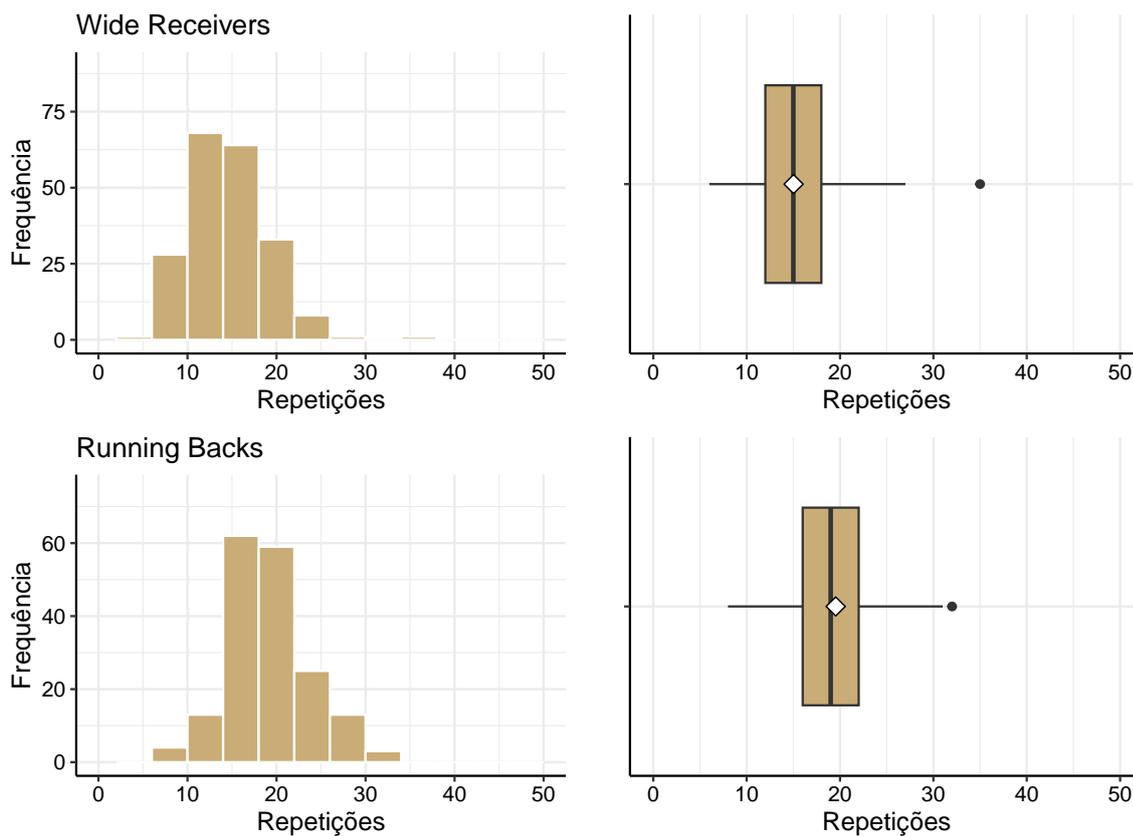


Figura 25: Histogramas e Boxplots para a variável número de repetições no supino, com 102kg por posição

Por último, para o único teste de força dos membros superiores, o supino com carga de 102kg para todos os atletas, percebe-se uma maior força nos corredores, em relação aos recebedores. Como já informado, os *Quarterbacks* da amostra não participaram do teste.

4.1.4 Análise bivariada

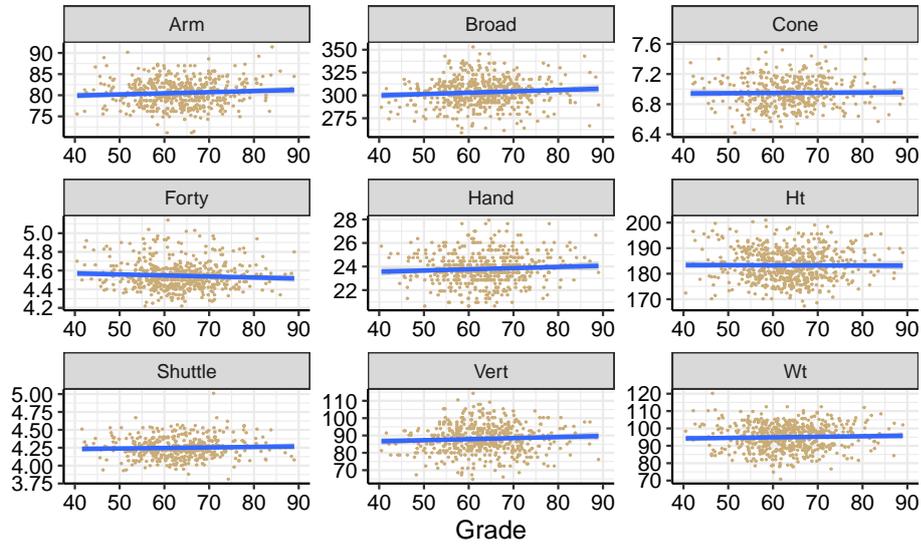


Figura 26: Dispersão da varável resposta pelas variáveis explicativas para as *Skill positions*

Assim, observando os gráficos de dispersão para as *skill positions*, pode-se notar uma leve correlação entre a variável resposta e as variáveis *Wt*, *Arm* e *Shuttle*, e uma aparente correlação mais severa com as variáveis *Forty* e *Broad*. Tal tendência será checada a partir do teste de correlação de Spearman.

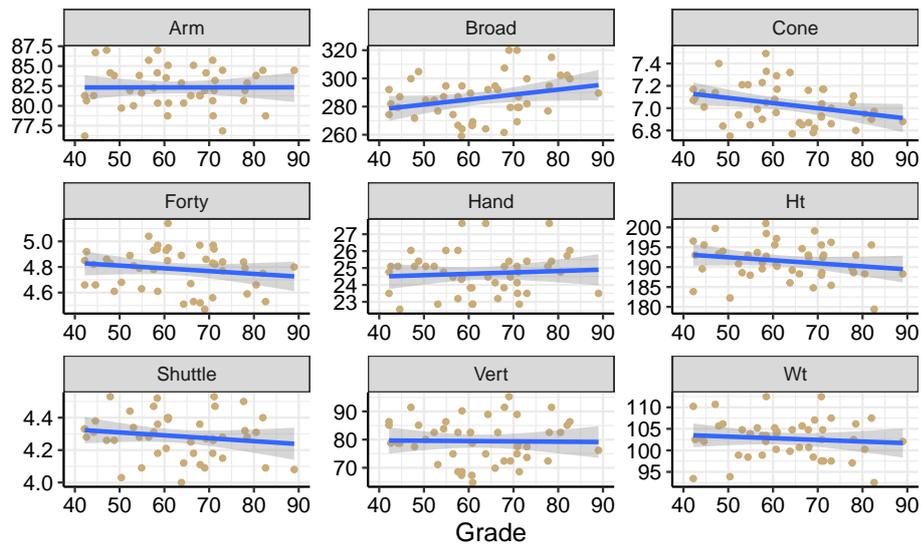


Figura 27: Dispersão da varável resposta pelas variáveis explicativas para os *Quarterbacks*

Para os passadores, percebe-se que há correlação aparentemente mais acentuada para *Cone*, *Shuttle* e *Broad*, e ainda para *Ht*, *Hand* e tiro de 40 jardas.

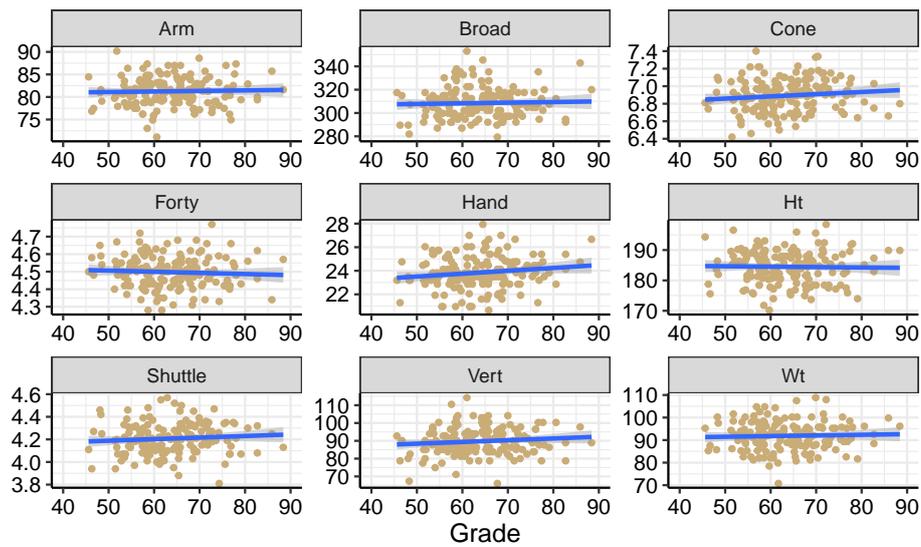


Figura 28: Dispersão da varável resposta pelas variáveis explicativas para os *Wide receivers*

Para os recebedores, as variáveis com aparente correlação são *Hand*, *Vert*, *Cone* e *Forty*, mas para essa posição apenas o tamanho da mão aparenta ser mais acentuado pela análise gráfica.

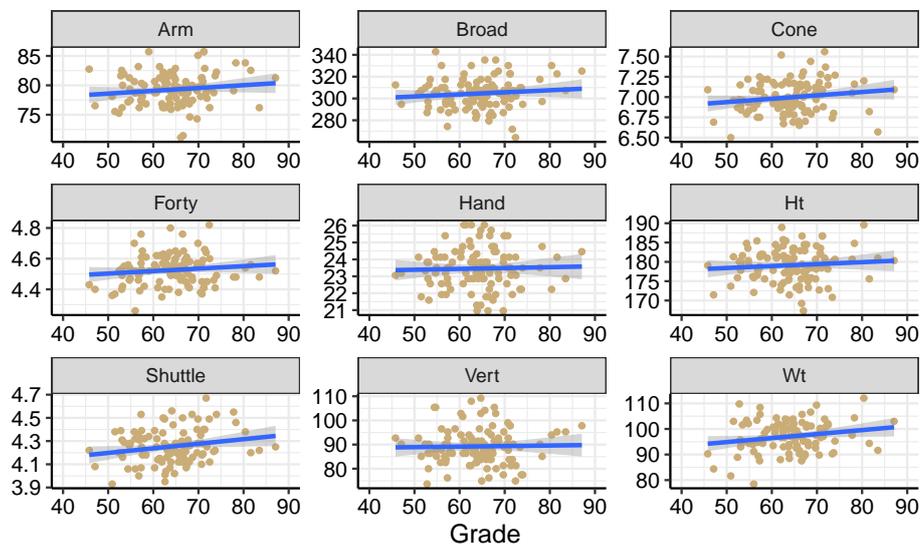


Figura 29: Dispersão da varável resposta pelas variáveis explicativas para os *Running backs*

Por último, na posição de *Running back*, há algumas variáveis que aparentam apresentar correlação com o desempenho do jogador, as variáveis *Shuttle*, *Cone* e *Wt* aparentam ter tendência mais acentuada, seguido do *Broad*

4.1.5 Testes de correlação entre as variáveis

A fim de garantir ou refutar as tendências analisadas graficamente, para cada duas variáveis e para todos os pares do banco, testes de correlação de Spearman, feitos sob hipóteses já mencionadas, que geraram os resultados que se seguiram:

- Variável dependente e independentes

Tabela 6: Tabela com p-valores para o teste de correlação de Spearman entre a variável resposta e as independentes

	Arm	Broad	Cone	Forty	Hand	Ht	Shuttle	Vert	Wt
Grade Skill positions	0.277	0.026	0.667	0.048	0.33	0.447	0.553	0.06	0.431
QB	0.651	0.068	0.231	0.497	0.906	0.335	0.547	0.666	0.577
WR	0.29	0.682	0.089	0.695	0.045	0.925	0.172	0.095	0.287
HB	0.172	0.052	0.673	0.63	0.178	0.295	0.515	0.516	0.105

Assim, pode-se perceber que, para as posições de habilidade, há correlação estatística entre o desempenho e as variáveis *Forty* e *Broad* apenas.

Dentre as posições individualmente, apenas *Hand* rejeitou para os recebedores.

- Variáveis independentes

Prosseguindo-se com a análise dos coeficientes de correlação, a partir dos mapas de calor a seguir, nos quais as cores indicam a força da correlação, e os valores aparentes no gráfico são o p-valor para o teste de correlação de Spearman, sob as hipóteses já citadas. Dessa forma, os resultados para ambas as amostras são as que se seguem:

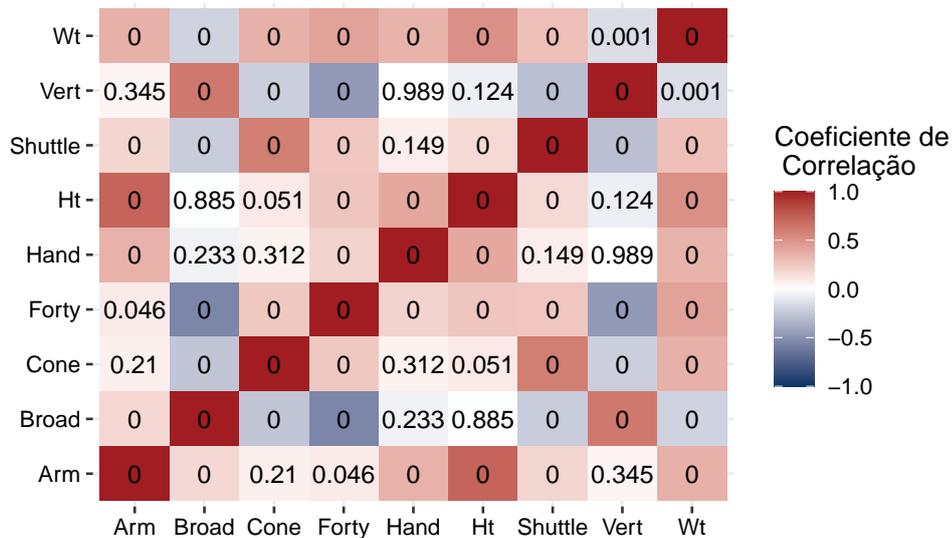


Figura 30: Mapa de calor com p-valores do teste de correlação de Spearman entre as variáveis independentes para Skill positions

Agora, ao apresentar as correlações entre os jogadores das posições da habilidade, verifica-se que para essa amostra, não se rejeita a hipótese de que não há correlação entre as variáveis, para os pares Altura e salto horizontal, altura e teste de três cones (por muito pouco), altura e salto vertical, cone com o tamanho dos braços e salto horizontal com mão e altura, fora estas variáveis, todas as outras análises 2 a 2 rejeitaram H0, ou seja, apresentaram correlação.

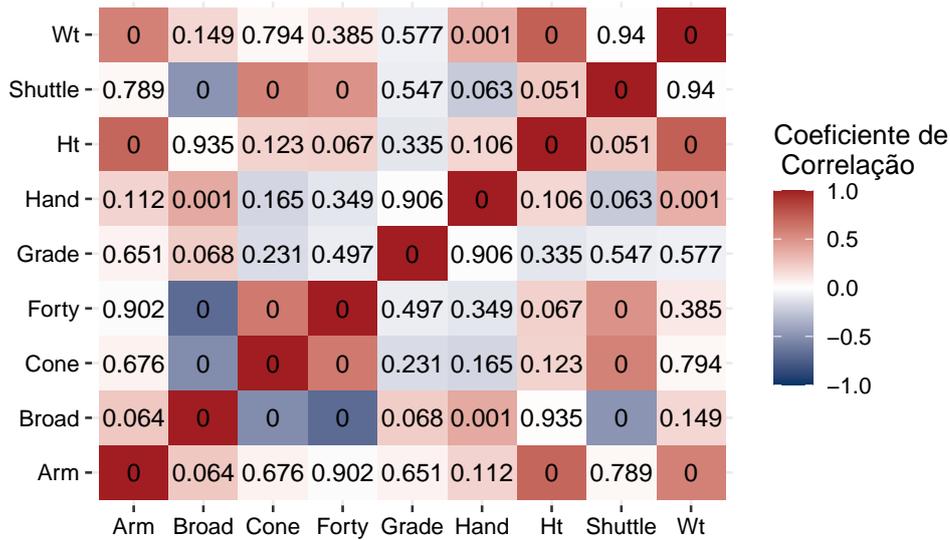


Figura 31: Mapa de calor com correlação entre as variáveis independentes para os *Quarterbacks*

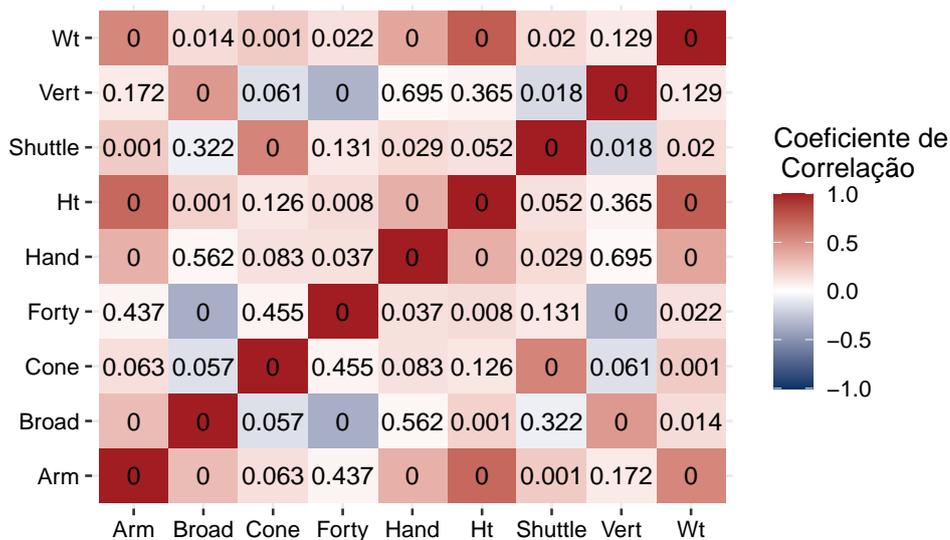


Figura 32: Mapa de calor com correlação entre as variáveis independentes para os *Wide receivers*

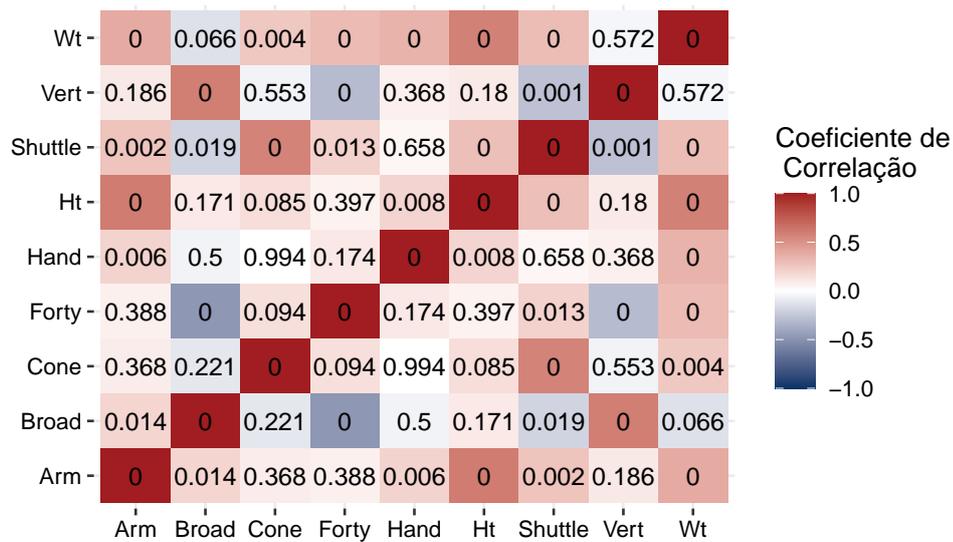


Figura 33: Mapa de calor com correlação entre as variáveis independentes para os *Running backs*

Pode-se perceber que para os passadores, a maioria dos pares de variáveis foi não significativo para correlação, de outra maneira, os recebedores, muitas das variáveis ainda apresentaram correlação, analogamente aos recebedores.

4.2 Modelo de regressão linear múltipla

Nessa seção, será feita a modelagem dos dados com um modelo de regressão linear múltipla dos jogadores de posições de habilidade, para assim, obter os resultados objetivados no trabalho. Além disso, a fim de entender como o modelo geral adequa-se aos dados, este será comparado com modelos das posições individualmente. Os resultados são os que se seguem:

4.2.1 Modelo das *Skill positions*

Modelo inicial

Para iniciar a modelagem dos dados, a fim de analisar a variável dependente à luz de todas as outras variáveis explicativas, o modelo tem o formato que se segue, sendo Y a variável resposta nota de desempenho:

$$\hat{Y} = \beta_0 + \beta_1 Ht + \beta_2 Wt + \beta_3 Hand + \beta_4 Arm + \beta_5 Forty + \beta_6 Vert + \beta_7 Broad + \beta_8 Cone + \beta_9 Shuttle + \epsilon.$$

A partir deste modelo, foram obtidos os resultados que se seguem:

Tabela 7: Análise do modelo inicial de regressão linear múltipla

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	-0.9679	45.7145	-0.02	0.9831
Ht	-0.2803	0.1520	-1.84	0.0663
Wt	0.0448	0.1139	0.39	0.6945
Hand	0.8601	0.4665	1.84	0.0664
Arm	0.2061	0.3008	0.69	0.4939
Forty	0.5792	6.0662	0.10	0.9240
Vert	-0.0357	0.0973	-0.37	0.7138
Broad	0.1128	0.0609	1.85	0.0652
Cone	2.8277	4.1165	0.69	0.4928
Shuttle	4.9023	5.3055	0.92	0.3564

Como pode-se observar, a partir de uma regressão completa, nenhuma variável explicativa nesse modelo inicial é significativa, a um $\alpha = 5\%$.

Seleção de variáveis

A partir do método *stepwise*, selecionando-se os modelos sob o critério de informação AIC, seguem os resultados abaixo.

Tabela 8: Seleção stepwise com critério de informação AIC

Parâmetros do modelo										AIC
Intercepto	Ht	Hand	Broad	Shuttle	Cone	Arm	Wt	Vert	Forty	
X	X	X	X	X	X	X	X	X	X	1148.59
X	X	X	X	X	X	X	X	X	-	1146.60
X	X	X	X	X	X	X	X	-	-	1144.75
X	X	X	X	X	X	X	-	-	-	1142.95
X	X	X	X	X	X	-	-	-	-	1141.42
X	X	X	X	X	-	-	-	-	-	1140.10

Dessa maneira, através do Stepwise, pode-se perceber que o menor AIC obtido foi através do modelo com as variáveis *Ht*, *Hand*, *Broad* e *Shuttle*, e assim, o novo modelo recebe:

$$\hat{Y} = \beta_0 + \beta_1 Ht + \beta_3 Hand + \beta_7 Broad + \beta_9 Shuttle + \epsilon.$$

E obteve os seguintes resultados:

Tabela 9: Análise do modelo final de regressão linear múltipla das *Skill positions*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	8.7970	27.2933	0.32	0.7475
Ht	-0.1881	0.1033	-1.82	0.0698
Hand	0.9538	0.4437	2.15	0.0326
Broad	0.1004	0.0373	2.69	0.0077
Shuttle	8.4188	4.1737	2.02	0.0447

Agora, tendo em vista o modelo final escolhido, pode-se perceber que, a $\alpha = 5\%$, as variáveis *Hand*, *Broad* e *Shuttle* rejeitam H_0 , de tal forma que há evidências de que os valores dos coeficientes $\beta \neq 0$ para essas variáveis, ou seja, elas influenciam na variável dependente *Grade*. Vale observar também, que para o nível de significância do estudo, a variável *Ht* não rejeita a hipótese nula, contudo, apresenta um valor muito baixo também, e, através do AIC, sabe-se que o melhor modelo ainda é o que a contém.

Análise ANOVA

Tabela 10: Tabela ANOVA do modelo final para a amostra geral

	GL	Soma de quadrados	Quadrado médio	Estatística F	P-valor
Ht	1	76.75	76.75	0.91	0.3409
Hand	1	460.68	460.68	5.47	0.0202
Broad	1	457.24	457.24	5.43	0.0206
Shuttle	1	342.93	342.93	4.07	0.0447
Residuals	251	21154.95	84.28		

Através da tabela, pode-se concluir que, adicionando-se as variáveis a um modelo contendo apenas o intercepto, a variável *Hand* é a que mais adiciona à soma de quadrados da regressão, e a *Broad*, levemente menos, contudo, adicionar apenas a variável *Ht* ao intercepto acrescenta apenas 76,75 na soma dos quadrados da regressão.

Além disso, pode-se perceber que a soma de quadrado dos erros é muito maior que a da regressão, o que mostra que o modelo apresenta baixo percentual de explicação na variável resposta, existem outras fontes de variações não incluídas nas variáveis explicativas selecionadas.

Análise dos coeficientes do modelo

Dessa forma, a partir do modelo de regressão linear das posições de habilidade, pode-se obter a função do modelo, que tem a seguinte forma:

$$\hat{Y}_i = 8,797 - 0,1881 \cdot Ht_i + 0,9538 \cdot Hand_i + 0,1004 \cdot Broad_i + 8,4188 \cdot Shuttle_i + \epsilon.$$

Assim, interpreta-se que, para cada desvio padrão da variável altura, o modelo indica uma queda de 1,0759 na variável resposta, mantendo-se o resto fixo. Para a variável mão, a variável resposta obtém aumento de 1,2017; e de 1,42 para cada desvio do salto horizontal. Por último, há um aumento de 1,17 unidades no desempenho para a variável *Shuttle*.

Ou seja, para o modelo obtido, os jogadores obtiveram melhor desempenho ao obter pior desempenho no teste *Shuttle* (demoraram mais para concluir), menor altura, e maior distância de salto vertical e tamanho de mão.

Tabela 11: Valores do modelo final geral de regressão linear múltipla

R^2	R^2 Ajustado	Estatística F	P-valor
0.05947	0.04448	3.968	0.003854

Por fim, vale ressaltar a baixíssima explicação de variação da variável resposta pelo modelo ajustado, antes vista na tabela da ANOVA, e confirmada pelo coeficiente de determinação do modelo, apesar disso, pela estatística F e pelo p-valor, aliados à análise da ANOVA anteriormente, observando que o modelo foi significativo, sob $\alpha = 5\%$.

Análise dos pressupostos

Após o ajuste do modelo de regressão linear múltipla, é necessário conduzir uma análise de resíduos, e avaliar os principais pressupostos associados a esse tipo de análise: normalidade, homocedasticidade e independência, então, abaixo prossegue-se tal análise.

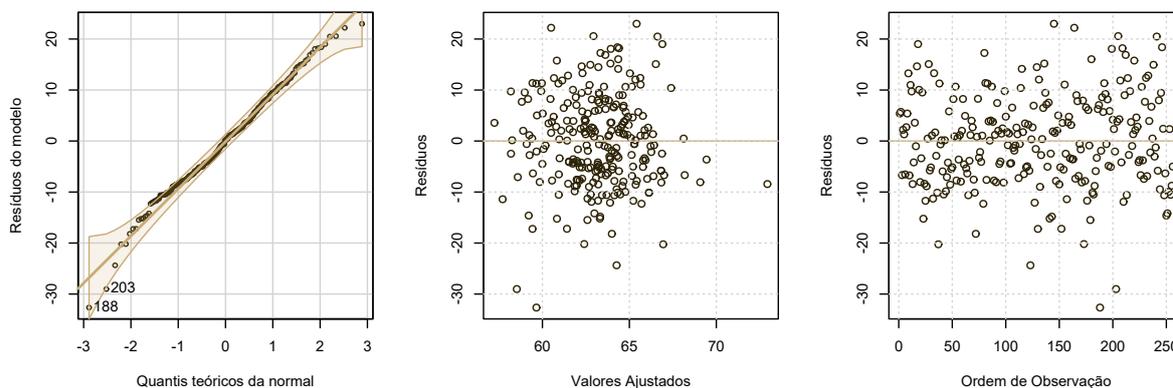


Figura 34: Pressupostos modelo geral

Dessa maneira, apresentando que os resíduos do modelo estão aparentemente seguindo distribuição normal, que é uma das premissas na análise dos pressupostos do modelo, a normalidade dos mesmos apenas será garantida com o teste de Shapiro-Wilk.

Visualizando os resíduos versus valores ajustados, pode-se analisar o pressuposto da homocedasticidade, dessa maneira, pelo gráfico, os resíduos aparentam ser homocedásticos, pois se comportam de maneira aleatória. A homocedasticidade deles será comprovada a partir do teste de Breusch-Pagan.

Após isso, checka-se a independência, e pode ser feita a partir do gráfico de resíduos versus ordem de observação, a partir da análise então, aparenta que eles são independentes, que será confirmado com um teste de Durbin-Watson.

Por último, como indicados nos gráficos, nenhum pressuposto aparenta ser rejeitado, contudo, para haver comprovação estatística, é necessário fazer os testes de hipótese, os pressupostos analisados, seus respectivos testes e p-valores se seguem na tabela a seguir:

Tabela 12: P-valores dos testes de pressupostos para modelo de regressão linear múltiplo final

Pressuposto	Teste	P-Valor
Normalidade	Shapiro-Wilk	0.1574
Homocedasticidade	Breusch-Pagan	0.1896
Independência	Durbin-Watson	0.7149

Análise do VIF

Como obtido a partir a análise da correlação de Spearman, foi observada correlação em muitas das variáveis explicativas, por isso, é importante fazer a verificação do Fator de Inflação de variância, para saber se no modelo ajustado houve multicolinearidade.

Para isso, seguem as variáveis e os respectivos valores de seus VIF's:

Tabela 13: Valores de VIF de cada variável explicativa

Variável	VIF_j
Ht	1.22
Hand	1.18
Broad	1.06
Shuttle	1.08

Assim, todas as variáveis independentes do modelo de regressão apresentaram um valor de VIF baixo, sendo o ideal abaixo de 5, logo, não há multicolinearidade no modelo ajustado.

Análise de pontos influentes

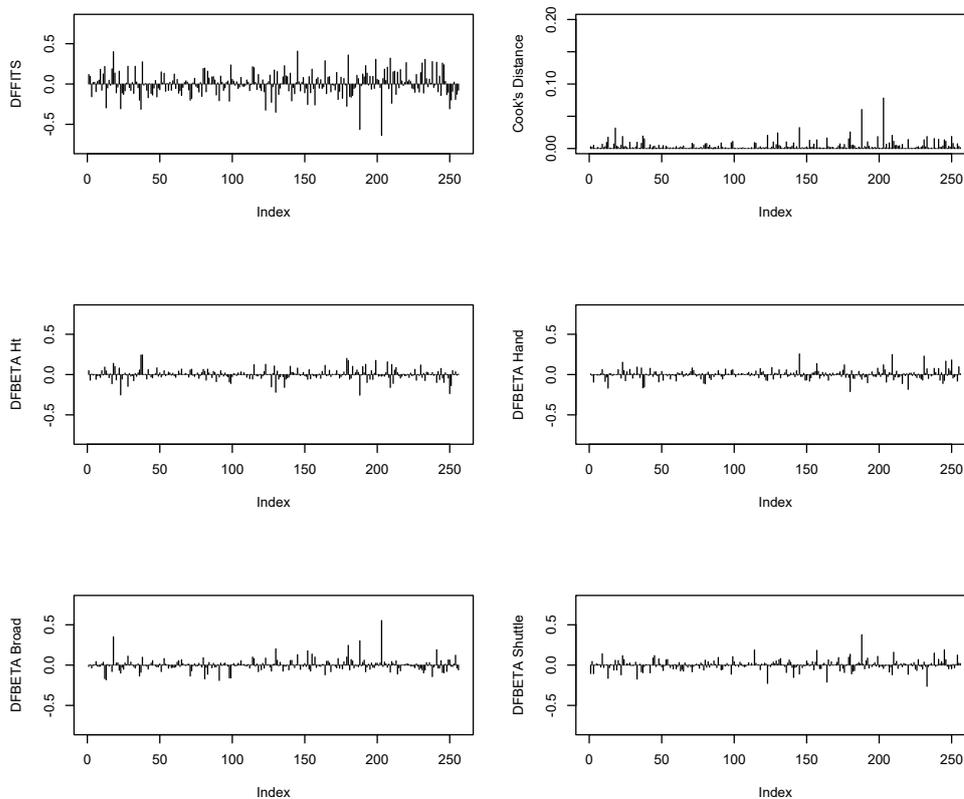


Figura 35: Análise gráfica de pontos influentes do modelo geral

Dessa maneira, a partir da análise das métricas, não há pontos de alavancagem na regressão, com regra básica de 0,8 para a distância de Cook e 1 para DFFITS e DFBETAS.

4.2.2 Modelos por posição

A fim de entender o comportamento dos mesmos dados de forma individual, a seção a seguir contém três modelos de regressão linear múltipla, cada um modelando uma posição entre as posição de habilidade, ou seja, modelando amostras nas quais contém apenas a posição desejada.

Modelo inicial

Como modelo geral, para iniciar, é feita a análise da variável dependente em relação a todas as outras variáveis explicativas. Os modelo têm, de forma análoga ao modelo geral, o formato que se segue, sendo Y a variável resposta nota de desempenho:

$$\hat{Y} = \beta_0 + \beta_1 Ht + \beta_2 Wt + \beta_3 Hand + \beta_4 Arm + \beta_5 Forty + \beta_6 Vert + \beta_7 Broad + \beta_8 Cone + \beta_9 Shuttle + \epsilon.$$

A partir deste modelo, foram obtidos os resultados que se seguem:

Quarterbacks

Tabela 14: Análise do modelo inicial de regressão linear múltipla para os *Quarterbacks*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	-16.7304	232.3007	-0.07	0.9431
Ht	-1.4445	0.9806	-1.47	0.1512
Wt	-0.0071	0.8003	-0.01	0.9930
Hand	2.0964	2.1733	0.96	0.3424
Arm	2.1465	1.3400	1.60	0.1197
Forty	16.8179	27.1839	0.62	0.5408
Vert	-1.0756	0.5068	-2.12	0.0422
Broad	0.6294	0.3238	1.94	0.0614
Cone	-15.1634	19.1263	-0.79	0.4341
Shuttle	13.4819	22.1253	0.61	0.5469

A tabela indica que, a um $\alpha = 5\%$, o salto vertical foi a única variável que foi significativa.

Running backsTabela 15: Análise do modelo inicial de regressão linear múltipla para os *Running backs*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	-40.4734	77.3599	-0.52	0.6023
Ht	-0.2932	0.3138	-0.93	0.3530
Wt	0.1287	0.2143	0.60	0.5498
Hand	-0.0906	0.8168	-0.11	0.9120
Arm	0.2433	0.4176	0.58	0.5618
Forty	8.7107	11.3580	0.77	0.4455
Vert	-0.0135	0.1617	-0.08	0.9336
Broad	0.1652	0.0866	1.91	0.0602
Cone	0.0649	5.6149	0.01	0.9908
Shuttle	9.1535	8.8686	1.03	0.3052

Da maneira que no modelo inicial da amostra geral, a partir de uma regressão completa, não houveram variáveis explicativas significativas nesse modelo, a um $\alpha = 5\%$.

Wide receiversTabela 16: Análise do modelo inicial de regressão linear múltipla para os *Wide receivers*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	-23.5425	60.1074	-0.39	0.6960
Ht	-0.0970	0.2215	-0.44	0.6623
Wt	0.0104	0.1691	0.06	0.9513
Hand	0.8785	0.5957	1.47	0.1429
Arm	0.1572	0.3475	0.45	0.6518
Forty	6.3564	8.7836	0.72	0.4707
Vert	0.0322	0.1085	0.30	0.7671
Broad	0.0387	0.0718	0.54	0.5906
Cone	5.3179	4.7596	1.12	0.2661
Shuttle	-2.3346	6.7271	-0.35	0.7292

Por último, analisando os recebedores, o modelo completo apresentou também 0 variáveis significantes.

Agora, é necessária a seleção de variáveis, através de um *stepwise*, sob o critério

de informação AIC, seus resultados foram os seguintes.

Quarterbacks

Tabela 17: Seleção stepwise com critério de informação AIC para os *Quarterbacks*

Parâmetros do modelo										AIC
Intercepto	Ht	Arm	Vert	Broad	Hand	Forty	Cone	Shuttle	Wt	
X	X	X	X	X	X	X	X	X	X	217.39
X	X	X	X	X	X	X	X	X	-	215.39
X	X	X	X	X	X	X	X	-	-	213.91
X	X	X	X	X	X	X	-	-	-	212.38
X	X	X	X	X	X	-	-	-	-	210.94
X	X	X	X	X	-	-	-	-	-	209.78

Dessa maneira, após a seleção de variáveis para os *Quarterbacks* através do *stepwise*, o menor AIC foi obtido através do modelo com as variáveis *Ht*, *Arm*, *Vert* e *Broad*, e assim, o novo modelo recebe:

$$\hat{Y}_i = \beta_0 + \beta_1 Ht + \beta_4 Arm + \beta_6 Vert + \beta_7 Broad + \epsilon.$$

E obteve os seguintes resultados:

Tabela 18: Análise do modelo final de regressão linear múltipla para os *Quarterbacks*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	40.3521	99.3916	0.41	0.6872
Ht	-1.3634	0.6977	-1.95	0.0587
Arm	2.2417	1.1963	1.87	0.0693
Vert	-1.0327	0.4232	-2.44	0.0199
Broad	0.6243	0.2185	2.86	0.0072

Sendo feita a análise da tabela anova, obtiveram-se os seguintes resultados:

Tabela 19: Tabela ANOVA do modelo final para os *Quarterbacks*

	GL	Soma de quadrados	Quadrado médio	Estatística F	P-valor
Ht	1	170.90	170.90	1.01	0.3210
Arm	1	1346.64	1346.64	7.98	0.0077
Vert	1	41.03	41.03	0.24	0.6249
Broad	1	1376.82	1376.82	8.16	0.0072
Residuals	35	5903.60	168.67		

Através da tabela, pode-se perceber que o modelo para a posição de QB é o que obteve maior explicação para a variação da variável dependente *Grade*.

Assim, os coeficientes do modelo para os *Quarterbacks* têm os seguintes valores:

$$\hat{Y}_i = 40,3521 - 1,3634 \cdot Ht_i + 2,2417 \cdot Arm_i - 1,0327 \cdot Vert_i + 0,6243 \cdot Broad_i + \epsilon.$$

A partir do modelo final escolhido, tem-se como resultado, a $\alpha = 5\%$, que as variáveis *Vert*, *Broad* rejeitam H_0 , apresentando evidências de que os valores dos coeficientes $\beta \neq 0$ para essas variáveis. Analogamente ao modelo geral, que para o nível de significância do estudo, as variáveis *Ht* e *Arm* não rejeitam a hipótese nula, apesar de seus p-valores estarem muito próximos ao α .

Dessa forma, observa-se que o modelo, de forma semelhante ao geral, indica que a variável altura é negativamente proporcional ao desempenho, além também da altura do salto vertical, pois a cada desvio padrão adicional nessas variáveis, os *Quarterbacks* apresentaram diminuição de 1,6769 e de 7,622 unidades na variável resposta respectivamente. Apesar disso, cada desvio acrescido no tamanho dos braços apresentou aumento de 5,936 unidades na nota, e de 9,4426 para o salto horizontal.

Running backsTabela 20: Seleção stepwise com critério de informação AIC para os *Running backs*

Parâmetros do modelo										AIC
Intercepto	Shuttle	Broad	Forty	Ht	Wt	Arm	Hand	Vert	Cone	
X	X	X	X	X	X	X	X	X	X	372.98
X	X	X	X	X	X	X	X	X	-	370.98
X	X	X	X	X	X	X	X	-	-	368.99
X	X	X	X	X	X	X	-	-	-	367.00
X	X	X	X	X	X	-	-	-	-	365.38
X	X	X	X	X	-	-	-	-	-	363.81
X	X	X	X	-	-	-	-	-	-	362.12
X	X	X	-	-	-	-	-	-	-	361.11

Assim, através do Stepwise, pode-se perceber que o menor AIC obtido foi através do modelo com as variáveis *Broad* e *Shuttle*, e assim, o novo modelo recebe:

$$\hat{Y} = \beta_0 + \beta_7 \text{Broad} + \beta_9 \text{Shuttle} + \epsilon.$$

E então, os resultados do modelo foram:

Tabela 21: Análise do modelo final de regressão linear múltipla para os *Running backs*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	-22.0655	33.0967	-0.67	0.5068
Broad	0.1200	0.0590	2.03	0.0454
Shuttle	11.7706	5.8884	2.00	0.0488

Tabela ANOVA dos *Running backs*:Tabela 22: Tabela ANOVA do modelo final para os *Running backs*

	GL	Soma de quadrados	Quadrado médio	Estatística F	P-valor
Broad	1	184.66	184.66	3.01	0.0865
Shuttle	1	245.18	245.18	4.00	0.0488
Residuals	84	5154.39	61.36		

Os coeficientes do modelo para os *Running backs* têm os seguintes valores:

$$\hat{Y}_i = -22,0655 + 0,12 \cdot Broad_i + 11,7706 \cdot Shuttle_i + \epsilon.$$

Portanto, para os *Running backs*, tanto a variável *Shuttle*, quanto a *Broad* rejeitaram H_0 , indicando um $\beta \neq 0$. Para essa posição, cada desvio padrão a mais no teste indicou um aumento em 1,6908 da variável resposta, o que pode representar que o contrário do que se espera, o teste de agilidade não é importante para a posição, visto que quanto mais rápido o resultado deste, menor a nota do jogador. Além disso, para o aumento de 1 desvio padrão no *Shuttle*, é apresentado 1,7655 de aumento na nota dos jogadores.

Por último, a seleção de modelo para os recebedores foi a seguinte:

Tabela 23: Seleção stepwise com critério de informação AIC para os *Wide receivers*

Parâmetros do modelo										AIC
Intercepto	Hand	Cone	Broad	Forty	Ht	Arm	Shuttle	Vert	Wt	
X	X	X	X	X	X	X	X	X	X	559.61
X	X	X	X	X	X	X	X	X	-	557.62
X	X	X	X	X	X	X	X	-	-	555.71
X	X	X	X	X	X	X	-	-	-	553.88
X	X	X	X	X	X	-	-	-	-	552.04
X	X	X	X	X	-	-	-	-	-	550.11
X	X	X	X	-	-	-	-	-	-	548.53
X	X	X	-	-	-	-	-	-	-	546.96
X	X	-	-	-	-	-	-	-	-	546.11

Dessa maneira, o menor AIC corresponde ao modelo que contém apenas a variável explicativa *Hand*, e assim, o modelo se trata de um modelo de regressão linear simples, de fórmula:

$$Y = \beta_0 + \beta_3 Hand + \epsilon.$$

Assim, após o ajuste, o modelo obteve:

Tabela 24: Análise do modelo final de regressão linear múltipla para os *Wide receivers*

	Estimativa	Erro padrão	Estatística T	P-valor
(Intercepto)	38.2720	12.7622	3.00	0.0033
Hand	1.0610	0.5337	1.99	0.0489

Analisando-se abaixo a tabela ANOVA para os recebedores:

Tabela 25: Tabela ANOVA do modelo final para os *Wide receivers*

	GL	Soma de quadrados	Quadrado médio	Estatística F	P-valor
Hand	1	251.58	251.58	3.95	0.0489
Residuals	129	8212.73	63.66		

Ao adicionar a variável *Hand* a uma regressão com o intercepto para os recebedores, adiciona-se 251,58 na soma de quadrados da regressão, o que ainda é muito distante da soma dos resíduos, que apresentou um valor de 8212,73.

Para o último modelo, dos *Wide receivers*, os coeficientes do modelo têm os seguintes valores:

$$\hat{Y}_i = 38,272 + 1,061 \cdot Hand_i + \epsilon.$$

Dessa forma, surpreendentemente, apenas uma variável foi significativa ao modelar os recebedores, e com isso, a cada desvio padrão a mais na medição das mãos, a nota apresentou um aumento de 37,93%.

Por último, vale analisar a estatística F, os p-valores e os R^2 de cada modelo, e foram obtidos os seguintes valores:

Tabela 26: Valores dos modelos finais de regressão linear múltipla

Modelo	R^2	R^2 Ajustado	Estatística F	P-valor
<i>Quarterbacks</i>	0.3321	0.2558	4.351	0.005831
<i>Running backs</i>	0.07697	0.055	3.502	0.03459
<i>Wide receivers</i>	0.02972	0.0222	3.952	0.04894

Assim, observa-se uma discrepância muito grande nos modelos, de forma que o modelo para os *Quarterbacks* apresentou uma explicação da variabilidade da variável desempenho muito maior que os demais, apesar de uma amostra pequena.

Análise dos pressupostos

Após o ajuste dos modelos, testa-se os seguintes pressupostos: normalidade, homocedasticidade e independência, a partir de análise gráfica e testes de hipótese abaixo.

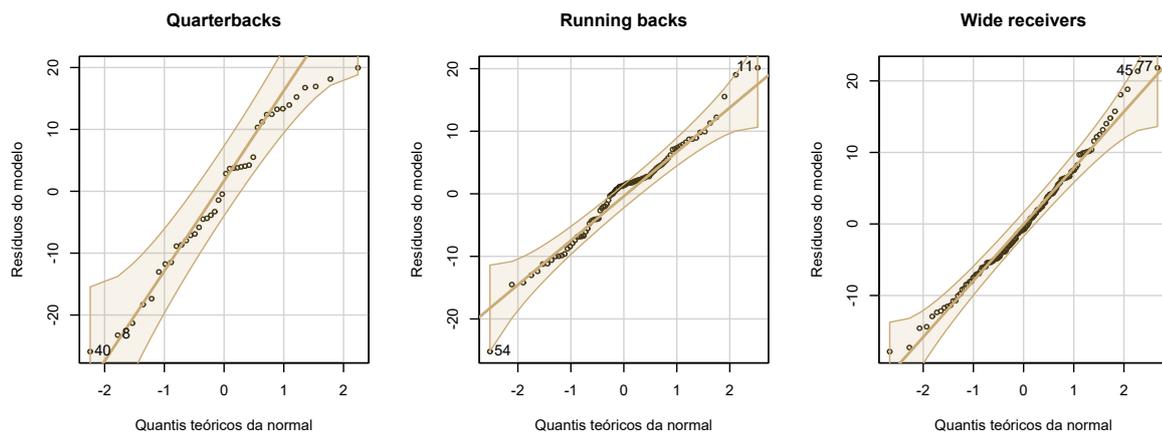


Figura 36: Normalidade

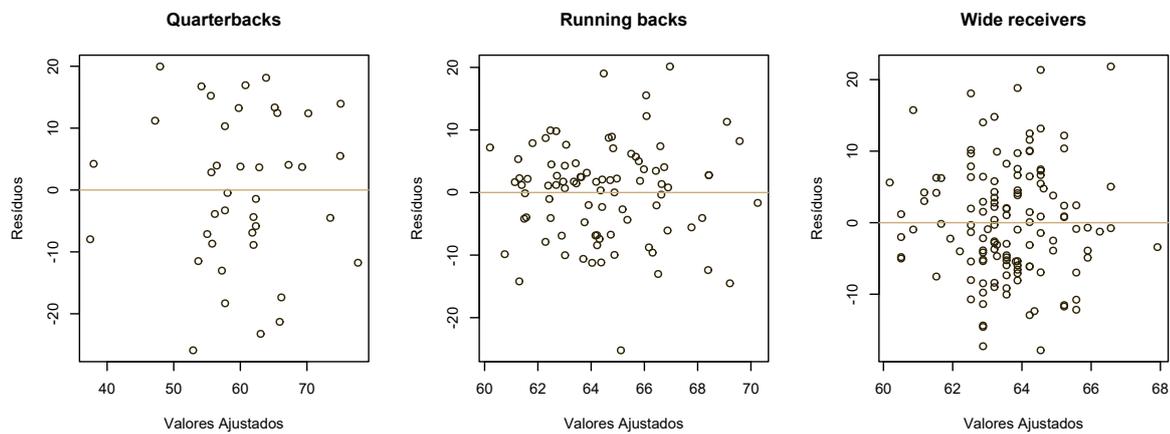


Figura 37: Homocedasticidade

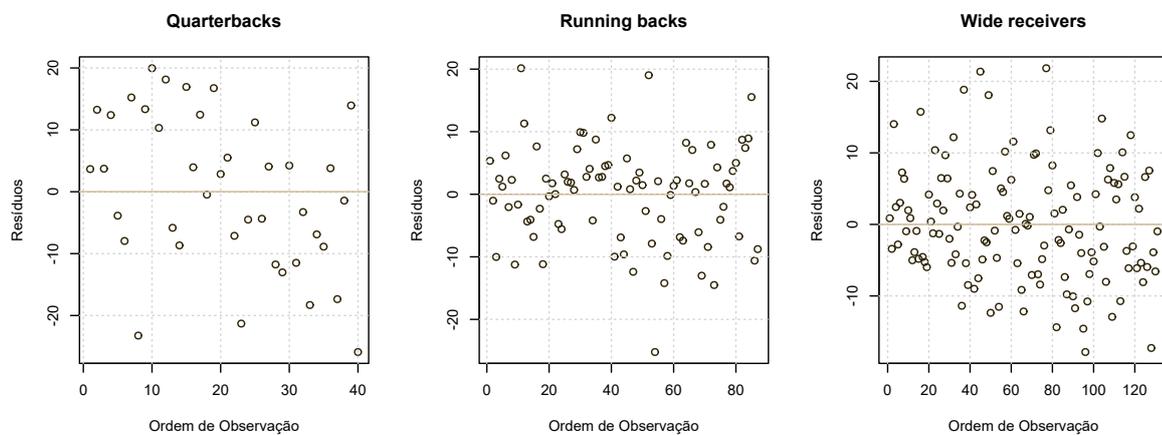


Figura 38: Independência

Pela análise gráfica dos valores dos resíduos vs. quantil da normal, aparente-

mente há normalidade. Visualizando os resíduos vs. valores ajustados, aparentemente há homocedasticidade. E pelos resíduos pelo índice, aparentemente há independência para as três amostras, contudo, todos os pressupostos serão testados com os testes de hipótese a seguir.

Tabela 27: P-valores dos testes de pressupostos para modelo de regressão linear múltiplo final

Modelo	P-Valores		
	Shapiro-Wilk	Breusch-Pagan	Durbin-Watson
Quarterbacks	0.2422	0.3286	0.4259
Running backs	0.2352	0.3476	0.2643
Wide receivers	0.4192	0.2457	0.7144

Análise do VIF

É importante fazer a verificação do Fator de inflação de variância, para saber se no modelo ajustado houve multicolinearidade, com exceção dos recebedores, por haver apenas uma variável explicativa. Para isso, seguem as variáveis e os respectivos valores de seus VIF's:

Tabela 28: Valores de VIF de cada variável explicativa para os *Quarterbacks*

Variável	VIF_j
Ht	2.32
Arm	2.43
Vert	2.27
Broad	2.46

Tabela 29: Valores de VIF de cada variável explicativa para os *Running backs*

Variável	VIF_j
Broad	1.03
Shuttle	1.03

Assim, todas as variáveis independentes do modelo de regressão apresentaram um valor de VIF baixo, sendo o ideal abaixo de 5, logo, não há multicolinearidade nos modelos ajustados para ambas as posições.

Análise de pontos influentes

Pode-se analisar se houveram pontos influentes, também conhecidos como pontos de alavancagem, para tal, usa-se as medidas distância de Cook, DFFITS e DFBETAS, e sua análise gráfica é a que se segue:

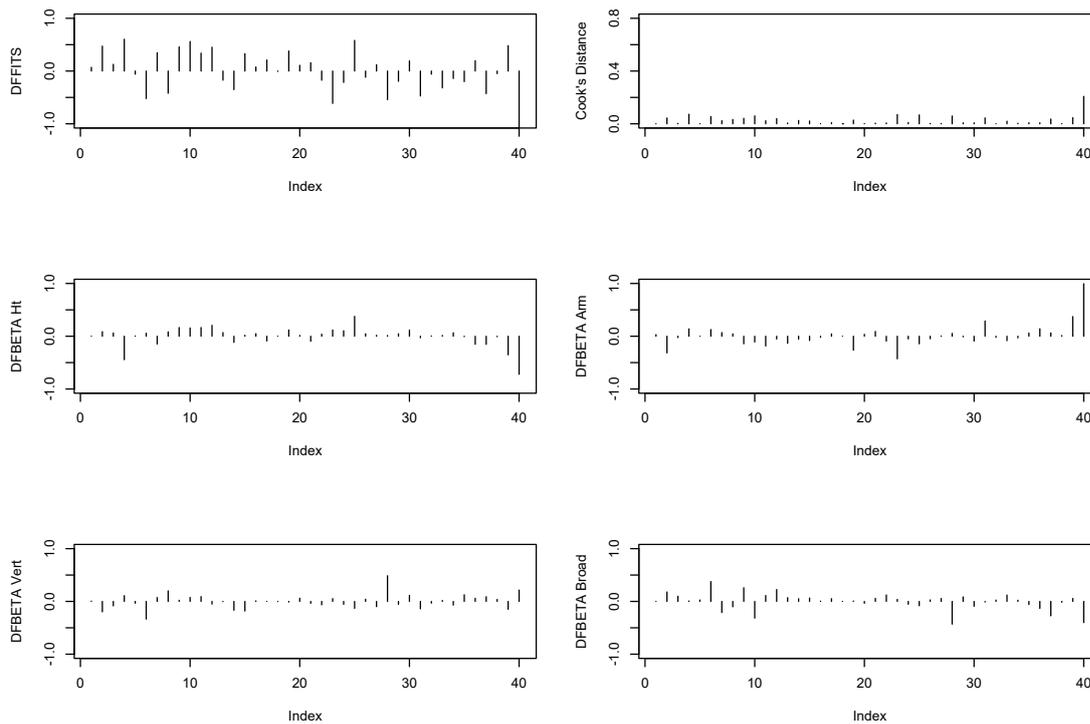


Figura 39: Análise gráfica de pontos influentes dos *Quarterbacks*

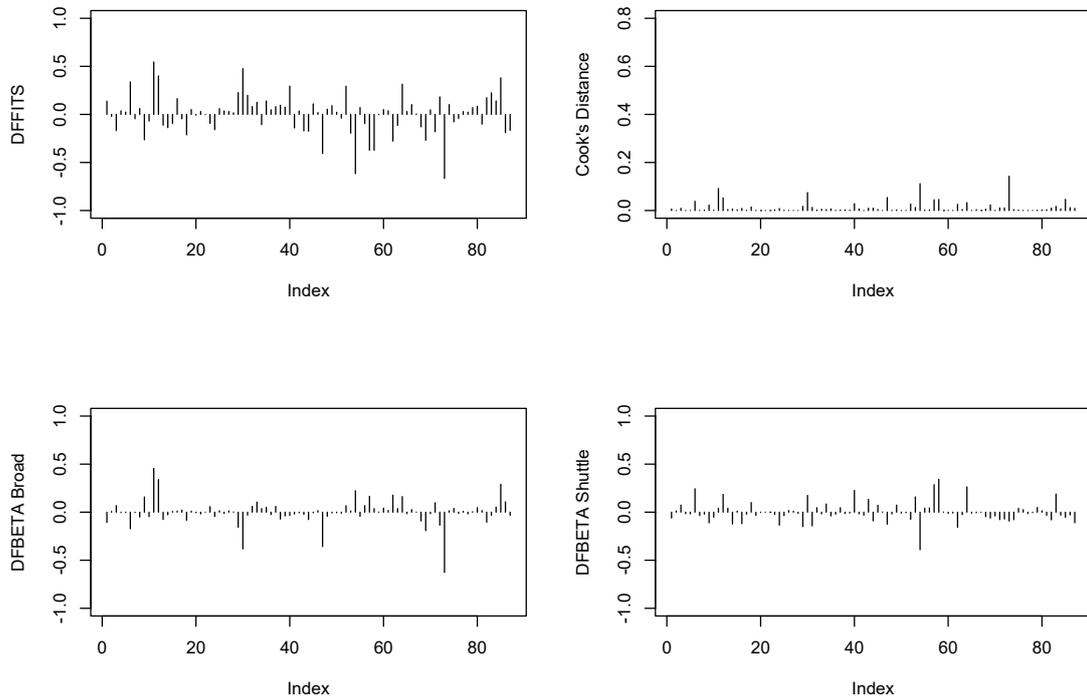


Figura 40: Análise gráfica de pontos influentes dos *Running backs*

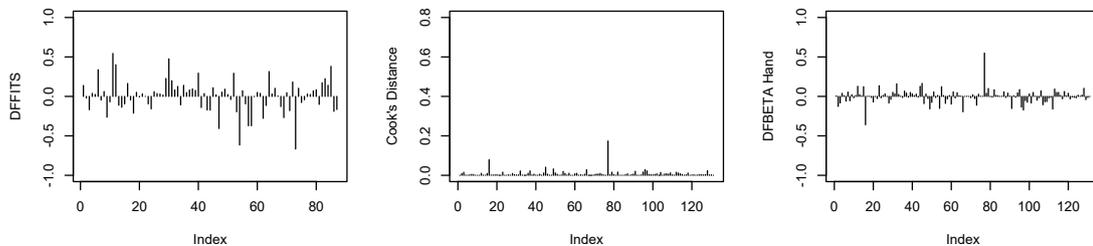


Figura 41: Análise gráfica de pontos influentes dos *Wide receivers*

Dessa maneira, a partir da análise das métricas, não há pontos de alavancagem na regressão, com regra básica de 0,8 para a distância de Cook e 1 para DFFITS e DFBETAS, exceto para o ponto de índice 40 na posição de *Quarterback*, pela métrica DFFITS, através da métrica DFBETAS e distância de Cook, contudo, esta observação não é um ponto de alavancagem.

4.2.3 Diagnóstico e comparação dos modelos

E então, os resultados dos modelos foram:

Tabela 30: Tabela de comparação de variáveis presentes nos modelos de regressão finais

Modelo	Variável	Estimativa	Erro padrão
<i>Skill Positions</i>	Ht	-0.1881	0.1033
	Hand	0.9538	0.4437
	Broad	0.1004	0.0373
	Shuttle	8.4188	4.1737
<i>Quarterbacks</i>	Ht	-1.3634	0.6977
	Arm	2.2417	1.1963
	Vert	-1.0327	0.4232
	Broad	0.6243	0.2185
<i>Running backs</i>	Broad	0.1200	0.0590
	Shuttle	11.7706	5.8884
<i>Wide receivers</i>	Hand	1.0610	0.5337

Pode-se perceber que, quando significativa, a altura foi inversamente proporcional ao desempenho do jogador. Os modelos dos *Running Backs* e dos *Wide Receivers* são sub-modelos do modelo de regressão geral, os primeiros, com a variável *Shuttle* consideravelmente mais influente no desempenho, e o *Broad* levemente. Para os recebedores, a variável mão foi levemente mais influente na *Grade* que para o modelo geral. Já os *Quarterbacks* obtiveram 2 variáveis em comum, ambas com mesmo sinal no β , contudo, para os passadores, ambas obtiveram maior influência na variável resposta. Estes últimos obtiveram também duas variáveis distintas dos outros modelos, que é o salto vertical, influente negativamente na nota, e o tamanho do braço, que é diretamente proporcional.

Tabela 31: Medidas de diagnóstico para as predições dos modelos finais

Modelo	PRESS	RMSE Treino	RMSE Teste
<i>Skill Positions</i>	5002.69	9.09	8.84
<i>Quarterbacks</i>	1012.98	12.148	10.61
<i>Running backs</i>	1132.60	7.697	7.34
<i>Wide receivers</i>	2559.02	7.917	8.94

Observa-se que os modelos de *Quarterbacks* e recebedores obtiveram maior percentual de explicação da variável resposta, contudo, a natureza da medida PRESS torna difícil a comparação das predições dos 3 modelos, pois ele não leva em consideração o tamanho amostral. Portanto, usando o erro quadrático médio, pode-se observar que o melhor desempenho na predição para o banco de teste foi o modelo de *Running Backs*. Os modelos posicionais, exceto pelo dos passadores, obtiveram resultados superiores ao modelo geral na predição dentro da amostra de teste, contudo, isso pode indicar apenas

uma variabilidade grande devido ao pequeno tamanho da amostra para essa posição, e não que o modelo ajustado foi pior, pois como já mostrado na análise dos modelos, apresentou um R^2 consideravelmente maior que os outros.

Por último, vale testar qual o melhor modelo entre eles, pelo critério AIC:

Tabela 32: AIC para os modelos de regressão

Modelo	AIC
<i>Skill Positions</i>	1140.10
<i>Quarterbacks</i>	209.78
<i>Running backs</i>	361.11
<i>Wide receivers</i>	546.11

Assim, confirma-se que o modelo para os *Quarterbacks*, apesar de ter tido previsões ruins no banco de teste, foi o melhor modelo pelo critério de Akaike e pelo R^2 , que mostrou um considerável aumento na explicação da variável desempenho em relação aos outros.

5 Conclusão

Este estudo abarcou uma amostra de jogadores de *Skill positions* que foram selecionados a partir do *Draft* da NFL entre 2007 e 2022, e além disso, que foram participantes do *NFL Combine* em suas respectivas edições. Com a finalidade de identificar se para as posições de *QB*, *WR*, *HB* existiu um padrão de se um ou mais testes, e quais, dentre os praticados nesse evento, apresentavam relação com o quão bem um jogador se torna na liga.

Como explicado mais cedo no trabalho, é de suma importância para todas as franquias da NFL alcançar um nível de conhecimento não só sobre a esfera de habilidade de um jogador com a bola, ou também de comprometimento e esforço, mas também um entendimento estatisticamente comprovado que, se porventura chegar a uma conclusão, irá tornar uma escolha no *Draft* muito mais valiosa, por apresentar menos erros, e selecionar realmente os jogadores que serão importantes para o time.

Assim, após a análise dos resultados, pôde-se notar que houve a significância de todos os modelos, tanto o geral quanto os posicionais, o que indica que há diferença nos desempenhos dos jogadores a partir dos dados dos testes físicos obtidos. Ao examinar os dados, também foi evidenciado que os *Quarterbacks* apresentaram um modelo mais explicativo do desempenho em comparação com os *Running backs* e *Wide receivers*. Embora os modelos tenham mostrado um ajuste intermediário ou até ruim, todos exibiram variáveis significativas, indicando que de certa maneira há relevância em analisar esses testes na avaliação do potencial de um jogador. A distinção observada entre os modelos ressalta a importância de analisar cuidadosamente e diferentemente jogadores de diferentes posições.

À luz dessas informações, é de se pensar que, ao se deparar com uma explicação tão baixa da variação do desempenho de um jogador de habilidade pelos modelos, sejam eles gerais ou posicionais, o que causaria tal variação? Existem muitos outros fatores que podem descrever a diferença entre um jogador que vai parar no *Hall* da fama da liga, ou um jogador que não se estabelece, talvez que nunca chegue realmente a jogar na liga após ser selecionado. Dentre eles pode-se pensar em uma evolução física tardia, ou em um jogador com habilidades intangíveis, não mensuráveis em nenhum teste, alguém muito dedicado, um bom companheiro de equipe e capitão, ou que tem uma ótima leitura de jogo, e que pode sempre estar a frente de seu adversário no entendimento das jogadas. Todos os fatores citados podem se sobrepor ao físico, na verdade, todos somados tornam cada jogador tão diferente e é o que traz muitos times a terem avaliação tão distintas sobre eles, o que cada dono, técnico ou companheiro avalia num jogador é totalmente pessoal.

Assim, dado o nível mínimo necessário para estar no *Combine*, e consequentemente ser candidato a estar na liga, os testes afinal não apresentaram tanta importância

de acordo com o estudo. Por isso, alguns estigmas físicos podem perdurar por muitos e muitos anos por não se ter conhecimento suficiente para refutá-los. Se um ou dois jogadores em um ano apresentam um padrão físico e não apresentam um bom desempenho, talvez no futuro algum time perca a oportunidade de selecionar uma estrela por preferir alguém semelhante aos anteriores. Quem sabe, alguma futura estrela prometida pode ter sua vida na liga reduzida após uma lesão, pois muitos fatores podem reduzir a durabilidade de um jogador numa liga tão exigente e perigosa.

Desse modo, em uma última análise, este estudo destaca a complexidade de se ler jogadores que ainda nem iniciaram sua carreira profissional, e reforça a importância de adotar uma abordagem abrangente ao considerar a seleção e o desenvolvimento de jogadores. Compreender as nuances das diferentes posições e reconhecer a variedade enorme de fatores que contribuem para o sucesso de um jogador é fundamental para que as franquias tomem, na noite do *Draft*, decisões certas.

Para estudos futuros, é sugerido analisar dados semelhantes utilizando diferentes técnicas estatísticas, como por exemplo modelos de resposta gradual. Além disso, é válido considerar modificações no escopo do estudo em relação aos anos analisados de um jogador, como avaliar o impacto dos testes físicos no ano de ingresso na liga, ou avaliar exclusivamente o auge de um jogador. Uma outra abordagem seria filtrar jogadores com um mínimo de jogos ou temporadas na *NFL*, entre outras estratégias para explorar esses dados de forma mais completa.

Referências

- AKAIKE, H. A new look at the statistical model identification. *IEEE transactions on automatic control*, Ieee, v. 19, n. 6, p. 716–723, 1974.
- BEAULIEU-JONES, B. R. et al. Epidemiology of injuries identified at the nfl scouting combine and their impact on performance in the national football league: evaluation of 2203 athletes from 2009 to 2015. *Orthopaedic journal of sports medicine*, SAGE Publications Sage CA: Los Angeles, CA, v. 5, n. 7, p. 2325967117708744, 2017.
- BERRI, D. J.; SIMMONS, R. Catching a draft: On the process of selecting quarterbacks in the national football league amateur draft. *Journal of Productivity Analysis*, Springer, v. 35, p. 37–49, 2011.
- CHAI, T.; DRAXLER, R. R. Root mean square error (rmse) or mean absolute error (mae)?—arguments against avoiding rmse in the literature. *Geoscientific model development*, Copernicus Publications Göttingen, Germany, v. 7, n. 3, p. 1247–1250, 2014.
- COOK, J. et al. The relationship between the national football league scouting combine and game performance over a 5-year period. *The Journal of Strength & Conditioning Research*, LWW, v. 34, n. 9, p. 2492–2499, 2020.
- HEDLUND, D. P. Performance of future elite players at the national football league scouting combine. *The Journal of Strength & Conditioning Research*, LWW, v. 32, n. 11, p. 3112–3118, 2018.
- KOZ, D.; FRASER-THOMAS, J.; BAKER, J. Accuracy of professional sports drafts in predicting career potential. *Scandinavian journal of medicine & science in sports*, Wiley Online Library, v. 22, n. 4, p. e64–e69, 2012.
- KUZMITS, F. E.; ADAMS, A. J. The nfl combine: does it predict performance in the national football league? *The Journal of strength & conditioning research*, LWW, v. 22, n. 6, p. 1721–1727, 2008.
- MONTGOMERY, D. C.; PECK, E. A.; VINING, G. G. *Introduction to linear regression analysis*. [S.l.]: John Wiley & Sons, 2021.
- MORETTIN, P. A.; BUSSAB, W. O. *Estatística básica*. [S.l.]: Saraiva Educação SA, 2017.
- NETER, J. et al. *Applied linear statistical models*. Irwin Chicago, 1996.
- PPF. *NFL Players Grade PFF*. 2023. Disponível em: <<https://www.pff.com>>.
- ROBBINS, D. W. The national football league (nfl) combine: does normalized data better predict performance in the nfl draft? *The Journal of Strength & Conditioning Research*, LWW, v. 24, n. 11, p. 2888–2899, 2010.
- SUNDAY, S. V. *NFL Combine Data*. 2019. Disponível em: <<https://data.world/sportsvizsunday/nfl-combine-data/workspace/file?filename=NFL+Combine+Data.xlsx>>.

TERAMOTO, M.; CROSS, C. L.; WILLICK, S. E. Predictive value of national football league scouting combine on future performance of running backs and wide receivers. *The Journal of Strength & Conditioning Research*, LWW, v. 30, n. 5, p. 1379–1390, 2016.

TUKEY, J. W. et al. *Exploratory data analysis*. [S.l.]: Reading, MA, 1977. v. 2.