


```
In [45]: # Calculating league average Medium Danger Shots per game
league_mdshots_per_game = (Monte_carlo_df['MDSH%'].sum()/(team_df_21['GP'].sum()))

Out[45]: 6.32259064516129

In [46]: # Calculating league average Low Danger Shots per game
league_ldshots_per_game = (Monte_carlo_df['LDSH%'].sum()/(team_df_21['GP'].sum()))

Out[46]: 10.099078341013926

In [47]: # Calculating league average Medium Danger Shooting Percentage
league_avg_md_shoot_pct = (round((Monte_carlo_df['MDSH%']).mean(),2))
league_avg_md_shoot_pct

Out[47]: 9.15

In [48]: # Calculating league average High Danger Shooting Percentage
league_avg_hd_shoot_pct = (round((Monte_carlo_df['HDSH%']).mean(),2))
league_avg_hd_shoot_pct

Out[48]: 18.0

In [49]: # Calculating league average Low Danger Shooting Percentage
league_avg_ld_shoot_pct = (round((Monte_carlo_df['LDSH%']).mean(),2))
league_avg_ld_shoot_pct

Out[49]: 2.7
```

Now I want to do a little verification that changing the index to the team names will allow me to easily select the statistics from each team inside the new function.

```
In [50]: # Using .loc to select all the stats from the Sharks
Monte_carlo_df.loc['San Jose Sharks']

Out[50]:
SHR      8.05
SV%      90.37
SF       1305.00
SA       1391.00
LDSF     489.00
LDSA     595.00
LDSV%    96.13
LDSH%     3.07
MDSF     368.00
MDSA     371.00
MDSV%    90.03
MDSH%     7.61
HDSF     403.00
HDSA     367.00
HDSV%    79.84
HDSH%    15.38
Name: San Jose Sharks, dtype: float64

In [51]: # Again using .loc for the stats from Vegas.
Monte_carlo_df.loc['Vegas Golden Knights']

Out[51]:
SHR      8.58
SV%      92.15
SF       1409.00
SA       1235.00
LDSF     553.00
LDSA     490.00
LDSV%    97.55
LDSH%     3.62
MDSF     411.00
MDSA     319.00
MDSV%    91.22
MDSH%    10.95
HDSF     386.00
HDSA     353.00
HDSV%    83.85
HDSH%    17.88
Name: Vegas Golden Knights, dtype: float64

In [52]: # Now parsing the stats and dividing by the number of games played to see how many Low Danger Shots For Vegas <
Monte_carlo_df.loc['Vegas Golden Knights']['LDSH']/56

Out[52]: 9.875

In [53]: # Calculating how much worse the Sharks Medium Save Percentage was than the league average
Monte_carlo_df.loc['San Jose Sharks']['MDSV%'] - league_avg_mdsv_pct

Out[53]: -0.81999999999999932

In [54]: # Calling Vegas' Medium Danger Shooting Percentage
Monte_carlo_df.loc['Vegas Golden Knights']['MDSH%']

Out[54]: 10.95

In [55]: # By subtracting the Medium Danger Shooting Percentage from 100
# I am able to calculate the Medium Danger Save Percentage against Vegas
100 - Monte_carlo_df.loc['Vegas Golden Knights']['MDSH%']

Out[55]: 89.05

In [56]: # Calculating how much better Vegas' Medium Danger Shooting Percentage is than league average.
Monte_carlo_df.loc['Vegas Golden Knights']['MDSH%'] - league_avg_md_shoot_pct

Out[56]: 1.7999999999999999

In [57]: # Calculating how much better Vegas shoots than league average on Medium Danger shots.
(100-Monte_carlo_df.loc['Vegas Golden Knights']['MDSH%']) - (100-league_avg_md_shoot_pct)

Out[57]: -1.7999999999999992
```

Building a More Robust Monte Carlo Simulation:

Now that I have my new dataframe it can be used to calculate the expectations within the new function.

```
In [58]: # Creating the function, the only require information will be the team names.
def monte_carlo_matchup(team1, team2):
    # Entering an assert statement to make sure the team name entered is a valid NHL team.
    assert team1 in Monte_carlo_df.index, 'Team1 is not a valid NHL team name. Please enter a valid NHL team.'
    assert team2 in Monte_carlo_df.index, 'Team2 is not a valid NHL team name. Please enter a valid NHL team.'

    # Setting up the connection between the entered team name and the data in the dataframe
    team1 = Monte_carlo_df.loc[team1]
    team2 = Monte_carlo_df.loc[team2]

    # Setting up a connection to be used in the calculations so the data can be extracted from the team name
    # with a .loc
    team1_name = team1.name
    team2_name = team2.name

    # Calculating the league average statistics
    league_avg_ldsv_pct = (round((Monte_carlo_df['LDSV%'].mean()),2))
    league_avg_mdsv_pct = (round((Monte_carlo_df['MDSV%'].mean()),2))
    league_avg_ldsv_pct = (round((Monte_carlo_df['LDSV%'].mean()),2))
    # Here I used the sum of GP from my original dataframe
    league_mdshots_per_game = (Monte_carlo_df['MDSF'].sum()/(team_df_21['GP'].sum()))
    league_ldshots_per_game = (Monte_carlo_df['LDSF'].sum()/(team_df_21['GP'].sum()))
    league_avg_md_shoot_pct = (round((Monte_carlo_df['MDSH%']).mean(),2))
    league_avg_ld_shoot_pct = (round((Monte_carlo_df['LDSH%']).mean(),2))

    # Calculating the expected stats for each team based on their opponent and vice versa
    team1_expected_ldsv_pct = ((Monte_carlo_df.loc[team1_name]['LDSV%'] - league_avg_ldsv_pct) +
                               ((100-Monte_carlo_df.loc[team1_name]['LDSH%']) - (100-league_avg_md_shoot_pct)))
    team2_expected_ldsv_pct = ((Monte_carlo_df.loc[team2_name]['LDSV%'] - league_avg_ldsv_pct) +
                               ((100-Monte_carlo_df.loc[team2_name]['LDSH%']) - (100-league_avg_md_shoot_pct)))

    team1_expected_mdsv_pct = ((Monte_carlo_df.loc[team1_name]['MDSV%'] - league_avg_mdsv_pct) +
                               ((100-Monte_carlo_df.loc[team1_name]['MDSH%']) - (100-league_avg_md_shoot_pct)))
    team2_expected_mdsv_pct = ((Monte_carlo_df.loc[team2_name]['MDSV%'] - league_avg_mdsv_pct) +
                               ((100-Monte_carlo_df.loc[team2_name]['MDSH%']) - (100-league_avg_md_shoot_pct)))

    team1_expected_hdsv_pct = ((Monte_carlo_df.loc[team1_name]['HDSV%'] - league_avg_hdsv_pct) +
                               ((100-Monte_carlo_df.loc[team1_name]['HDSH%']) - (100-league_avg_hd_shoot_pct)))
    team2_expected_hdsv_pct = ((Monte_carlo_df.loc[team2_name]['HDSV%'] - league_avg_hdsv_pct) +
                               ((100-Monte_carlo_df.loc[team2_name]['HDSH%']) - (100-league_avg_hd_shoot_pct)))

    team1_expected_ldshots_game = (round((((Monte_carlo_df.loc[team1_name]['LDSF']/56) - league_ldshots_per_game) +
                                             league_ldshots_per_game),0)).astype(int)
    team2_expected_ldshots_game = (round((((Monte_carlo_df.loc[team2_name]['LDSF']/56) - league_ldshots_per_game) +
                                             league_ldshots_per_game),0)).astype(int)

    team1_expected_mdshots_game = (round((((Monte_carlo_df.loc[team1_name]['MDSF']/56) - league_mdshots_per_game) +
                                             league_mdshots_per_game),0)).astype(int)
    team2_expected_mdshots_game = (round((((Monte_carlo_df.loc[team2_name]['MDSF']/56) - league_mdshots_per_game) +
                                             league_mdshots_per_game),0)).astype(int)

    team1_expected_hdshots_game = (round((((Monte_carlo_df.loc[team1_name]['HDSF']/56) - league_hdshots_per_game) +
                                             league_hdshots_per_game),0)).astype(int)
    team2_expected_hdshots_game = (round((((Monte_carlo_df.loc[team2_name]['HDSF']/56) - league_hdshots_per_game) +
                                             league_hdshots_per_game),0)).astype(int)

    # Creating tallies for all of the stats I want to track
    total_goals_team1 = 0
    goals_scored_onhd_team1 = 0
    goals_scored_onmd_team1 = 0
    goals_scored_onld_team1 = 0
    goals_scored_onhd_team2 = 0
    goals_scored_onmd_team2 = 0
    goals_scored_onld_team2 = 0
    team1_victories = 0
    team2_victories = 0
    overtime = 0
    # Creating a for loop that runs 100,000 simulations of the game
    for i in range(0, 100000):
        # Goals scored for each team
        goals_scored_team1 = 0
        goals_scored_team2 = 0
        # A for loop for team1 shot attempts and outcomes
        for k in range(0, team1_expected_ldshots_game + 1):
            # using np.random.random to generate random floats between 1 and 0
            # If the number hits above the save percentage threshold then a goal is scored
            # otherwise the shot is saved by the goalie.
            r = np.random.random()
            if(r>=team1_expected_ldsv_pct):
                # Adding a goal if the shot falls below the set shooting percentage
                goals_scored_team1 += 1
                goals_scored_onld_team1 += 1
                total_goals_team1 += 1
            else:
                # Shot is saved, no goal is added
                goals_scored_team1 += 0

        for k in range(0, team1_expected_mdshots_game + 1):
            # using np.random.random to generate random floats between 1 and 0
            # If the number hits above the save percentage threshold then a goal is scored
            # otherwise the shot is saved by the goalie.
            r = np.random.random()
            if(r>=team1_expected_mdsv_pct):
                # Adding a goal if the shot falls below the set shooting percentage
                goals_scored_team1 += 1
                goals_scored_onmd_team1 += 1
                total_goals_team1 += 1
            else:
                # Shot is saved, no goal is added
                goals_scored_team1 += 0

        for k in range(0, team1_expected_hdshots_game + 1):
            # using np.random.random to generate random floats between 1 and 0
            # If the number hits above the save percentage threshold then a goal is scored
            # otherwise the shot is saved by the goalie.
            r = np.random.random()
            if(r>=team1_expected_hdsv_pct):
                # Adding a goal if the shot falls below the set shooting percentage
                goals_scored_team1 += 1
                goals_scored_onhd_team1 += 1
                total_goals_team1 += 1
            else:
                # Shot is saved, no goal is added
                goals_scored_team1 += 0

        # Running the same for loop for team 2
        r = np.random.random()
        if(r>=team2_expected_ldsv_pct):
            goals_scored_team2 += 1
            goals_scored_onld_team2 += 1
            total_goals_team2 += 1
        else:
            goals_scored_team2 += 0

        for k in range(0, team2_expected_mdshots_game + 1):
            # Running the same for loop for team 2
            r = np.random.random()
            if(r>=team2_expected_mdsv_pct):
                goals_scored_team2 += 1
                goals_scored_onmd_team2 += 1
                total_goals_team2 += 1
            else:
                goals_scored_team2 += 0

        for k in range(0, team2_expected_hdshots_game + 1):
            # Running the same for loop for team 2
            r = np.random.random()
            if(r>=team2_expected_hdsv_pct):
                goals_scored_team2 += 1
                goals_scored_onhd_team2 += 1
                total_goals_team2 += 1
            else:
                goals_scored_team2 += 0

    # Calculating which team won the game
    if team1_scored more, they get a victory
    if goals_scored_team1 > goals_scored_team2:
        # The victory is added to their total victories
        team1_victories += 1
    # The same is done for team 2
    if goals_scored_team2 > goals_scored_team1:
        team2_victories += 1
    # the score is tied, the game goes into overtime
    if goals_scored_team1 == goals_scored_team2:
        overtime += 1

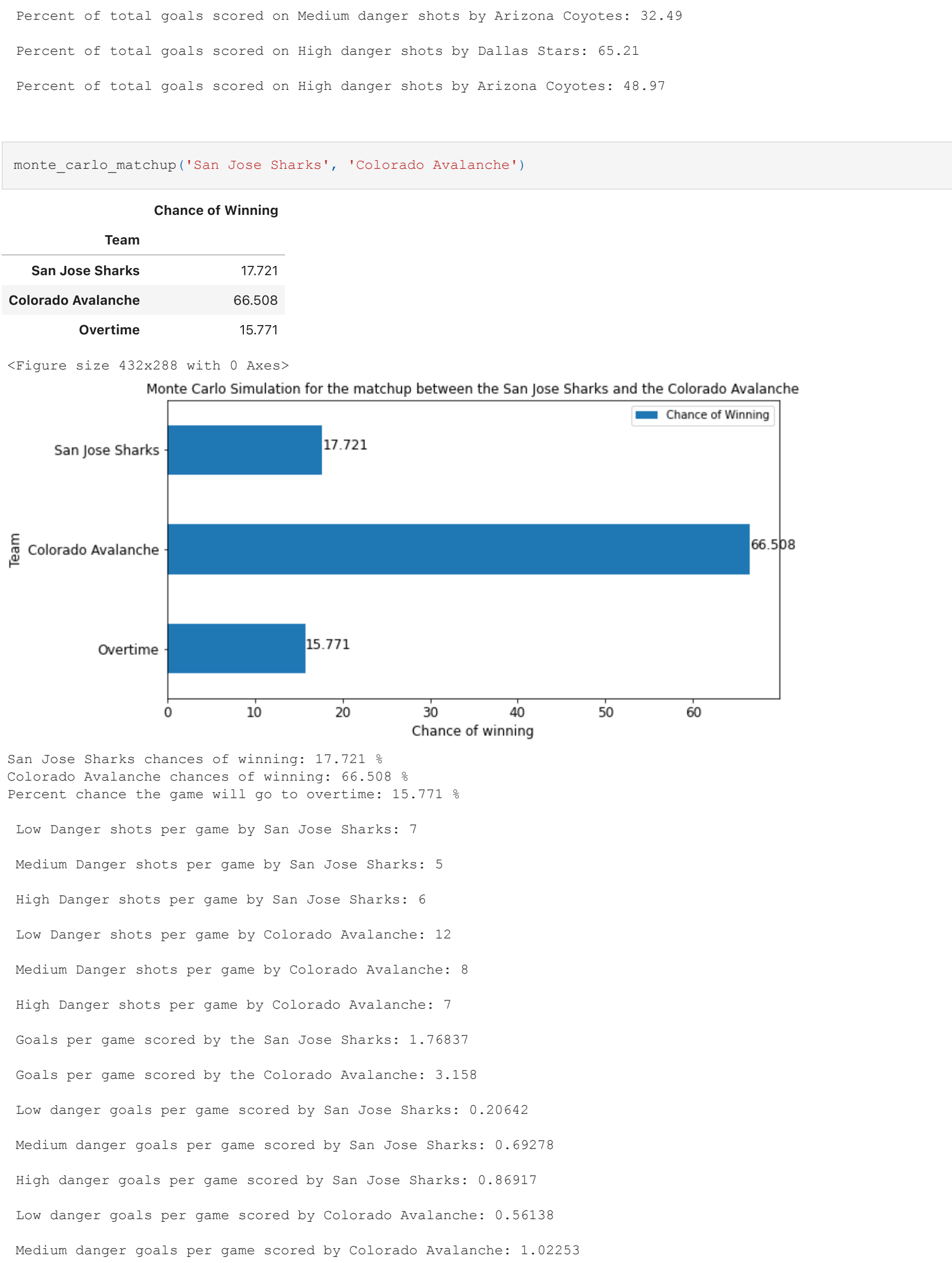
    # Creating a dataframe to make the odds of winning easier to see for each team
    team1_chance = team1_victories/1000
    team2_chance = team2_victories/1000
    overtime_chance = overtime/1000
    # Setting up the columns and rows with the data from the simulation
    odds = ['Team1', [team1_name, team2_name, 'Overtime'], 'Chance of Winning'], [team1_chance, team2_chance, overtime_chance]
    # Creating a dataframe from the new data
    odds_df = pd.DataFrame(data = odds)
    # Setting the index as the Team column
    odds_df = odds_df.set_index('Team')
    # Displaying the dataframe
    display(odds_df)

    # Creating a plot using a horizontal bar chart so the results are easy to process
    plt.figure()
    odds_df.loc[0:3].plot(kind='bar', figsize=(10,3), fontsize=12, xticks = (np.arange(0, 70, step=10)))
    plt.grid(True,linestyle='solid')
    plt.xlabel('Chance of winning', fontsize=12)
    plt.ylabel('Team', fontsize=12)
    for index, value in enumerate(odds_df['Chance of Winning']):
        plt.text(value, index, str(value), fontsize=12)
    plt.title(f'Monte Carlo Simulation for the matchup between the {team1_name} and the {team2_name}')
    plt.show()

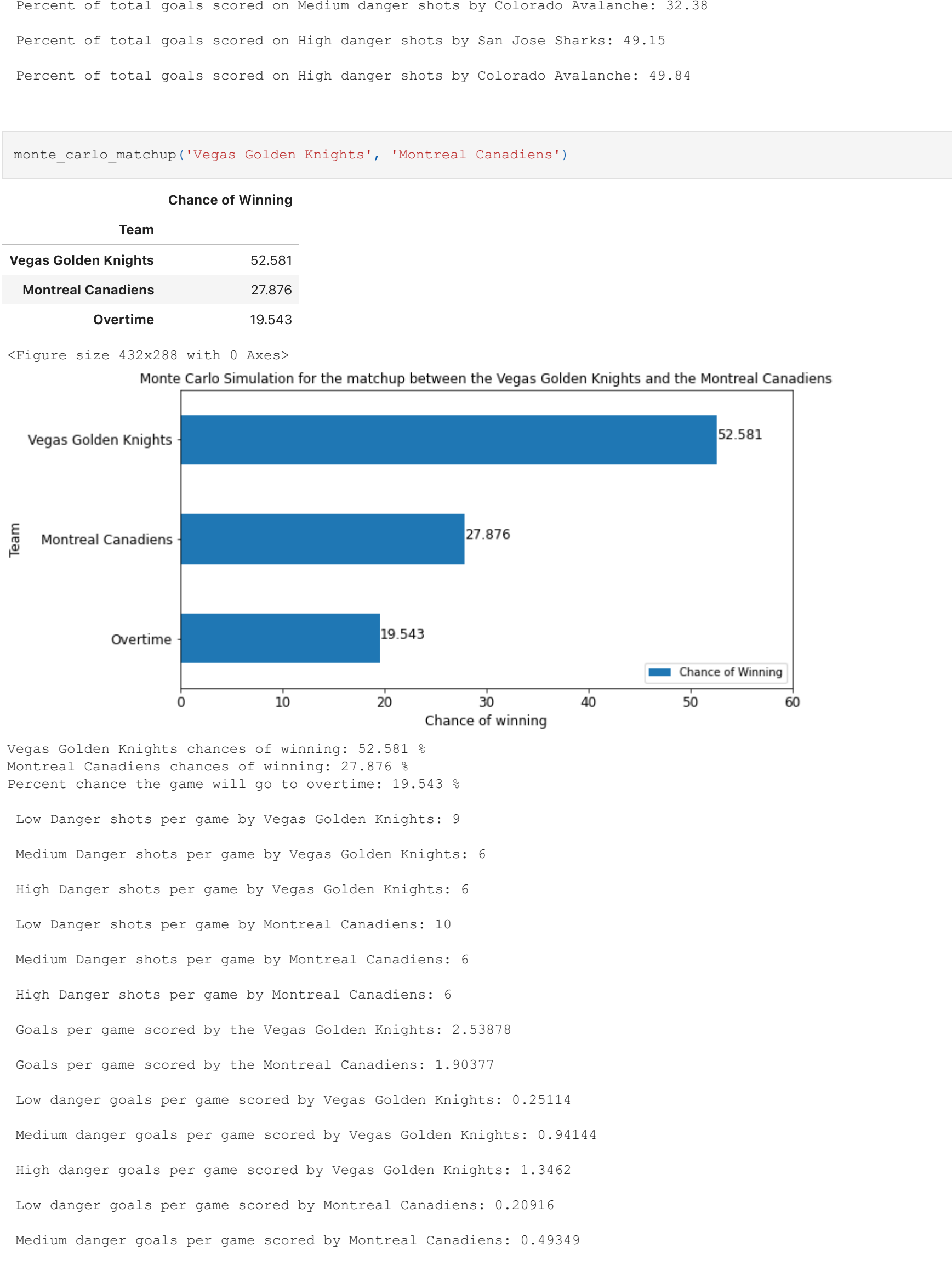
    # Printing the percent chance of team 1 winning, team 2 winning, and the game going to overtime
    as well as the odds of winning, team 1 winning, team 2 winning, the goals scored in the different danger areas per
    # and the percentage of goals scored by danger location
    print(f'{team1_name} chances of winning: {team1_victories/1000:.1%}', '\n',
          f'{team2_name} chances of winning: {team2_victories/1000:.1%}', '\n',
          f'Percent chance the game will go to overtime: {overtime/1000:.1%}', '\n',
          '\n',
          f'Low Danger shots per game by {team1_name}': team1_expected_ldshots_game, '\n', '\n',
          f'Medium Danger shots per game by {team1_name}': team1_expected_mdshots_game, '\n', '\n',
          f'High Danger shots per game by {team1_name}': team1_expected_hdshots_game, '\n', '\n',
          f'Low Danger shots per game by {team2_name}': team2_expected_ldshots_game, '\n', '\n',
          f'Medium Danger shots per game by {team2_name}': team2_expected_mdshots_game, '\n', '\n',
          f'High Danger shots per game by {team2_name}': team2_expected_hdshots_game, '\n', '\n',
          f'Goals per game scored by the {team1_name}': total_goals_team1/100000, '\n', '\n',
          f'Goals per game scored by the {team2_name}': total_goals_team2/100000, '\n', '\n',
          f'Low danger goals per game scored by {team1_name}': goals_scored_onld_team1/100000, '\n', '\n',
          f'Medium danger goals per game scored by {team1_name}': goals_scored_onmd_team1/100000, '\n', '\n',
          f'High danger goals per game scored by {team1_name}': goals_scored_onhd_team1/100000, '\n', '\n',
          f'Low danger goals per game scored by {team2_name}': goals_scored_onld_team2/100000, '\n', '\n',
          f'Medium danger goals per game scored by {team2_name}': goals_scored_onmd_team2/100000, '\n', '\n',
          f'High danger goals per game scored by {team2_name}': goals_scored_onhd_team2/100000, '\n', '\n',
          f'Percent of total goals scored on Low danger shots by {team1_name}': round((goals_scored_onld_team1/total_goals_team1), 2),
          f'Percent of total goals scored on Low danger shots by {team2_name}': round((goals_scored_onld_team2/total_goals_team2), 2),
          f'Percent of total goals scored on Medium danger shots by {team1_name}': round((goals_scored_onmd_team1/total_goals_team1), 2),
          f'Percent of total goals scored on Medium danger shots by {team2_name}': round((goals_scored_onmd_team2/total_goals_team2), 2),
          f'Percent of total goals scored on High danger shots by {team1_name}': round((goals_scored_onhd_team1/total_goals_team1), 2),
          f'Percent of total goals scored on High danger shots by {team2_name}': round((goals_scored_onhd_team2/total_goals_team2), 2))
```

Now that we've built the function, let's test it by running it with several different teams.

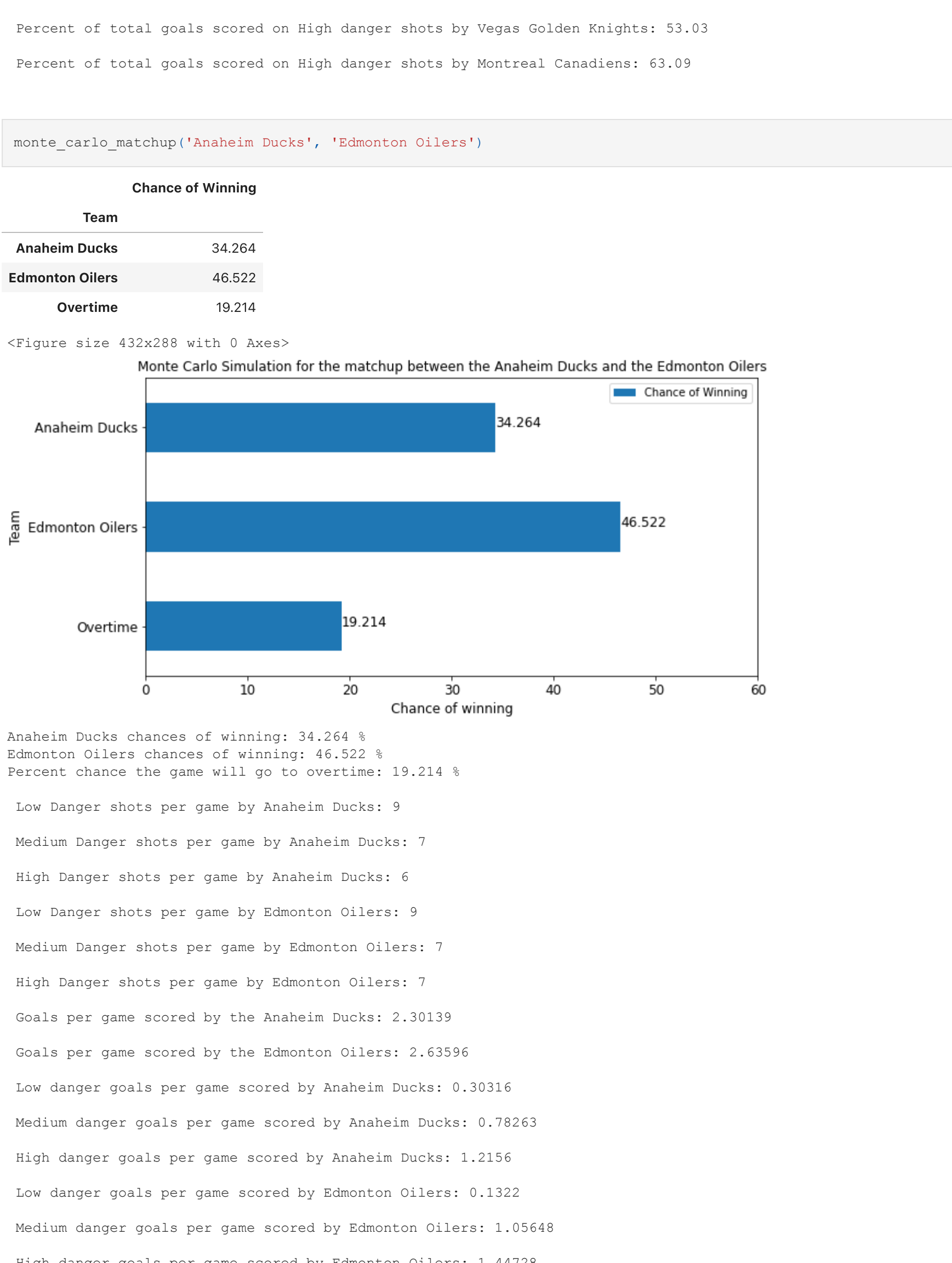
```
In [59]: monte_carlo_matchup('Washington Capitals', 'Boston Bruins')
```



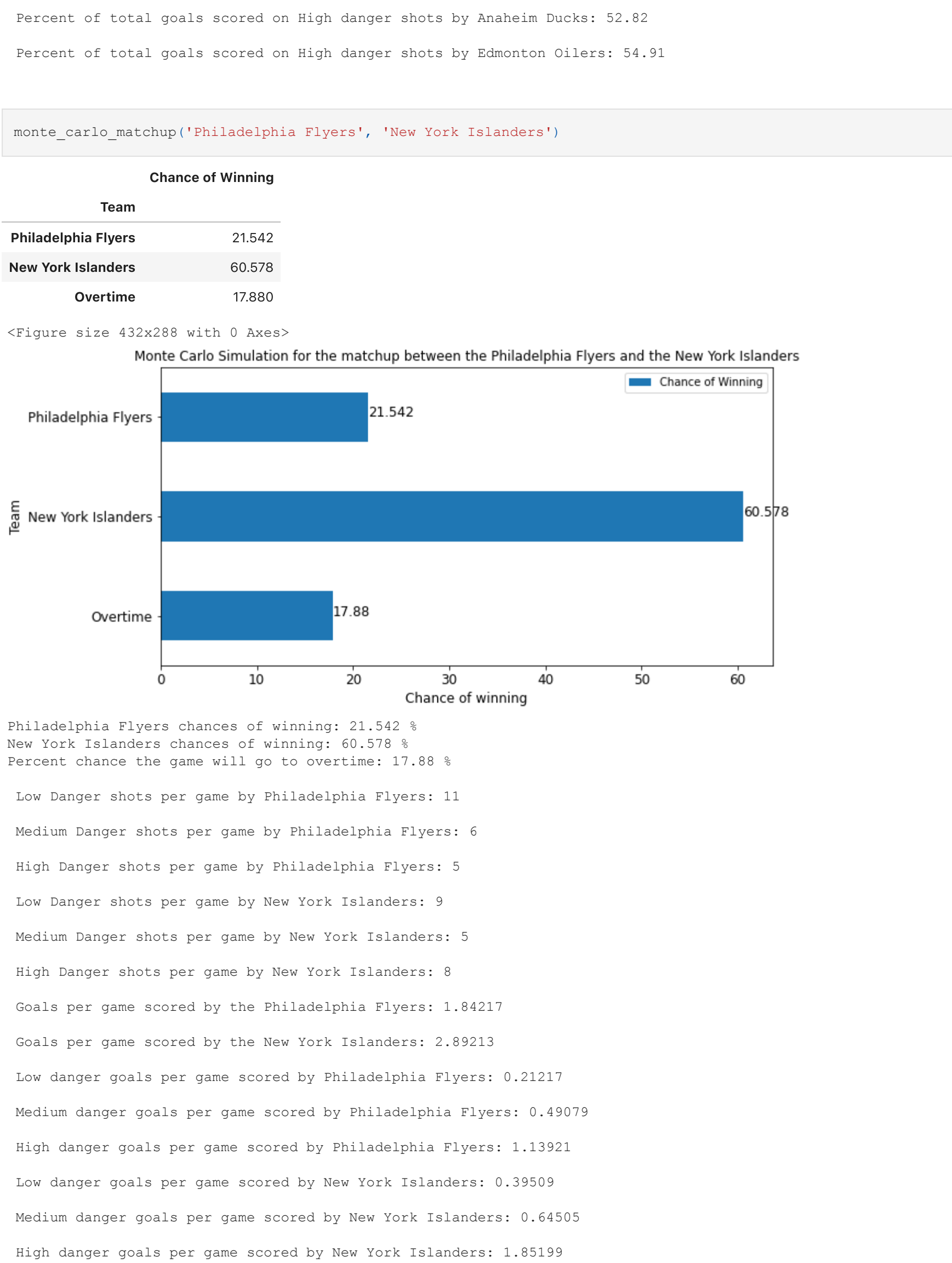
```
In [60]: monte_carlo_matchup('Dallas Stars', 'Arizona Coyotes')
```



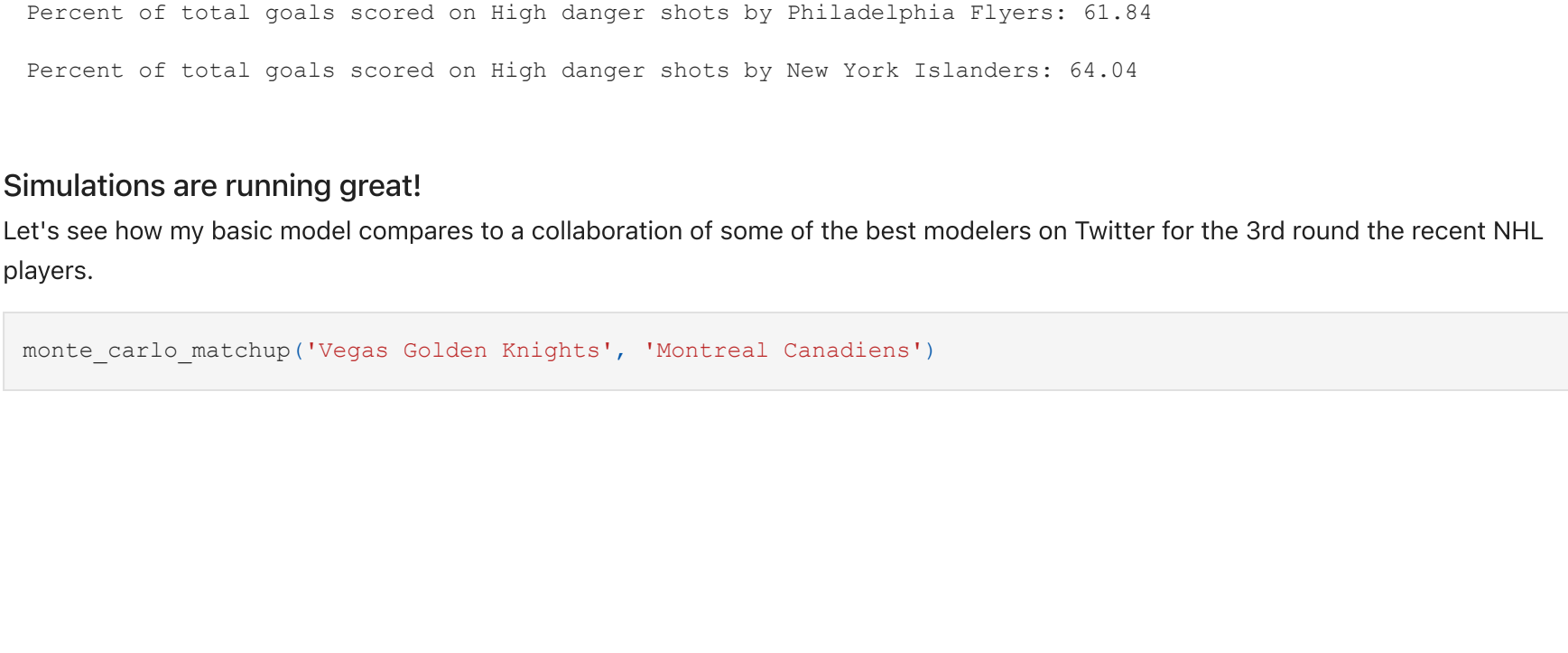
```
In [61]: monte_carlo_matchup('San Jose Sharks', 'Colorado Avalanche')
```



```
In [62]: monte_carlo_matchup('Vegas Golden Knights', 'Montreal Canadiens')
```



```
In [63]: monte_carlo_matchup('Anaheim Ducks', 'Edmonton Oilers')
```

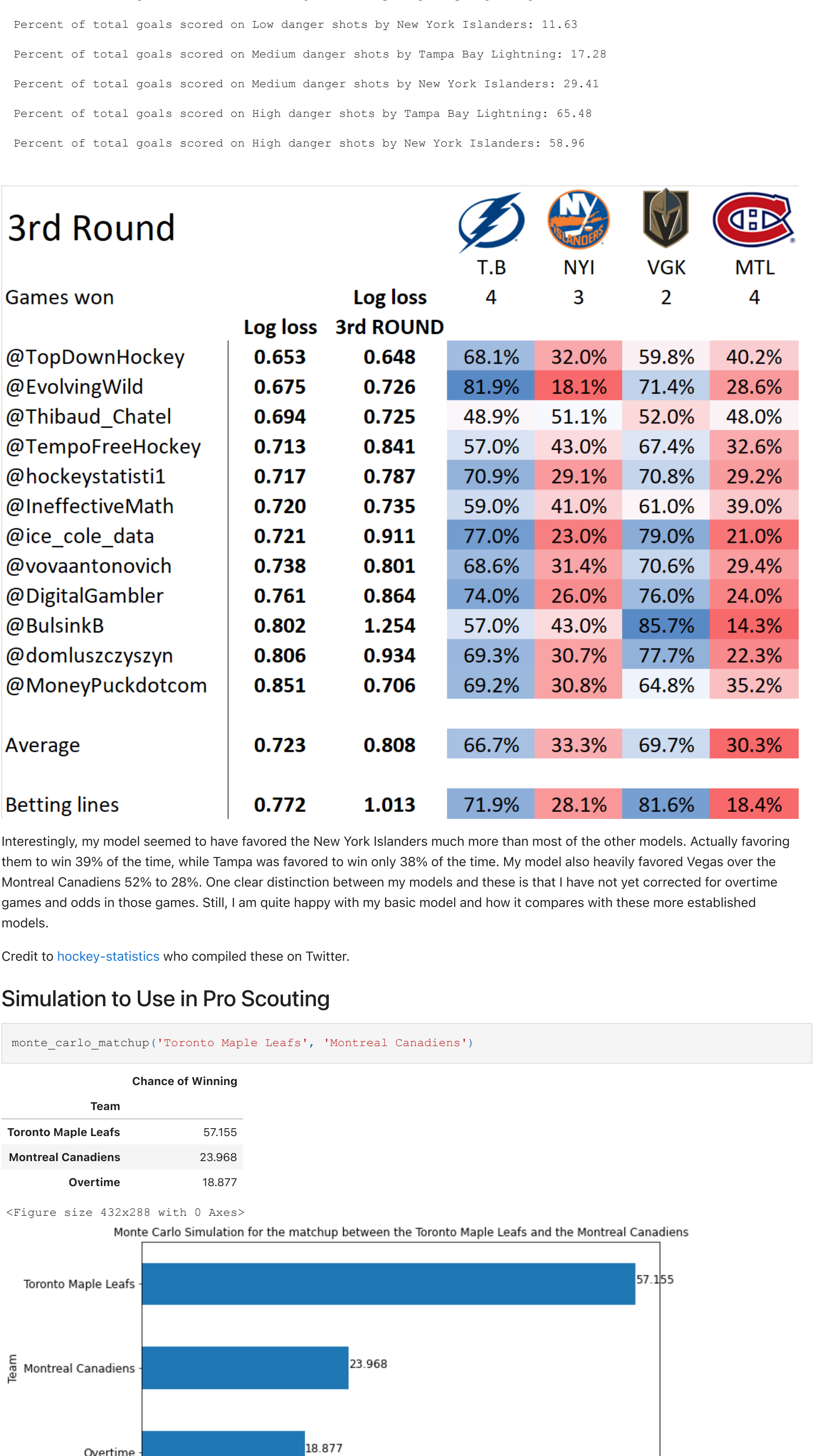
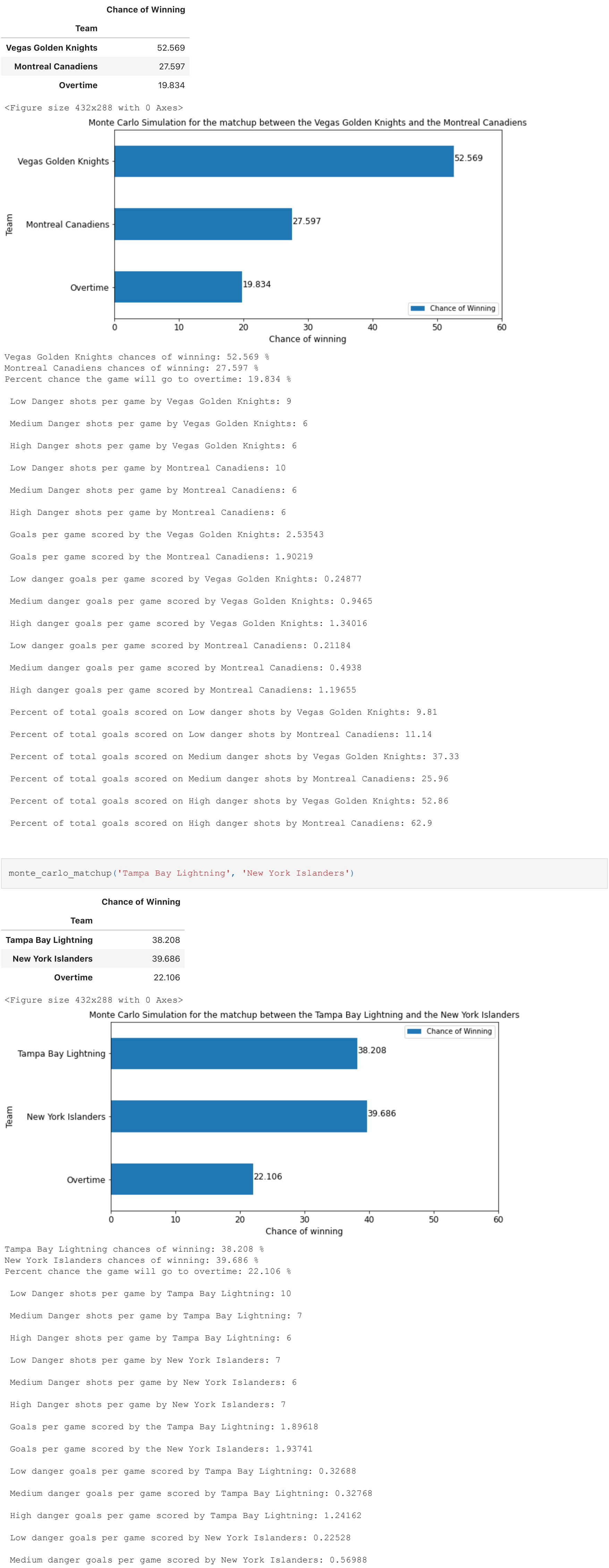


```
In [64]: monte_carlo_matchup('Philadelphia Flyers', 'New York Islanders')
```


Simulations are running great!

Let's see how my basic model compares to a collaboration of some of the best modelers on Twitter for the 3rd round the recent NHL players.

```
In [65]: monte_carlo_matchup('Vegas Golden Knights', 'Montreal Canadiens')
```

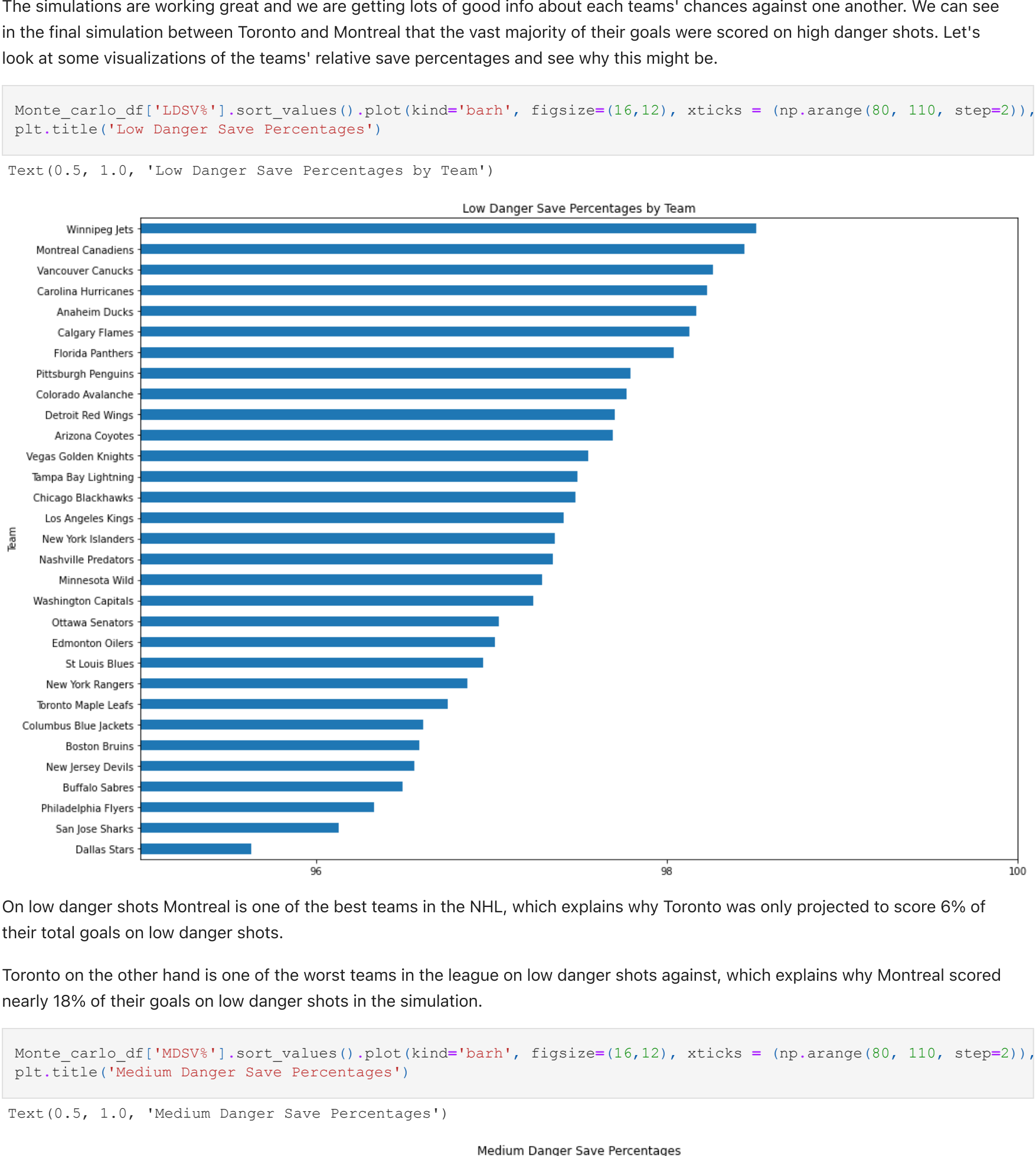



	0	10	20	30	40	50	60
Chance of winning							
Toronto Maple Leafs chances of winning:	57.155	4					
Montreal Canadiens chances of winning:	23.968	3					
Percent chance the game will go to overtime:	18.877	4					
Low Danger shots per game by Toronto Maple Leafs:	8						
Medium Danger shots per game by Toronto Maple Leafs:	7						
High Danger shots per game by Toronto Maple Leafs:	7						
Low Danger shots per game by Montreal Canadiens:	11						
Medium Danger shots per game by Montreal Canadiens:	6						
High Danger shots per game by Montreal Canadiens:	5						
Goals per game scored by the Toronto Maple Leafs:	2.69001						
Goals per game scored by the Montreal Canadiens:	1.82274						
Low danger goals per game scored by Toronto Maple Leafs:	0.18384						
Medium danger goals per game scored by Toronto Maple Leafs:	0.99719						
High danger goals per game scored by Toronto Maple Leafs:	1.50898						
Low danger goals per game scored by Montreal Canadiens:	0.32552						
Medium danger goals per game scored by Montreal Canadiens:	0.38858						

Interestingly, my model seemed to have favored the New York Islanders much more than most of the other models. Actually favoring them to win 39% of the time, while Tampa was favored to win only 38% of the time. My model also heavily favored Vegas over the Montreal Canadiens 52% to 28%. One clear distinction between my models and these is that I have not yet corrected for overtime games and odds in those games. Still, I am quite happy with my basic model and how it compares with these more established models.

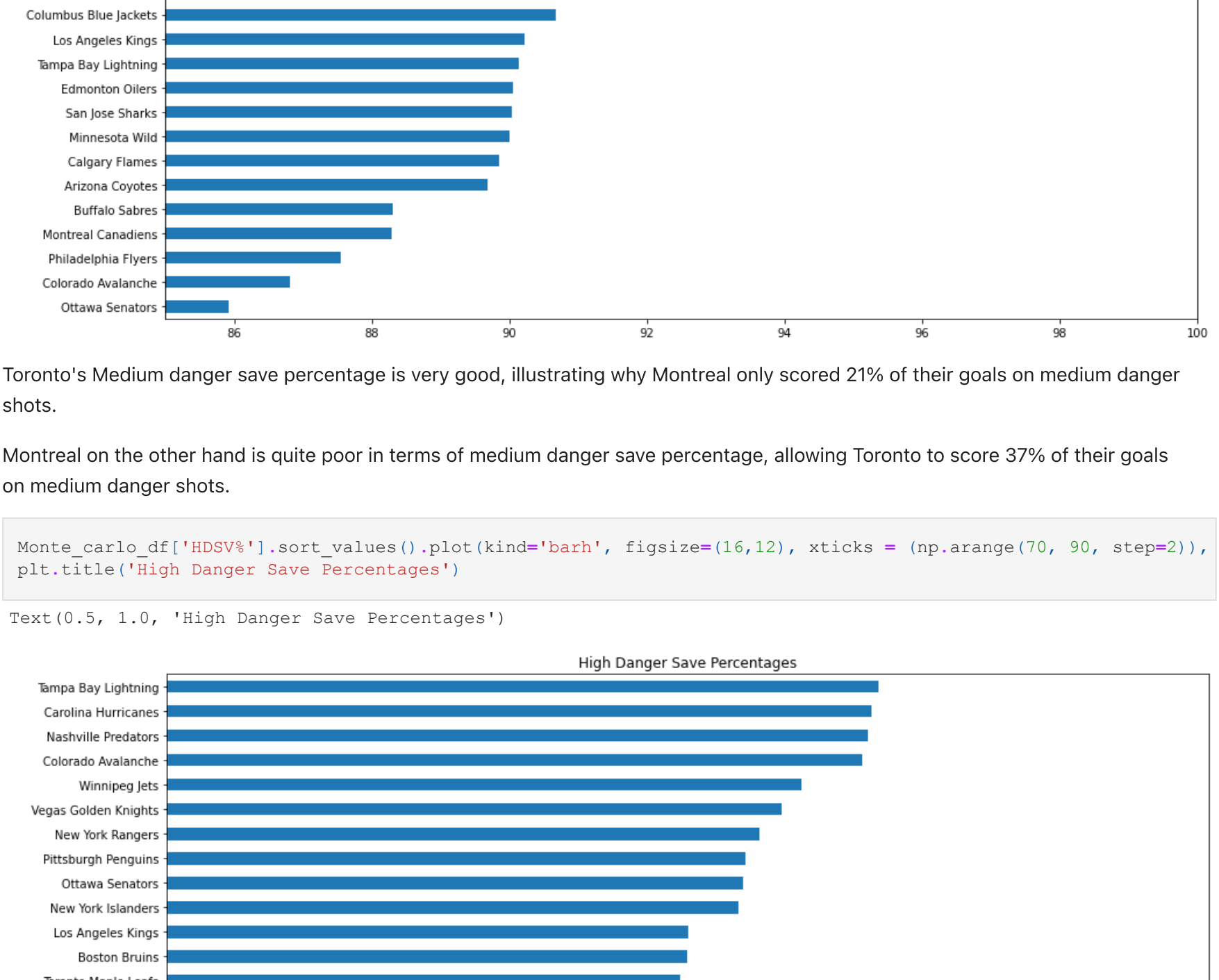
Credit to [hockey-statistics](#) who compiled these on Twitter.

Simulation to Use in Pro Scouting



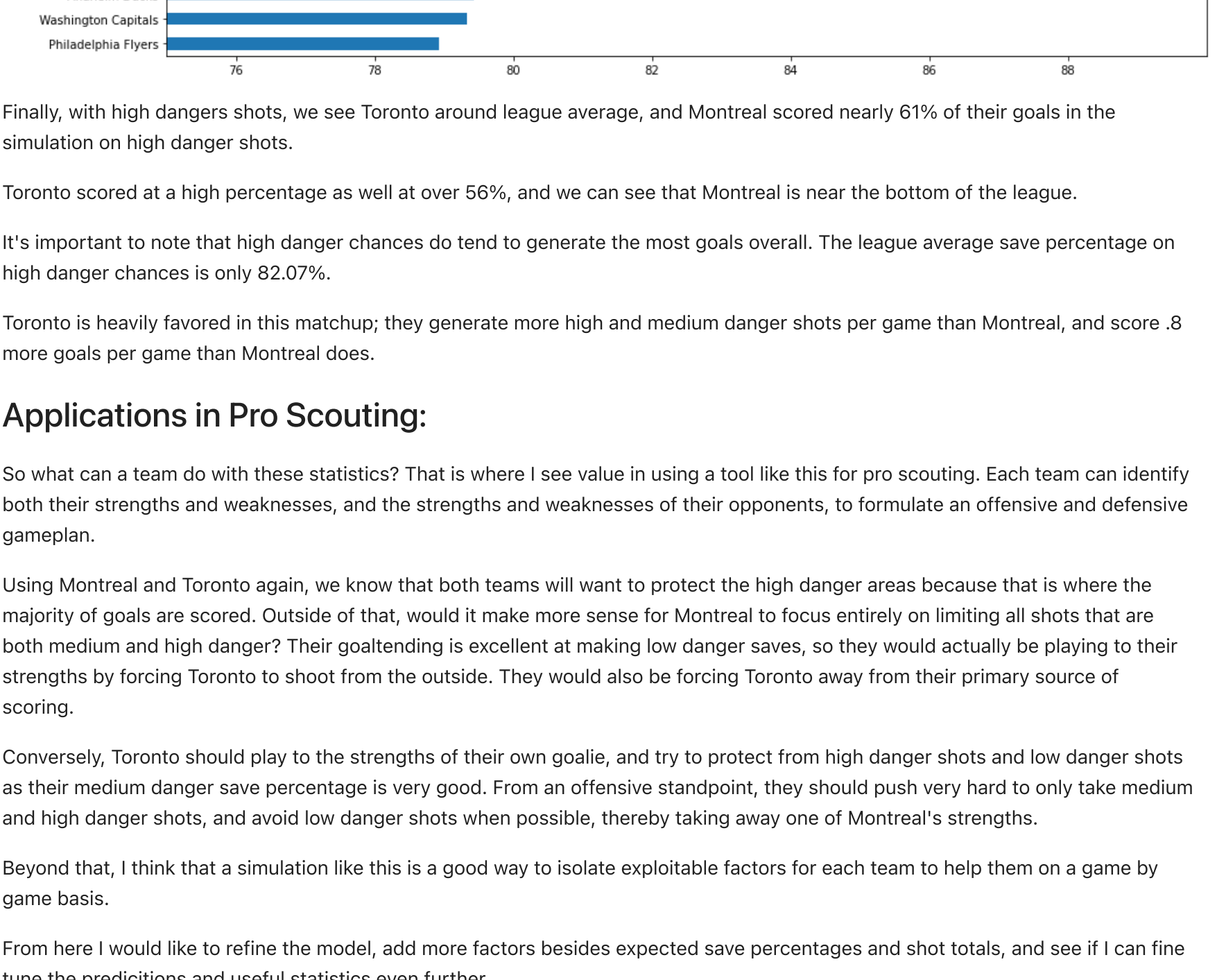
On low danger shots Montreal is one of the best teams in the NHL, which explains why Toronto was only projected to score 6% of their total goals on low danger shots.

Toronto on the other hand is one of the worst teams in the league on low danger shots against, which explains why Montreal scored nearly 18% of their goals on low danger shots in the simulation.



Toronto's Medium danger save percentage is very good, illustrating why Montreal only scored 21% of their goals on medium danger shots.

Montreal on the other hand is quite poor in terms of medium danger save percentage, allowing Toronto to score 37% of their goals on medium danger shots.



Finally, with high dangers shots, we see Toronto around league average, and Montreal scored nearly 61% of their goals in the simulation on high danger shots.

Toronto scored at a high percentage as well at over 56%, and we can see that Montreal is near the bottom of the league.

It's important to note that high danger chances do tend to generate the most goals overall. The league average save percentage on high danger chances is only 82.07%.

Toronto is heavily favored in this matchup; they generate more high and medium danger shots per game than Montreal, and score .8 more goals per game than Montreal does.

Applications in Pro Scouting:

So what can a team do with these statistics? That is where I see value in using a tool like this for pro scouting. Each team can identify both their strengths and weaknesses, and the strengths and weaknesses of their opponents, to formulate an offensive and defensive gameplan.

Using Montreal and Toronto again, we know that both teams will want to protect the high danger areas because that is where the majority of goals are scored. Outside of that, would it make more sense for Montreal to focus entirely on limiting all shots that are both medium and high danger? Their goaltending is excellent at making low danger saves, so they would actually be playing to their strengths by forcing Toronto to shoot from the outside. They would also be forcing Toronto away from their primary source of scoring.

Conversely, Toronto should play to the strengths of their own goalie, and try to protect from high danger shots and low danger shots as their medium danger save percentage is very good. From an offensive standpoint, they should push very hard to only take medium and high danger shots, and avoid low danger shots when possible, thereby taking away one of Montreal's strengths.

Beyond that, I think that a simulation like this is a good way to isolate exploitable factors for each team to help them on a game by game basis.

From here I would like to refine the model, add more factors besides expected save percentages and shot totals, and see if I can fine tune the predictions and useful statistics even further.

Next Steps:

- As I continue to refine this model, I would also like to build a player valuation model. The world of hockey statistics is growing and there will always be new stats added that can hopefully create more accurate evaluations. After building a player valuation model I hope to add the individual player impacts on each game into the Monte Carlo simulations to grow the intricacies even further. Along with that I'd like to add home ice advantage, days of rest, and add weights for each of those. I'd also like to build functions for tracking daily changes in all stat and alterations in lineups and rosters. Finally, I would really like to build a line generator that creates the best 3 forward line combination.
- Conclusions:
 - There is a lot of data out there that can be used to help with the analysis of hockey.
 - Using a basic Monte Carlo simulation created a relatively competitive model compared to others on Twitter.
 - There is a lot more to build to create a highly effective model.
 - Building a detailed function from scratch is a fun and challenging endeavor.
 - Even the most basic hockey statistics can be useful for evaluation of teams.
 - This has solidified my belief that I would enjoy working in sports analytics.
 - I will focus my job search to try to work in this field.

This was an incredibly fun introduction into the world of hockey analytics. I am so excited to dig deeper and learn more every day!