

Applied Statistics and NHL Data

Alex Diaz-Papkovich
Carleton University
alexdiaz.ca
@adp_diaz

Introduction

Introduction

- Universal truths:

Introduction

- Universal truths:
 - Math is difficult

Introduction

- Universal truths:
 - Math is difficult
 - Hockey is fun

Introduction

- Universal truths:
 - Math is difficult
 - Hockey is fun
- My motivation: Using NHL data to study statistics

NHL data as a learning tool

NHL data as a learning tool

- Toyed with it in undergrad

NHL data as a learning tool

- Toyed with it in undergrad
- Same datasets often get used - beetles, flower petals, common public datasets

NHL data as a learning tool

- Toyed with it in undergrad
- Same datasets often get used - beetles, flower petals, common public datasets
- Using a dataset you enjoy enhances what you learn

NHL data as a learning tool

- Toyed with it in undergrad
- Same datasets often get used - beetles, flower petals, common public datasets
- Using a dataset you enjoy enhances what you learn
 - Provides understandable, concrete connections

The process

The process

- Stages of statistical analysis for a grad student?

The process

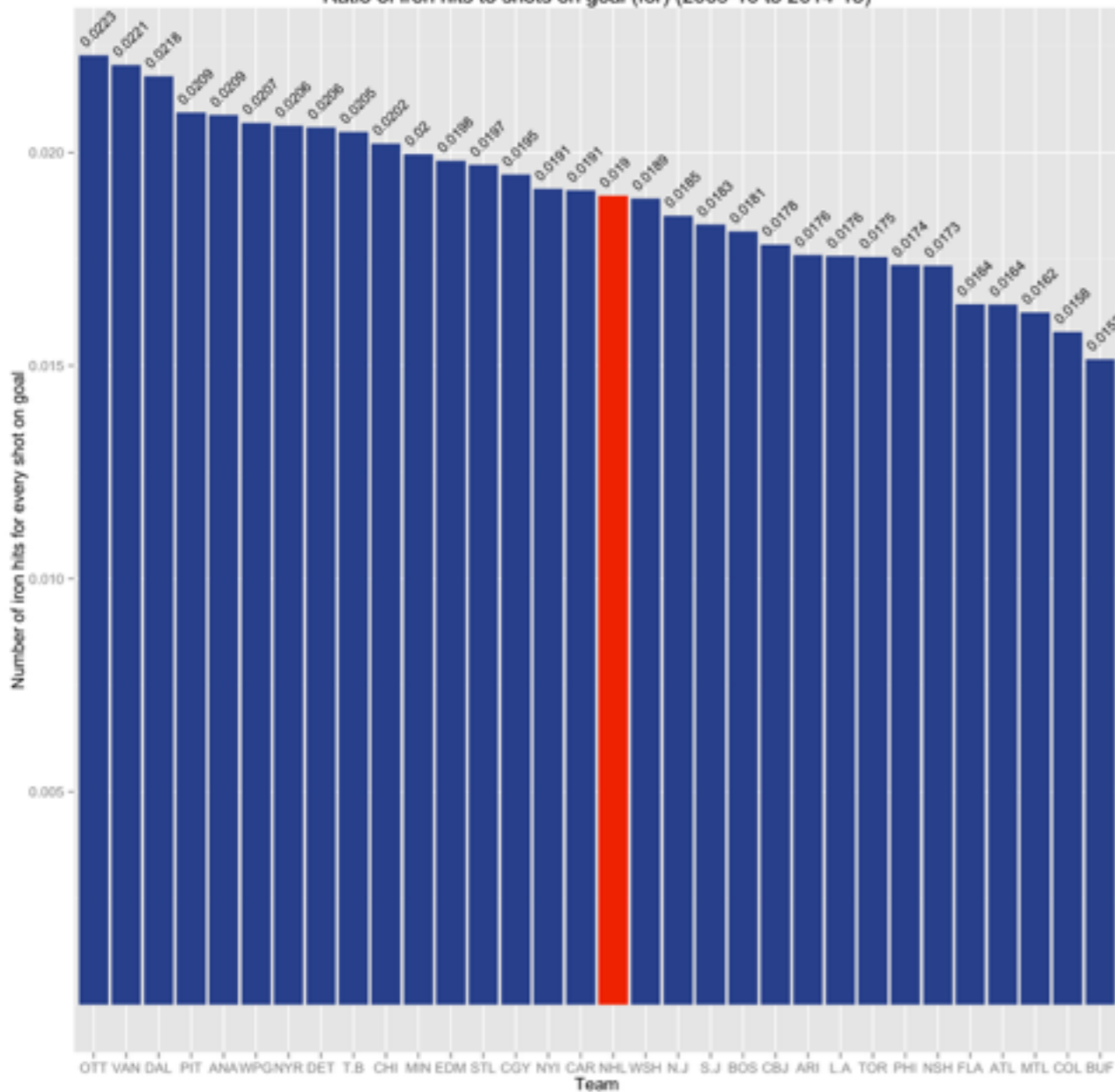
- Stages of statistical analysis for a grad student?
- Denial, anger, bargaining, depression, acceptance?

The process

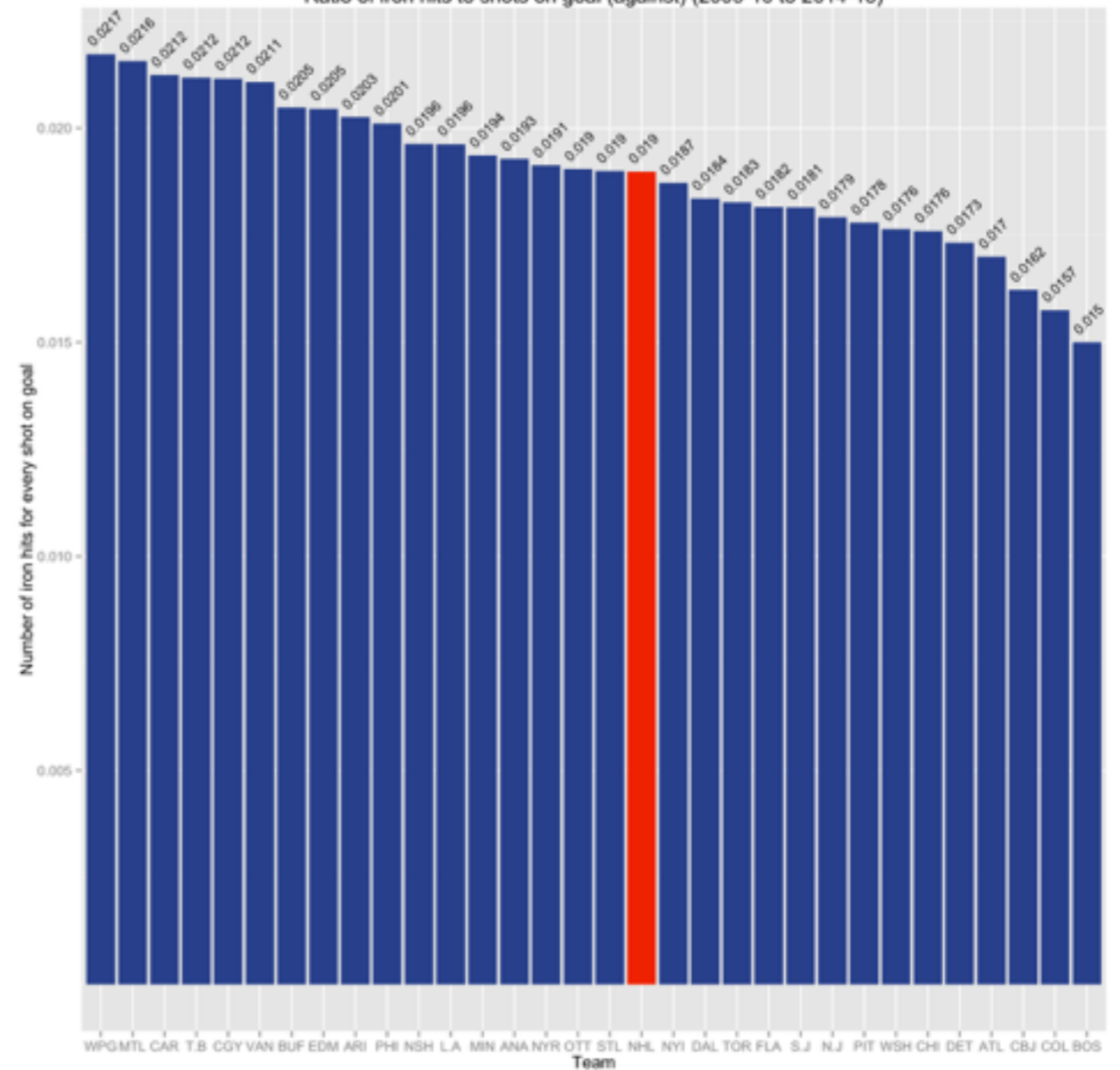
- Stages of statistical analysis for a grad student?
 - Denial, anger, bargaining, depression, acceptance?
 - Broadly: data preparation, exploration, modelling, analysis, evaluation

Who pays the iron price?

Ratio of iron hits to shots on goal (for) (2009-10 to 2014-15)



Ratio of iron hits to shots on goal (against) (2009-10 to 2014-15)





Data preparation

Data preparation

- Often 80%+ of the work

Data preparation

- Often 80%+ of the work
- Stages of data preparation:

Data preparation

- Often 80%+ of the work
- Stages of data preparation:
 - Anger, anger, anger, anger, acceptance

Data preparation

- Often 80%+ of the work
- Stages of data preparation:
 - Anger, anger, anger, anger, acceptance
 - “What do you mean the RTSS doesn’t track which team went offside?”

Data preparation

- Often 80%+ of the work
- Stages of data preparation:
 - Anger, anger, anger, anger, acceptance
 - “What do you mean the RTSS doesn’t track which team went offside?”
 - “There are blank player names??”

Data preparation

- Often 80%+ of the work
- Stages of data preparation:
 - Anger, anger, anger, anger, acceptance
 - “What do you mean the RTSS doesn’t track which team went offside?”
 - “There are blank player names??”
 - “Why does the 2nd period end at -16:0-1???”

Modelling

Modelling

- Make use of common techniques

Modelling

- Make use of common techniques
 - Linear regression

Modelling

- Make use of common techniques
 - Linear regression
 - Logistic regression

Modelling

- Make use of common techniques
 - Linear regression
 - Logistic regression
 - Hypothesis testing

Modelling

- Make use of common techniques
 - Linear regression
 - Logistic regression
 - Hypothesis testing
- Good idea to challenge some assumptions

Modelling

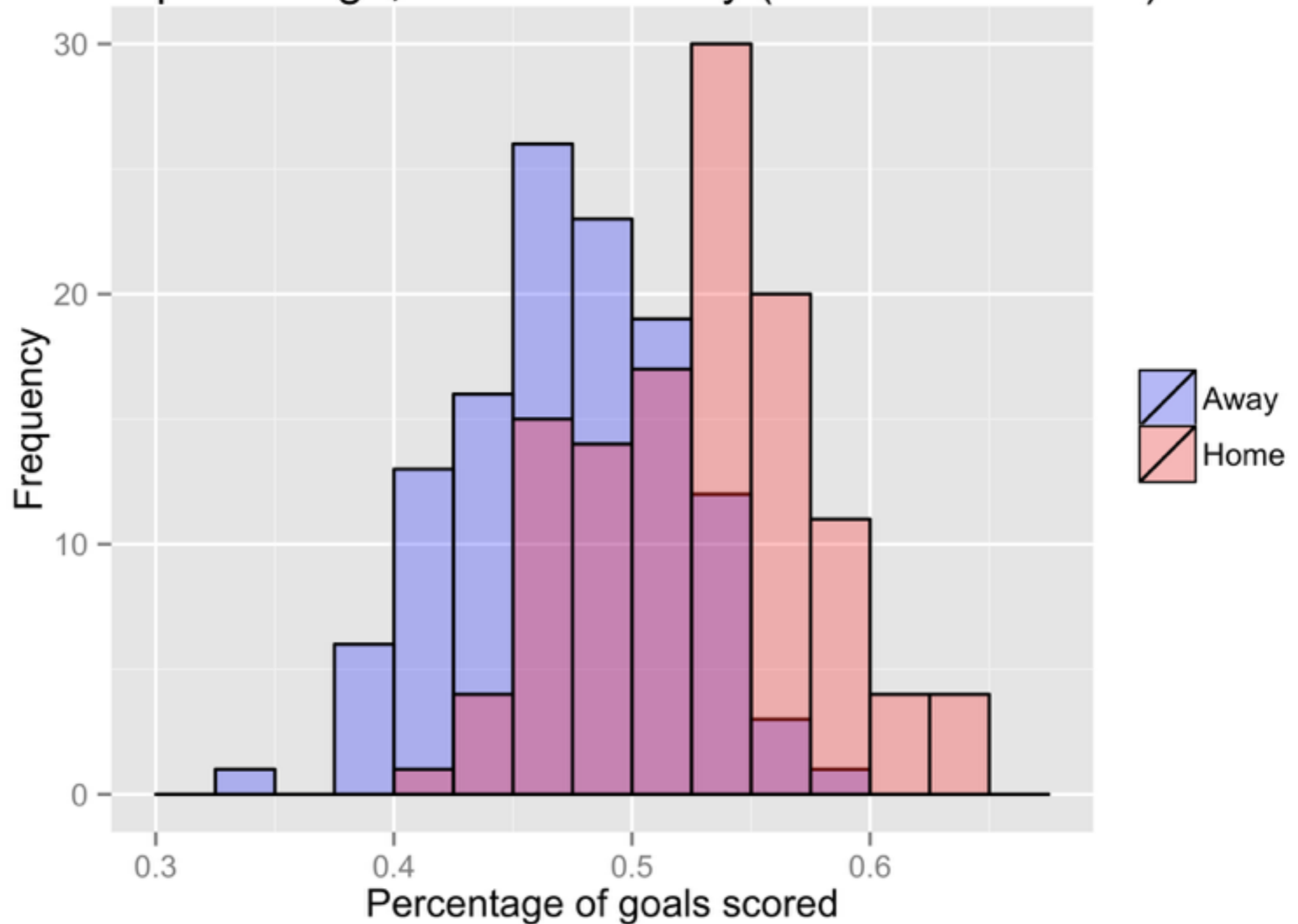
- Make use of common techniques
 - Linear regression
 - Logistic regression
 - Hypothesis testing
- Good idea to challenge some assumptions
 - Is Corsi *actually* a useful measure?

Modelling

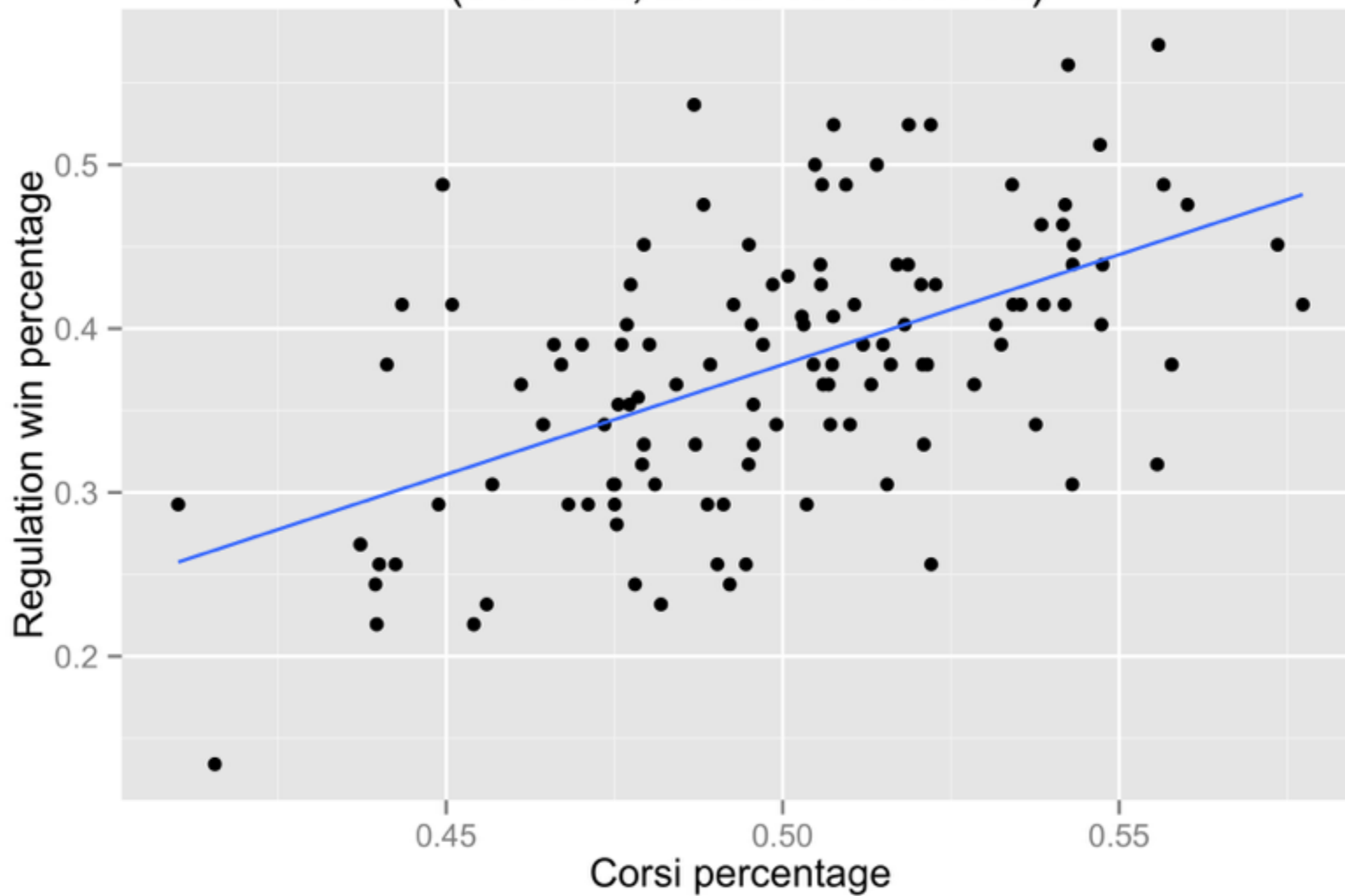
- Make use of common techniques
 - Linear regression
 - Logistic regression
 - Hypothesis testing
- Good idea to challenge some assumptions
 - Is Corsi *actually* a useful measure?
 - Does home ice advantage exist?

Home ice

Goal percentage, home and away (2012-13 excluded)



Regulation win percentage vs Corsi percentage
(5v5 tied, 2012-13 removed)



Comparing shot metrics

	<i>Dependent variable:</i>		
	Win%		
	(1)	(2)	(3)
Shot% (Std. Err.)	1.347*** (0.193)		
Fenwick% (Std. Err.)		1.304*** (0.192)	
Corsi% (Std. Err.)			1.342*** (0.187)
Constant (Std. Err.)	-0.296*** (0.097)	-0.274*** (0.096)	-0.293*** (0.094)
Observations	120	120	120
R ²	0.291	0.282	0.303
Adjusted R ²	0.285	0.276	0.297
Residual Std. Error (df = 118)	0.070	0.070	0.069
F Statistic (df = 1; 118)	48.533***	46.293***	51.372***

Note:

*p<0.1; **p<0.05; ***p<0.01

(1), (2), and (3) are model IDs for each metric

Other factors: PP vs PK

	<i>Dependent variable:</i>			
	Win%			
	(1)	(2)	(3)	(4)
PP goals for (Std. Err.)	0.908*** (0.193)		0.885*** (0.177)	1.978* (1.012)
PK goals against (Std. Err.)		-0.807*** (0.176)	-0.785*** (0.161)	0.232 (0.941)
PP/PK interaction (Std. Err.)				-4.897 (4.463)
Constant (Std. Err.)	0.181*** (0.042)	0.554*** (0.039)	0.358*** (0.053)	0.130 (0.214)
Observations	120	120	120	120
R ²	0.157	0.151	0.300	0.307
Adjusted R ²	0.150	0.144	0.288	0.289
Residual Std. Error	0.076	0.076	0.070	0.070
Residual Std. Error DF	118	118	117	116
F Statistic	22.025***	20.956***	25.084***	17.153***
F Statistic DF	(1;118)	(1;118)	(2;117)	(3; 116)

Note:

*p<0.1; **p<0.05; ***p<0.01

(1), (2), (3), and (4) are model IDs for special teams and interactions

Thinking about theory

Thinking about theory

- What does a shot attempt represent? Why is it useful to measure?

Thinking about theory

- What does a shot attempt represent? Why is it useful to measure?
- Used as a proxy for possession, but what else does it mean?

Thinking about theory

- What does a shot attempt represent? Why is it useful to measure?
- Used as a proxy for possession, but what else does it mean?
 - Control the game, forcing opposition to react

Thinking about theory

- What does a shot attempt represent? Why is it useful to measure?
- Used as a proxy for possession, but what else does it mean?
 - Control the game, forcing opposition to react
- Ultimately, justify decisions in further analysis

Association rule learning

Association rule learning

- Data mining technique that works on large binary databases

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)
- $SUPP(X)$ = % of transactions with itemset X

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)
- $SUPP(X)$ = % of transactions with itemset X
= $Pr(X)$ = Probability X happens

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)
- $SUPP(X)$ = % of transactions with itemset X
= $Pr(X)$ = Probability X happens
- $LIFT(X,Y) = SUPP(X,Y)/[SUPP(X)*SUPP(Y)]$

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)
- $SUPP(X)$ = % of transactions with itemset X
= $Pr(X)$ = Probability X happens
- $LIFT(X,Y) = SUPP(X,Y)/[SUPP(X)*SUPP(Y)]$
= Measure of independence (farther from 1 => less independent)

Association rule learning

- Data mining technique that works on large binary databases
- Previewed at last year's conference
- Think of it as a large-scale With Or Without You (WOWY)
- $SUPP(X)$ = % of transactions with itemset X
= $Pr(X)$ = Probability X happens
- $LIFT(X,Y) = SUPP(X,Y)/[SUPP(X)*SUPP(Y)]$
= Measure of independence (farther from 1 => less independent)
- Events {A,B} are independent if $Pr(A \text{ and } B) = Pr(A) \times Pr(B)$

Example database

COLTON_ORR	JOHN_MICHAEL_LILES	JAMES_VAN_RIEMSDYK	SHOT_FOR	SHOT_AGAINST	MISS_FOR	MISS_AGAINST	BLOCK_FOR	BLOCK_AGAINST
0	0	0	0	0	1	0	0	0
0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	1	0
0	0	1	0	1	0	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0

Example database

COLTON_ORR	JOHN_MICHAEL_LILES	JAMES_VAN_RIEMSDYK	SHOT_FOR	SHOT_AGAINST	MISS_FOR	MISS_AGAINST	BLOCK_FOR	BLOCK_AGAINST
0	0	0	0	0	1	0	0	0
0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	1	0
0	0	1	0	1	0	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0

- Created a database of “transactions”

Example database

COLTON_ORR	JOHN_MICHAEL_LILES	JAMES_VAN_RIEMSDYK	SHOT_FOR	SHOT_AGAINST	MISS_FOR	MISS_AGAINST	BLOCK_FOR	BLOCK_AGAINST
0	0	0	0	0	1	0	0	0
0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	1	0
0	0	1	0	1	0	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0

- Created a database of “transactions”
- Players on the left, events on the right

Example database

COLTON_ORR	JOHN_MICHAEL_LILES	JAMES_VAN_RIEMSDYK	SHOT_FOR	SHOT_AGAINST	MISS_FOR	MISS_AGAINST	BLOCK_FOR	BLOCK_AGAINST
0	0	0	0	0	1	0	0	0
0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	1	0
0	0	1	0	1	0	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0

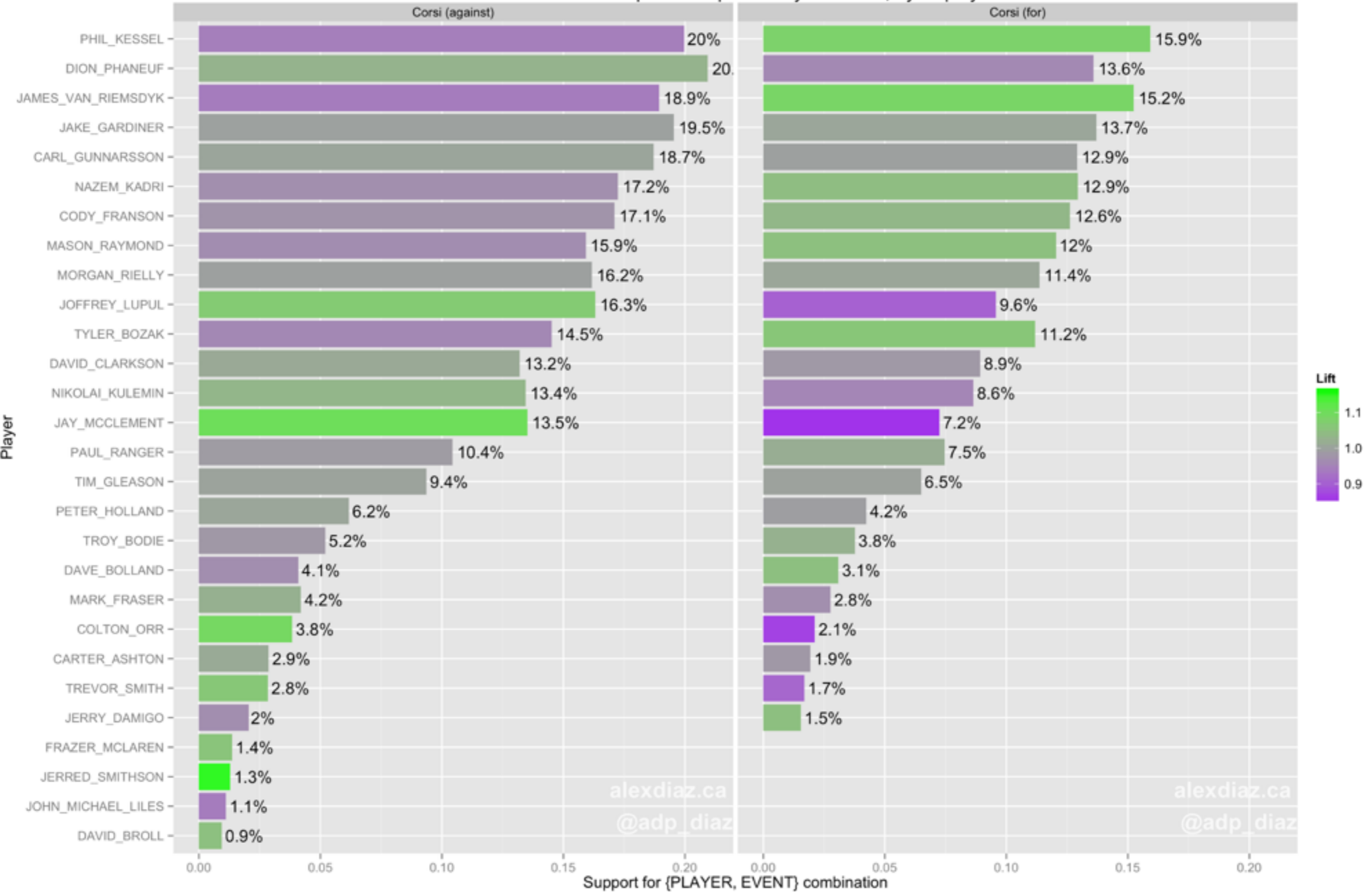
- Created a database of “transactions”
- Players on the left, events on the right
- Looking at Corsi events (EV, 5v5, tied)

Example database

COLTON_ORR	JOHN_MICHAEL_LILES	JAMES_VAN_RIEMSDYK	SHOT_FOR	SHOT_AGAINST	MISS_FOR	MISS_AGAINST	BLOCK_FOR	BLOCK_AGAINST
0	0	0	0	0	1	0	0	0
0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	1	0
0	0	1	0	1	0	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0

- Created a database of “transactions”
- Players on the left, events on the right
- Looking at Corsi events (EV, 5v5, tied)
- Basically try to figure out what players drive the play within a line, team, defense pairing, etc

2013-14 Toronto Maple Leafs probability of events, by all players



alexdiaz.ca
@adp_diaz

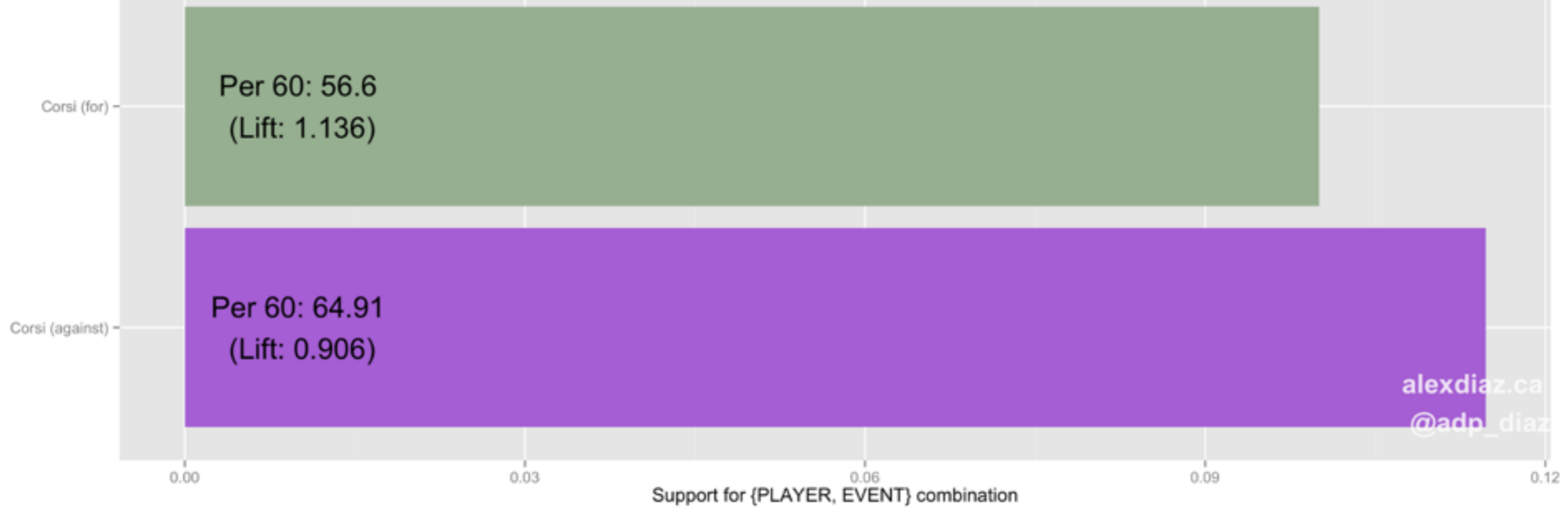
alexdiaz.ca
@adp_diaz

2013-14 Toronto Maple Leafs association rules, by forward lines

PHIL_KESSEL,NAZEM_KADRI,JAMES_VAN_RIEMSDYK



TYLER_BOZAK,PHIL_KESSEL,JAMES_VAN_RIEMSDYK



2014-15 Toronto Maple Leafs association rules, by forward lines

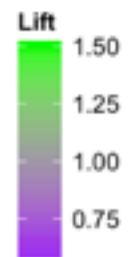
JOFFREY_LUPUL,DANIEL_WINNIK,NAZEM_KADRI



MIKE_SANTORELLI,DANIEL_WINNIK,NAZEM_KADRI



PHIL_KESSEL,TYLER_BOZAK,JAMES_VAN_RIEMSDYK



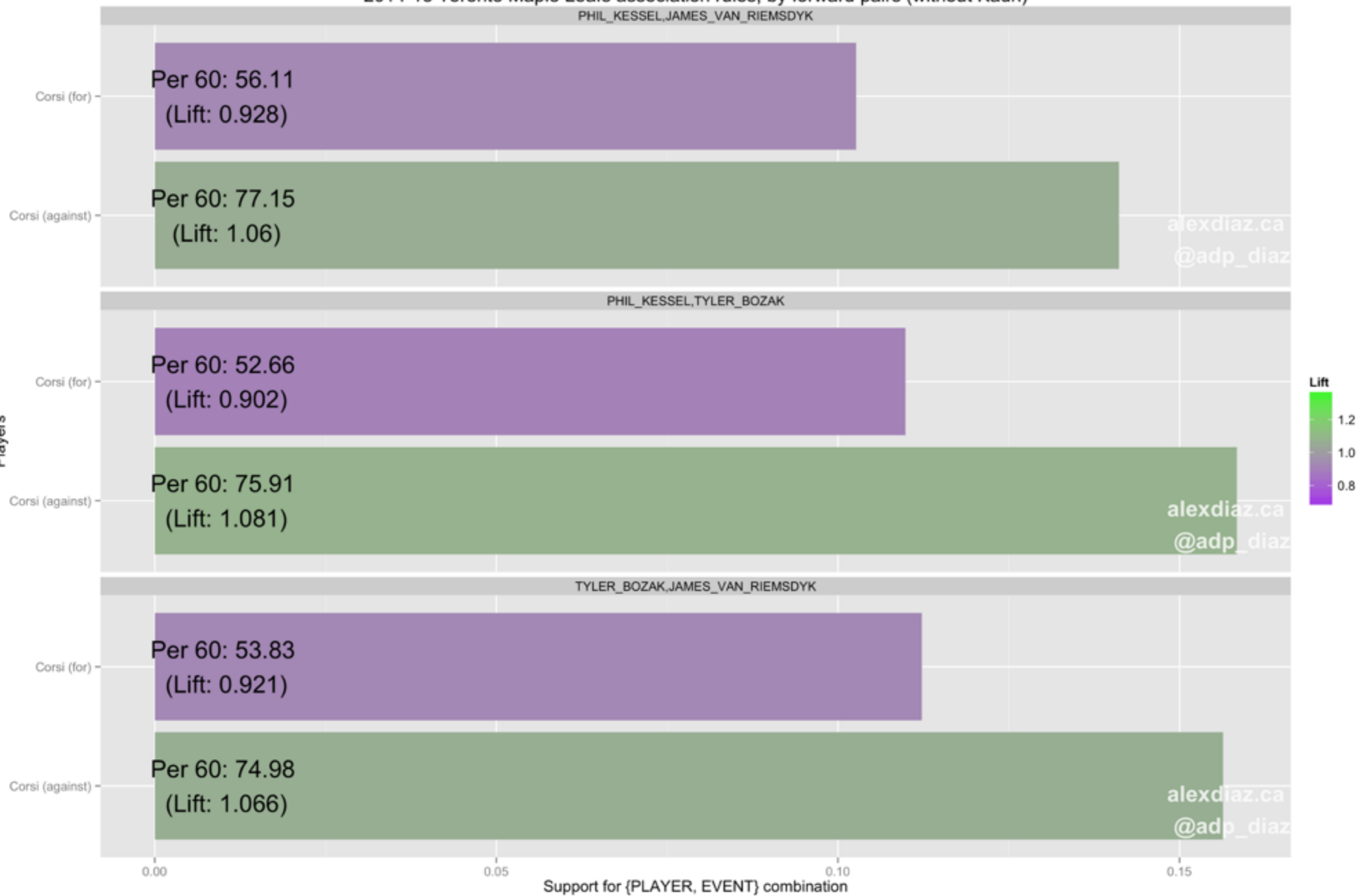
0.00 0.05 0.10
Support for {PLAYER, EVENT} combination

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

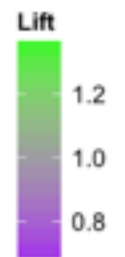
2014-15 Toronto Maple Leafs association rules, by forward pairs (without Kadri)



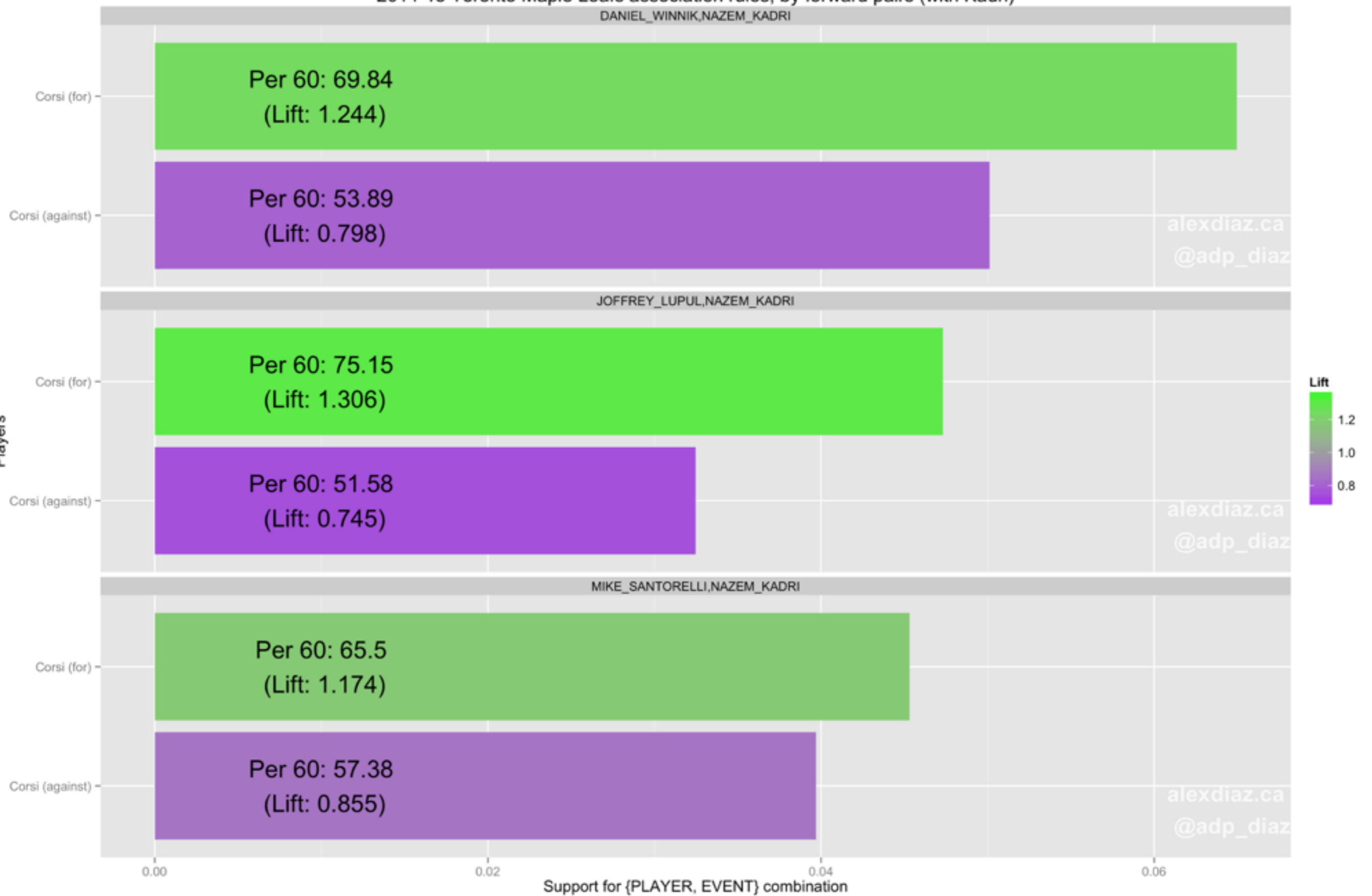
alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz



2014-15 Toronto Maple Leafs association rules, by forward pairs (with Kadri)



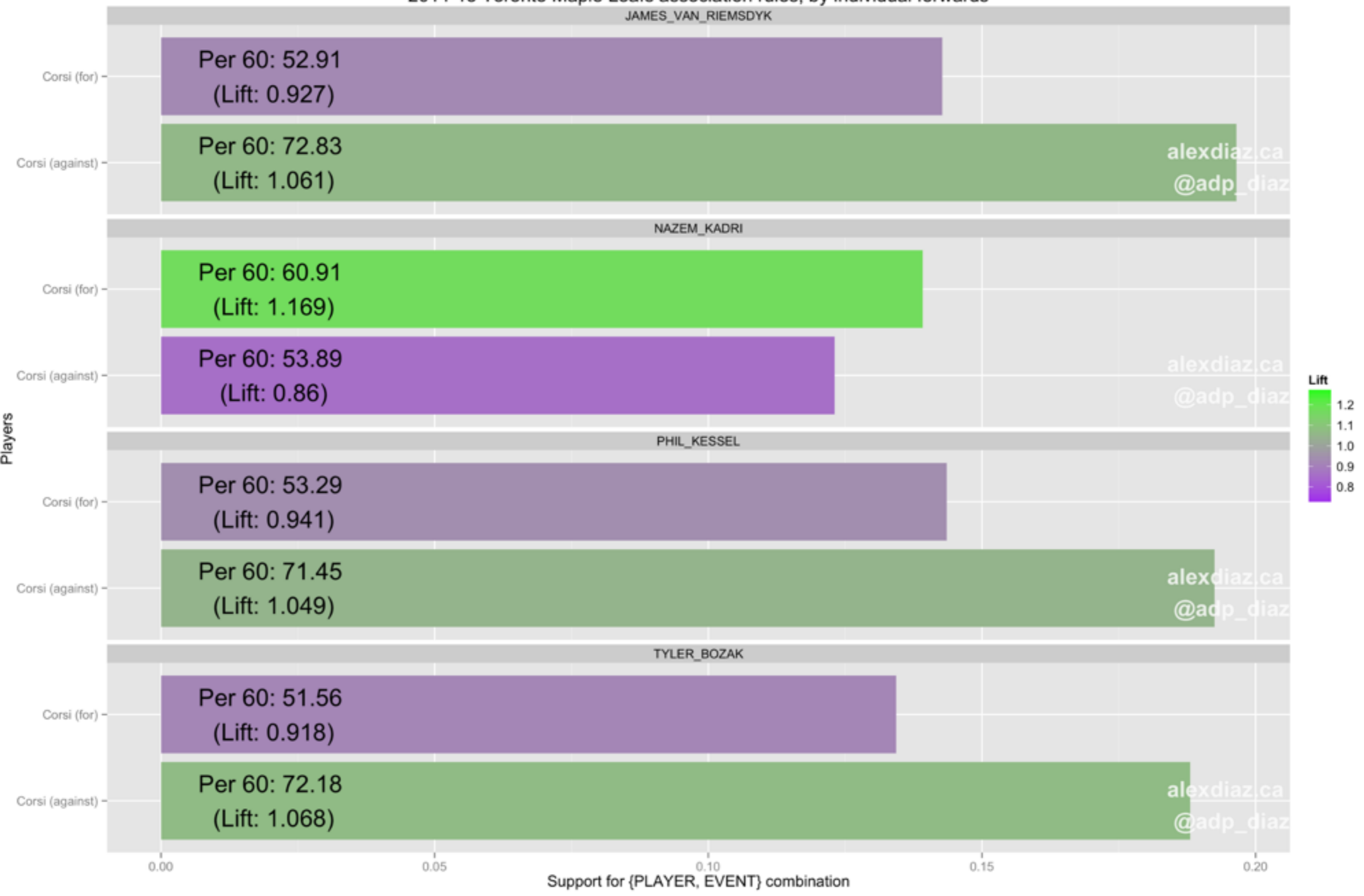
alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz



2014-15 Toronto Maple Leafs association rules, by individual forwards



2015-16 Chicago Blackhawks association rules, by forwards lines

ANDREW_SHAW,MARCUS_KRUGER,ANDREW_DESJARDINS

Corsi (for)
Per 60: 45.22
(Lift: 0.952)

Corsi (against)
Per 60: 46.59
(Lift: 1.052)

alexdiaz.ca
@adp_diaz

ANDREW_SHAW,MARIAN_HOSSA,JONATHAN_TOEWS

Corsi (for)
Per 60: 57.21
(Lift: 1.13)

Corsi (against)
Per 60: 40.57
(Lift: 0.86)

alexdiaz.ca
@adp_diaz

ARTEMI_PANARIN,PATRICK_KANE,ARTEM_ANISIMOV

Corsi (for)
Per 60: 60.75
(Lift: 1.116)

Corsi (against)
Per 60: 44.37
(Lift: 0.875)

alexdiaz.ca
@adp_diaz

TEUVO_TERAVAINEN,MARIAN_HOSSA,JONATHAN_TOEWS

Corsi (for)
Per 60: 61.39
(Lift: 0.996)

Corsi (against)
Per 60: 57.71
(Lift: 1.004)

alexdiaz.ca
@adp_diaz

0.00

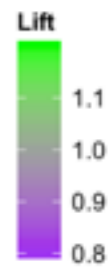
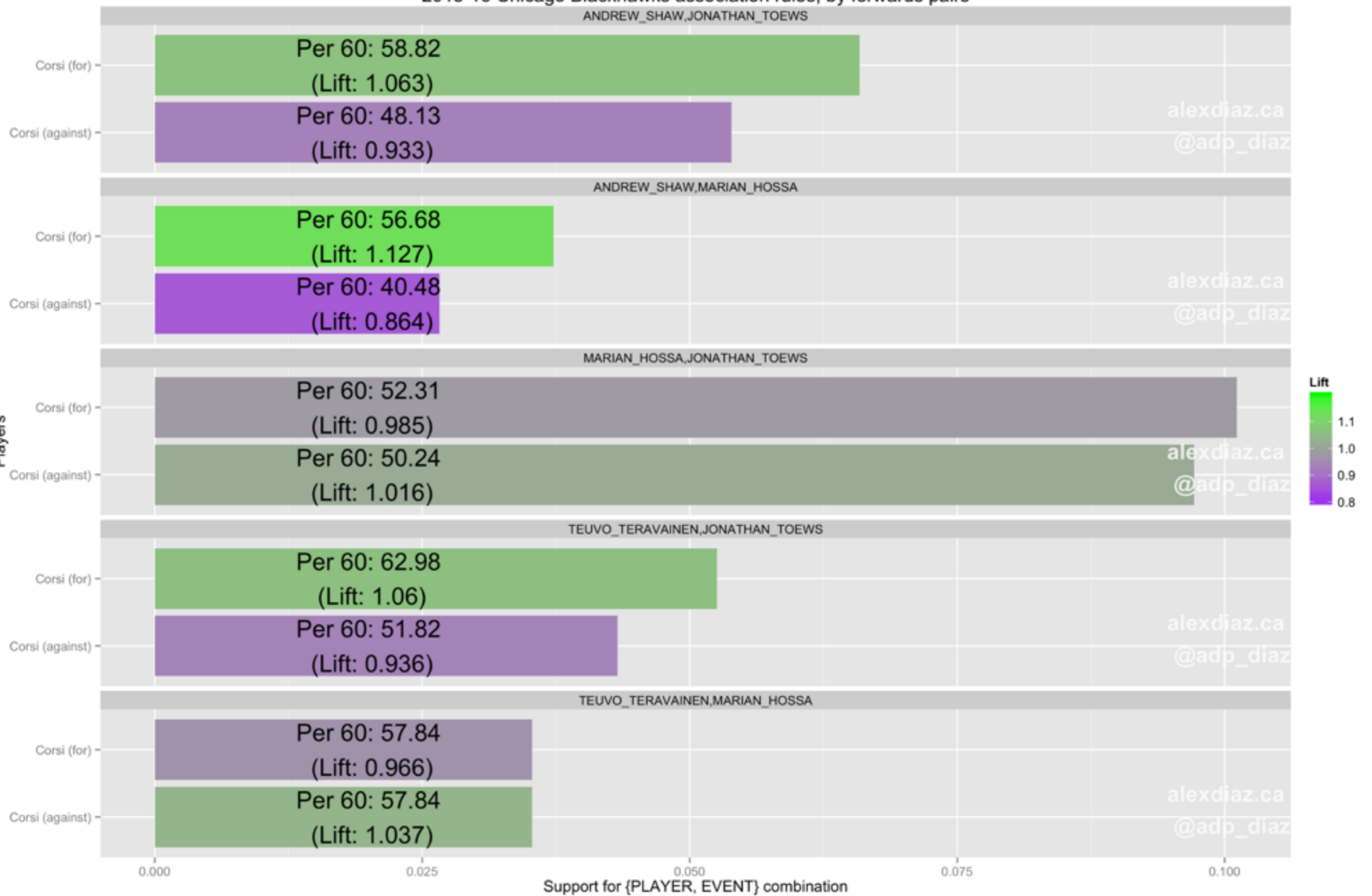
0.05

0.10

Support for {PLAYER, EVENT} combination



2015-16 Chicago Blackhawks association rules, by forwards pairs



alexdiaz.ca
@ado_diaz

alexdiaz.ca
@ado_diaz

alexdiaz.ca
@ado_diaz

alexdiaz.ca
@ado_diaz

alexdiaz.ca
@ado_diaz

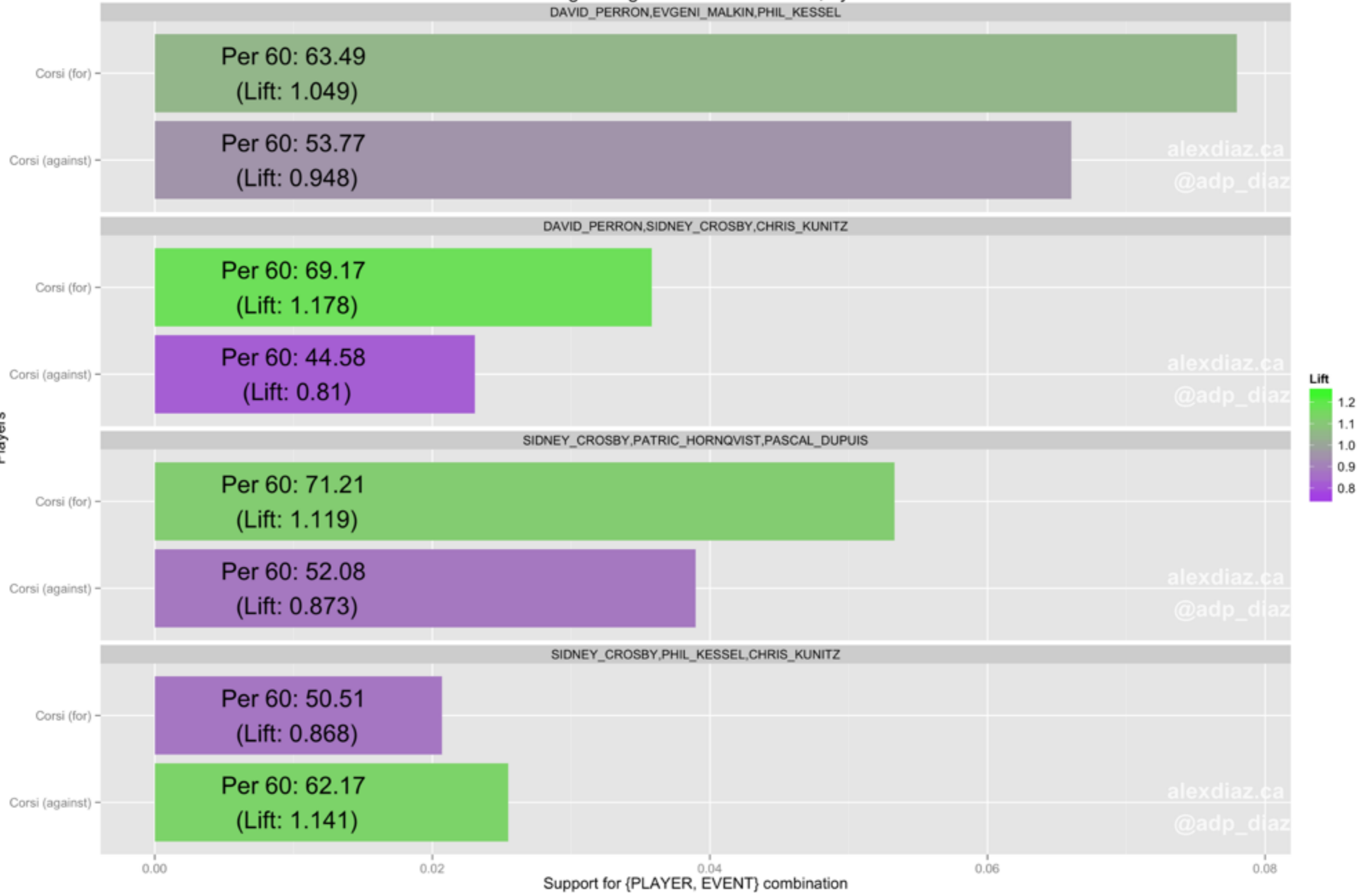
Players

Corsi (for)
Corsi (against)
Corsi (for)
Corsi (against)
Corsi (for)
Corsi (against)
Corsi (for)
Corsi (against)

0.000 0.025 0.050 0.075 0.100

Support for {PLAYER, EVENT} combination

2015-16 Pittsburgh Penguins association rules, by forward lines

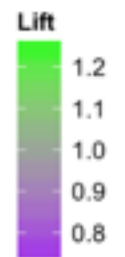


alexdiaz.ca
@adp_diaz

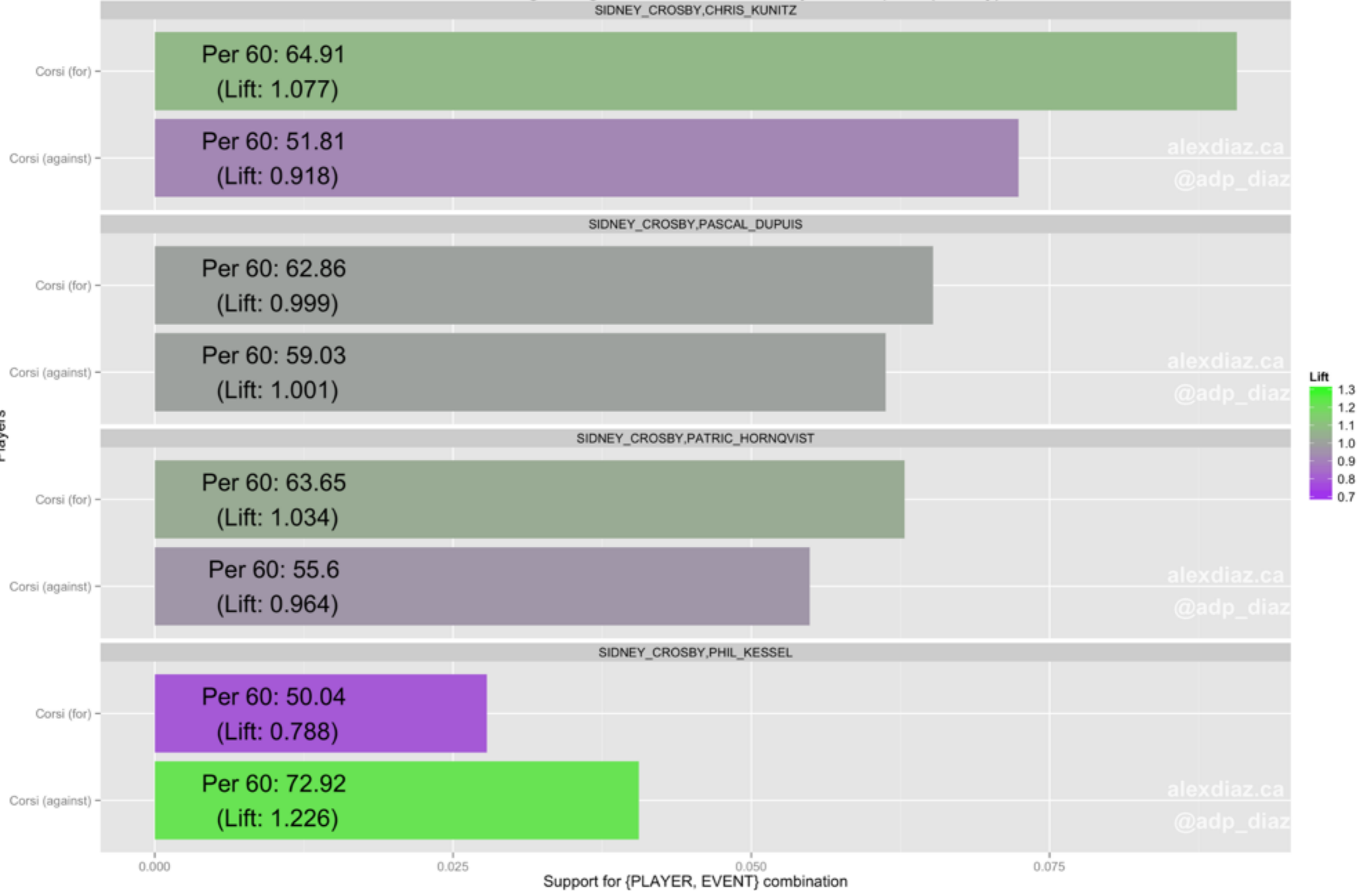
alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz



2015-16 Pittsburgh Penguins association rules, by forward pairs (Crosby)

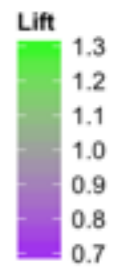


alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz

alexdiaz.ca
@adp_diaz



Evaluation

Evaluation

- Doesn't account for quality of opposition or zone starts

Evaluation

- Doesn't account for quality of opposition or zone starts
- Doesn't directly use time-on-ice

Evaluation

- Doesn't account for quality of opposition or zone starts
- Doesn't directly use time-on-ice
- No clear "thresholds"

Evaluation

- Doesn't account for quality of opposition or zone starts
- Doesn't directly use time-on-ice
- No clear "thresholds"
- Needs variance estimation, especially for lower sample sizes

Evaluation

- Doesn't account for quality of opposition or zone starts
- Doesn't directly use time-on-ice
- No clear "thresholds"
 - Needs variance estimation, especially for lower sample sizes
- Still useful and informative

Conclusions

Conclusions

- Useful for measuring player contributions within a team

Conclusions

- Useful for measuring player contributions within a team
- Potential uses for coaching, management

Conclusions

- Useful for measuring player contributions within a team
- Potential uses for coaching, management
 - Measure player performance within given context

Conclusions

- Useful for measuring player contributions within a team
- Potential uses for coaching, management
 - Measure player performance within given context
 - Create more effective lines

Conclusions

- Useful for measuring player contributions within a team
- Potential uses for coaching, management
 - Measure player performance within given context
 - Create more effective lines
 - Make decisions between players with similar roles

Questions/Comments/ Concerns?

- alexdiaz.ca
- @adp_diaz