# AFL – HOW TO WIN

A data driven analysis of winning factors in matches

James Hooi

# SPORT IS ABOUT WINNING

**Choosing personnel and tactics in sport used to be more art than science. In some sports, it arguably still is. Nowadays, the rich commercial returns of succeeding in sport through TV rights deals, sponsorships and merchandising is too much to leave to chance – teams need to maximise their chances of "making their own luck".**

The story of Moneyball and the Oakland A's signalled the beginnings of a data-driven revolution in sport

Elite sport is about winning and commercial success

- Commercial success often follows on-field success
- Data analysis is giving teams the competitive edge in elite competition

This analysis sets out to find a link between AFL player performance and team results

Player metrics **+** **Some formula** (explains) **=** Team results

The analysis uses free data from **afltables.com** and comprises:

**EXPLORATORY ANALYSIS**

Gathering player metric and team results data, creating team-wise totals for each match

Performing exploratory visualisation to look for relationships in the raw data

**INDEX STATISTICAL ANALYSIS**

Creating indices to measure team vs opponent, performing visual and statistical analysis

**MACHINE LEARNING**

Creating additional proportion measures and clustering data for modelling

Creating Decision Tree models, testing predictive accuracy non-clustered and clustered

Creating Logistic Regression models, testing predictive accuracy non-clustered and clustered

**CONCLUSIONS**

Documenting findings, opportunities for further analysis and wrap up

**Data visualisation provides a simple and digestible method to get an overview of the data.**

Below are plots showing team totals of player metrics in wins vs losses

Illustrated by year for last 15 years

Annotated with statistical confidence interval of win / loss difference being positive – test of whether metric is significant to winning



Behinds — 65.9% confident win/loss difference positive

Running Bounces — 57.1% confident win/loss difference positive

Clangers — 58.7% confident win/loss difference negative

Clearances — 59.9% confident win/loss difference positive

Contested marks — 62.2% confident win/loss difference positive

Contested possessions — 62.9% confident win/loss difference positive

Disposals — 69.9% confident win/loss difference positive

Free kicks awarded — 52.4% confident win/loss difference positive

Free kicks given away — 52.0% confident win/loss difference negative

Goal assists — 81.3% confident win/loss difference positive

Goals — 87.1% confident win/loss difference positive

Handballs — 57.5% confident win/loss difference positive

Hit-outs — 53.6% confident win/loss difference positive

Inside 50s — 77.9% confident win/loss difference positive

Kicks — 79.1% confident win/loss difference positive

Marks (catches) — 68.1% confident win/loss difference positive

Marks inside 50 — 76.1% confident win/loss difference positive

1 percenters — 53.2% confident win/loss difference positive

Tackles — 53.6% confident win/loss difference positive

Uncontested possessions — 65.9% confident win/loss difference positive

While some relationships are evident, they aren't useful

The raw data tells us to score, by possessing the ball, from winning it in contests (obvious)

Deeper analysis may uncover more useful insights

**In sports analytics, it's critical to remember the game is adversarial – opponents will also influence the result.**

Players are not in total control of their performance outcomes due to opponents' actions

Indices can be used to compare a team versus its opponent in key outcome measures (log used to standardise results)

### Win contested ground ball
idx_win_ground_ball = log

$\dfrac{\text{contested possessions}}{\text{opposition contested possessions}}$

### Win contested aerial ball
idx_win_aerial_ball = log

$\dfrac{\text{contested marks}}{\text{opposition contested marks}}$

### Assisted goal
idx_goal_assist = log

$\dfrac{\text{goal assists/goals}}{\text{opposition goal assists/opp goals}}$

### Maintain possession
idx_mark_kick = log

$\dfrac{\text{marks/kicks}}{\text{opposition marks/opp kicks}}$

### Accurate 50m entry
idx_50m_entry = log

$\dfrac{\text{marks inside 50/inside 50s}}{\text{opposition marks i50/opp i50s}}$

### Less clangers
idx_less_clangers = log

$\dfrac{\text{opposition clangers}}{\text{clangers}}$

### Tackling
idx_tackle = log

$\dfrac{\text{tackles}}{\text{opposition tackles}}$

### Clear ball to advantage
idx_clear_ball = log

$\dfrac{\text{clearances}}{\text{opposition clearanced}}$

### One percenters
idx_one_pct = log

$\dfrac{\text{one percenters}}{\text{opposition one percenters}}$

### Give away less frees
idx_less_frees = log

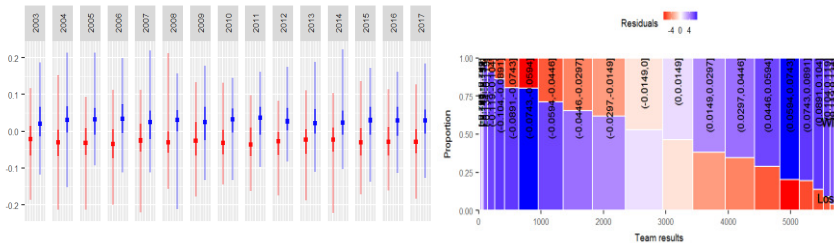$\dfrac{\text{opposition frees against}}{\text{frees against}}$

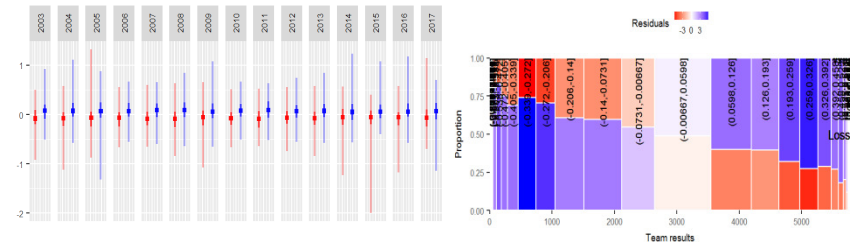Below are plots showing team indices in wins vs losses

Next to the line plots are mosaics coloured by Chi-squared residuals: blue for over-represented, red for under-represented

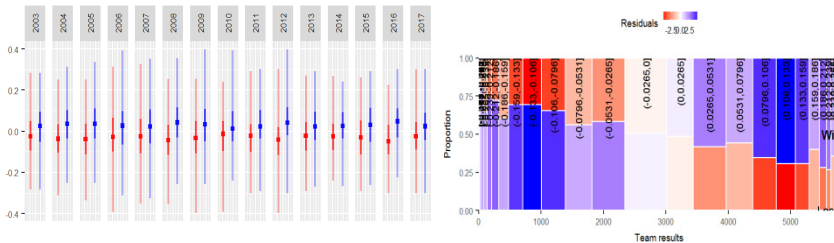Stronger slopes and shades indicate a more significant relationship to wins
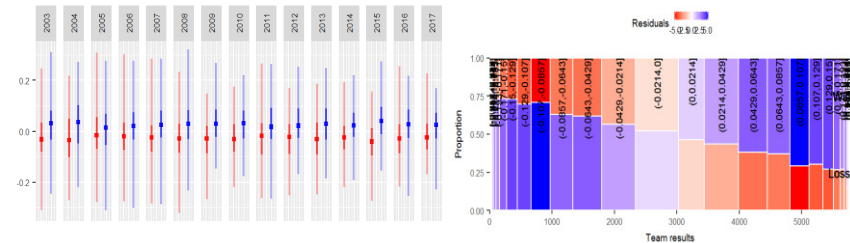


idx_win_ground_ball: Win ground ball

idx_win_aerial_ball: Win aerial ball

idx_clear_ball: Win first use

idx_less_clangers: Make less errors

idx_goal_assist: Teamwork for goal

idx_mark_kick: Maintain possession

idx_50m_entry: Accurate 50m entry

idx_tackle: Win ball back (tackles)

idx_one_pct: Exceptional effort (1 percenters)

idx_less_frees: Discipline

8

Distribution histograms of index values by win vs loss help confirm observations from index and mosaic plots

The more distinctly the win curve (blue) differs from the loss curve (red), the more confidently we can visually confirm the relationship of index to wins

**Quantifying the results statistically with a significance test helps give them more meaning.**

The null hypothesis is an index has no bearing on a win or loss result, so the true mean of the index is zero

If sport is random, a team should win 50% of the time

The calculated z-score from the population mean and standard deviation corresponds to probability of a win if the team scores higher on the index than its opponent
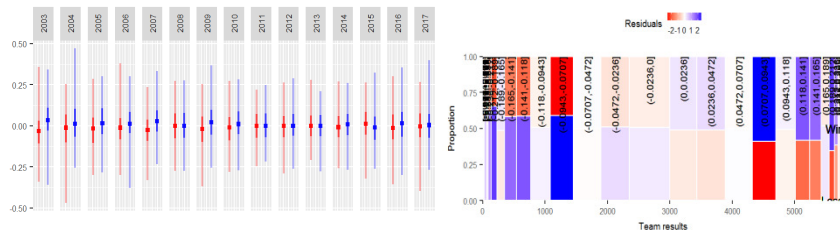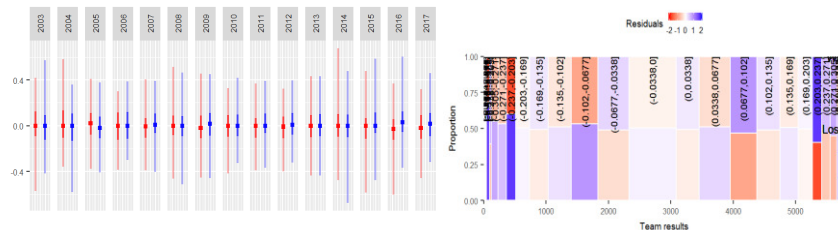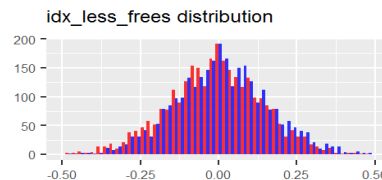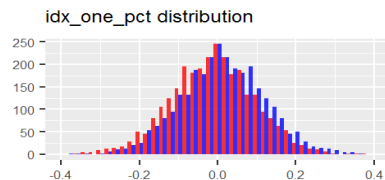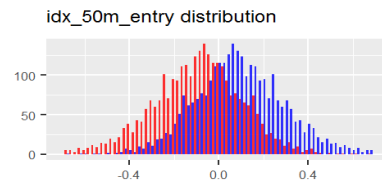
| Index | Z-score | Significance test result interpretation |
|---|---|---|
| idx_win_ground_ball | 0.52 | The team leading this index wins approximately **70%** of the time **(Strong influence)** |
| idx_win_aerial_ball | 0.33 | The team leading this index wins approximately **63%** of the time **(Fair influence)** |
| idx_clear_ball | 0.27 | The team leading this index wins approximately **61%** of the time **(Fair influence)** |
| idx_less_clangers | 0.30 | The team leading this index wins approximately **62%** of the time **(Fair influence)** |
| idx_goal_assist | 0.12 | The team leading this index wins approximately **55%** of the time **(Weak influence)** |
| idx_mark_kick | 0.19 | The team leading this index wins approximately **58%** of the time **(Weak influence)** |
| idx_50m_entry | 0.45 | The team leading this index wins approximately **67%** of the time **(Fair influence)** |
| idx_tackle | 0.16 | The team leading this index wins approximately **56%** of the time **(Weak influence)** |
| idx_one_pct | 0.10 | The team leading this index wins approximately **54%** of the time **(Weak influence)** |
| idx_less_frees | 0.06 | The team leading this index wins approximately **52%** of the time **(Weak influence)** |

The exploratory and statistical analyses show relationships between indices derived from player metrics and team results. The next step is to model these relationships seeking to predict match results from player performance, and more importantly interpreting the models to gain insights. When modelling, the data used to "train" the model must be quarantined from the data used to "test" the model's accuracy.

## Calculated proportion of top ten players belonging to each team for key stats

- E.g. Of players with top ten most inside 50 entries, team A had six (60%)
- Top ten expanded for ties, e.g. if players ranked 10 and 11 for a stat were the same, the proportion was of the top 11

## Created six data clusters, thought of as ladder position (H / M / L) and game plan (fast / slow)

- Future analyses could explore different cluster stratification

*Ball movement*

|  | Fast | Slow |
|---|---|---|
| High | | |
| Med | | |
| Low | | |

*Ladder finish*

## Tested for variable impact on each other (multicollinearity), which affects model validity

- Correlations below indicate some care in modelling is needed



Pearson Correlation scale: 1.0, 0.5, 0.0, -0.5, -1.0

**Before detailing the models, a quick point about how the models will be measured:**

## Data up to 2016 used to train the model and 2017 matches used to test it

- This is the usual set up for testing predictive accuracy of machine learning models
- The models will aim to classify whether a win was achieved or not, therefore accuracy is measured by how many wins and non-wins were correctly predicted

## Realistically, 2017 player performance data would not be available so true predictive accuracy will be tested

- 2017 player performance data will itself be predicted by taking the players' 2016 match averages
- These averages will be used to aggregate team totals and indices, then the models applied to predict the match result
- The models will again be measured on classification accuracy

Training the model: 2003 – 2016 matches

| Player performance data | Team results |

Testing the model: 2017 matches

Measures how well the player performance metrics **explain** the match result

| Player performance data | Team results |

Testing the model: 2017 matches

Measures how well the models actually **predict** match results when 2017 metrics are redacted

| 2016 average player performance data | Team results |

**A Decision Tree model iterates through each index to "split" by the most deterministic variable at each point. The tree below is "tuned" to balance over-fitting the model to the training data, predictive accuracy, and interpretability.**

## Each split checks if a condition is met

- LHS for Yes
- RHS for No

## Each "leaf" at the bottom shows the predicted match result and how many losses vs wins were in the training data

- 0 for predicted loss, 1 for predicted win
- LHS: number of actual losses in leaf
  RHS: number of actual wins in leaf

The Decision Tree explained results well in testing (**316** of **414** matches classified correctly, 2017 player performance available)

However, prediction was less accurate (**267** of **414** matches classified correctly, 2017 player performance estimated)



ROC Curve for Decision Tree - Tuned
AUC = 0.83
Class accuracy = 77%



ROC Curve for Decision Tree - Tuned
AUC = 0.63
Class accuracy = 64%

### Decision Tree diagram

Inside_50s < 0.49 — yes / no

- idx_50m_entry < 0.045
- idx_50m_entry < -0.026
- idx_win_ground_ball < -0.011
- idx_win_ground_ball < 0.017
- idx_win_ground_ball < -0.016
- idx_less_clangers < -0.011
- Inside_50s < 0.54
- Inside_50s < 0.42

Leaves:

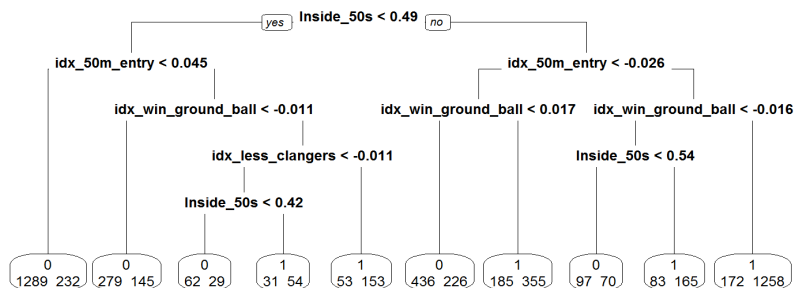| Node | Losses | Wins |
|------|--------|------|
| 0 | 1289 | 232 |
| 0 | 279 | 145 |
| 0 | 62 | 29 |
| 1 | 31 | 54 |
| 1 | 53 | 153 |
| 0 | 436 | 226 |
| 1 | 185 | 355 |
| 0 | 97 | 70 |
| 1 | 83 | 165 |
| 1 | 172 | 1258 |

*Classification accuracy = proportion of correct predictions against total matches. The benchmark to beat is flipping a coin (50% accuracy).*

*AUC = Area Under the Curve, a measure of how often a win is predicted for an actual win versus a win is predicted for an actual loss. The benchmark to beat is random (AUC = 0.50).*

Using one model per cluster, prediction marginally improved

| Cluster | Games | Classification accuracy | AUC |
|---------|-------|-------------------------|------|
| 1 | 11 | 100% | 1.00 |
| 2 | 17 | 65% | 0.50 |
| 3 | 43 | 67% | 0.65 |
| 4 | 133 | 65% | 0.65 |
| 5 | 146 | 64% | 0.63 |
| 6 | 64 | 66% | 0.66 |

Weighted avg accuracy  65%
Weighted avg AUC        0.65

**A Logistic Regression model seeks a linear relationship between indices and winning. The model below ignores non-significant indices to improve predictive accuracy, and checks for multicollinearity to avoid over-fitting.**

## Each variable has a coefficient ("estimate" below) which is a measure of influence on winning

- Note: estimate units are non-standardised so cannot be interpreted verbatim (see next page for variable significance)

```
##           term      estimate  std.error    statistic       p.value
## 1   (Intercept)   -3.7024821  0.1930328   -19.180587  5.376952e-82
## 2 idx_win_aerial_ball  0.6103372  0.1742489     3.502674  4.606129e-04
## 3   idx_less_frees   -4.3877527  0.3366227   -13.034632  7.774111e-39
## 4 idx_win_ground_ball 15.6620305  0.8361181    18.731841  2.723356e-78
## 5  idx_less_clangers   9.1791529  0.5693165    16.123111  1.755452e-58
## 6    idx_50m_entry    4.2322507  0.2110606    20.052300  1.927247e-89
## 7      Inside_50s     7.4055732  0.3788730    19.546319  4.433787e-85
```

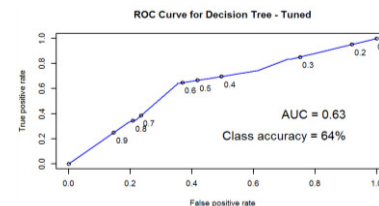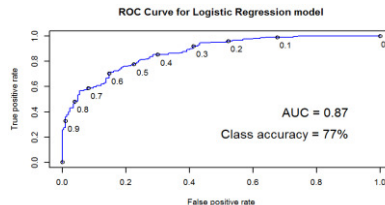Very small p values indicate the variable is significant to prediction of win/loss result

**Classification accuracy** = proportion of correct predictions against total matches. The benchmark to beat is flipping a coin (50% accuracy).

**AUC** = Area Under the Curve, a measure of how often a win is predicted for an actual win versus a win is predicted for an actual loss. The benchmark to beat is random (AUC = 0.50).

**Note:** * Further tuning is available by choosing threshold value (currently probability > 50%)

The Logistic Regression model explained results well in testing (**316** of **414** matches classified correctly, 2017 player performance available) – same accuracy as Decision Tree but better AUC

ROC Curve for Logistic Regression model

AUC = 0.87
Class accuracy = 77%

Prediction was less accurate (**259** of **414** matches classified correctly, 2017 player performance estimated) – less accurate than Decision Tree but better AUC

ROC Curve for Logistic Regression model

AUC = 0.65
Class accuracy = 63%

Using one model per cluster, prediction was unchanged*

| Cluster | Games | Classification accuracy | AUC |
|---------|-------|-------------------------|-----|
| 1 | 11 | 91% | 0.94 |
| 2 | 17 | 65% | 0.76 |
| 3 | 43 | 70% | 0.60 |
| 4 | 133 | 60% | 0.62 |
| 5 | 146 | 60% | 0.57 |
| 6 | 64 | 69% | 0.75 |

Weighted avg accuracy 63%
Weighted avg AUC     0.63

The relative significance of variables in the Logistic Regression model can be tested by leaving the variable out of the model and measuring the amount of AUC "lost"; i.e. the AUC lift the variable provided.

idx_50m_entry and idx_win_ground_ball both lift AUC strongly

idx_less_clangers and Inside_50s show some AUC lift

idx_less_frees had a negligible impact to AUC, while idx_win_aerial_ball actually dropped AUC, probably due to multicollinearity with inside 50 measures

| Variable | AUC lift | Interpretation |
|---|---|---|
| idx_50m_entry | +0.04 | Marking more 50m entries than opponents is very important |
| idx_win_ground_ball | +0.03 | Winning more contested ground ball than opponents is very important |
| idx_less_clangers | +0.01 | Making less unforced errors than opponents is somewhat important |
| Inside_50s | +0.01 | Having a high proportion of top ten players who get the ball inside 50m is somewhat important |
| idx_less_frees | ±0.00 | Giving away less free kicks than opponents is not really important at all |
| idx_win_aerial_ball | –0.01 | *We ignore this finding due to moderate correlation between idx_win_aerial_ball and idx_50m_entry* |

The most significant variables in the Logistic Regression model align with those in the Decision Tree – this is expected and confirms their relative importance in this analysis

**This analysis has been an insightful proof-of-concept to the value of data to AFL. Beyond the usual sports focus areas of fan engagement, direct marketing and ticket sales, it shows that player performance data can be used as a predictor to match results and more importantly to define what factors are vital to winning.**

## The analysis confirmed conventional wisdom about how to win AFL matches

- Coaches and talent managers should focus on finding players who are excellent at winning contested ground balls
- They should also seek accuracy in entering the forward 50m arc
- It helps to have more inside 50m contributors generally

## Player performance data is a surprisingly strong predictor of match outcomes

- This is surprising due to the adversarial nature of sport where one would expect there are many ways to win
- Decision Trees and Logistic Regression models were both reasonable prediction methods, including when 2017 player data was redacted

## Further predictive accuracy could be found… but probably isn't valuable

- Clustering showed promise as a way of boosting accuracy
- Trying different numbers of clusters might find improvement, as could setting individual threshold probability values for each Logistic Regression model
- A more robust method of predicting player metrics for the predicted season (than prior season averages) is probably needed
- However, finding the $n^{th}$ degree of accuracy in prediction may not be valuable – the results must be interpretable to be useful to the football department (coach or talent manager)

**Where to from here?**

**There are almost unlimited opportunities to expand this analysis to keep seeking deeper insights and greater competitive advantage.**

## Combining multiple, proprietary data sources could lead to huge tactical gains

- This analysis only used free data – much more detail is in Champion Data
- Combining this with player GPS data could lead to invaluable insights into play-by-play outcomes
- What if we could go from understanding we need to win the contested ball, to knowing how to win it with the right running patterns and team plans?

## Going back to sports analytics roots: Moneyball version 2

- An understanding of how players can successfully combine to create winning plays could lead to the next revolution in Moneyball
- The relative value of players in the draft and trade period probably has more to do with a team's specific needs than taking the best midfielders and big forwards on offer
- Are we sure we've taken the lessons of Moneyball in AFL?

## This approach can be applied to any sport with rich data

- The use of data analysis in short-form cricket (T20) is already defining competitive advantage
- It is evident that teams have specific plans by game phase, by match-up and by player