

# Equine Injury Database – models, risk factors and prediction

Tim.Parkin@Glasgow.ac.uk  
@ThoroughbredHN



## Introduction

- EID since 2008
- Raw descriptive statistics
- Modelling to identify risk factors
- Testing the predictive ability of the models
- The next 12 months

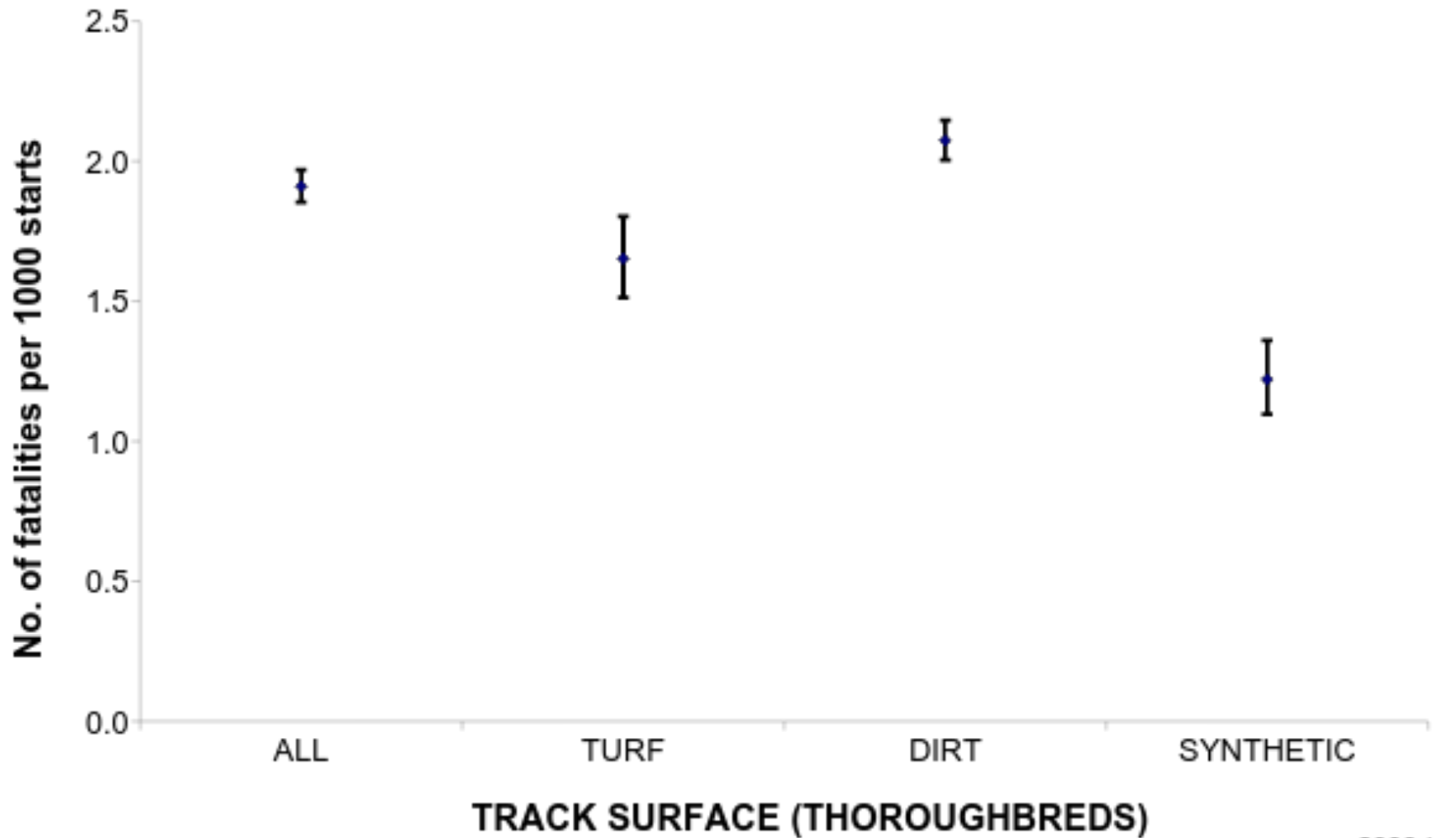
## Definitions of race day fatalities

- Within 72 hours of race
- Estimates now by calendar year
- Point estimates and 95% confidence intervals
- Now producing multivariable models that account for inter-relationships between variables



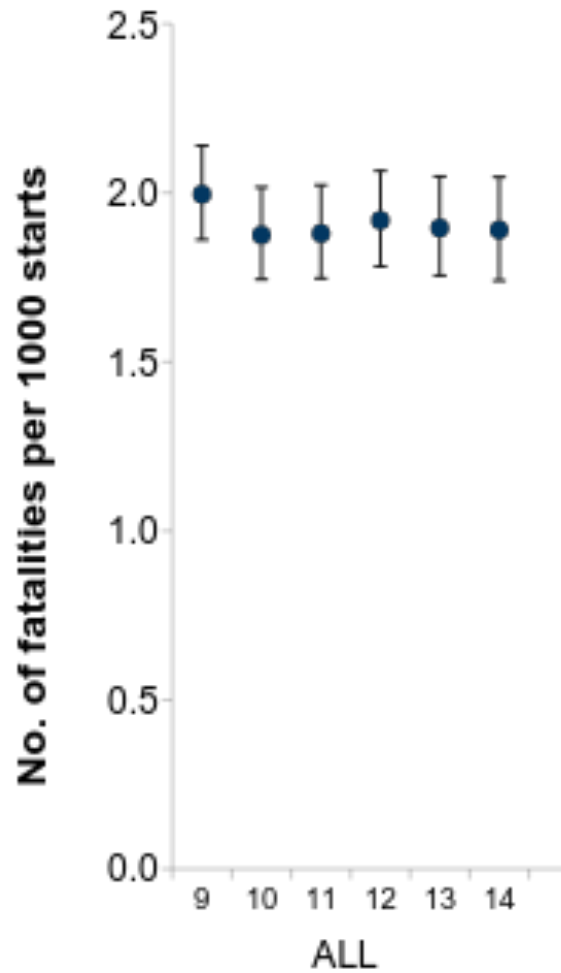


## By surface type





## By surface type 2009 - 2014

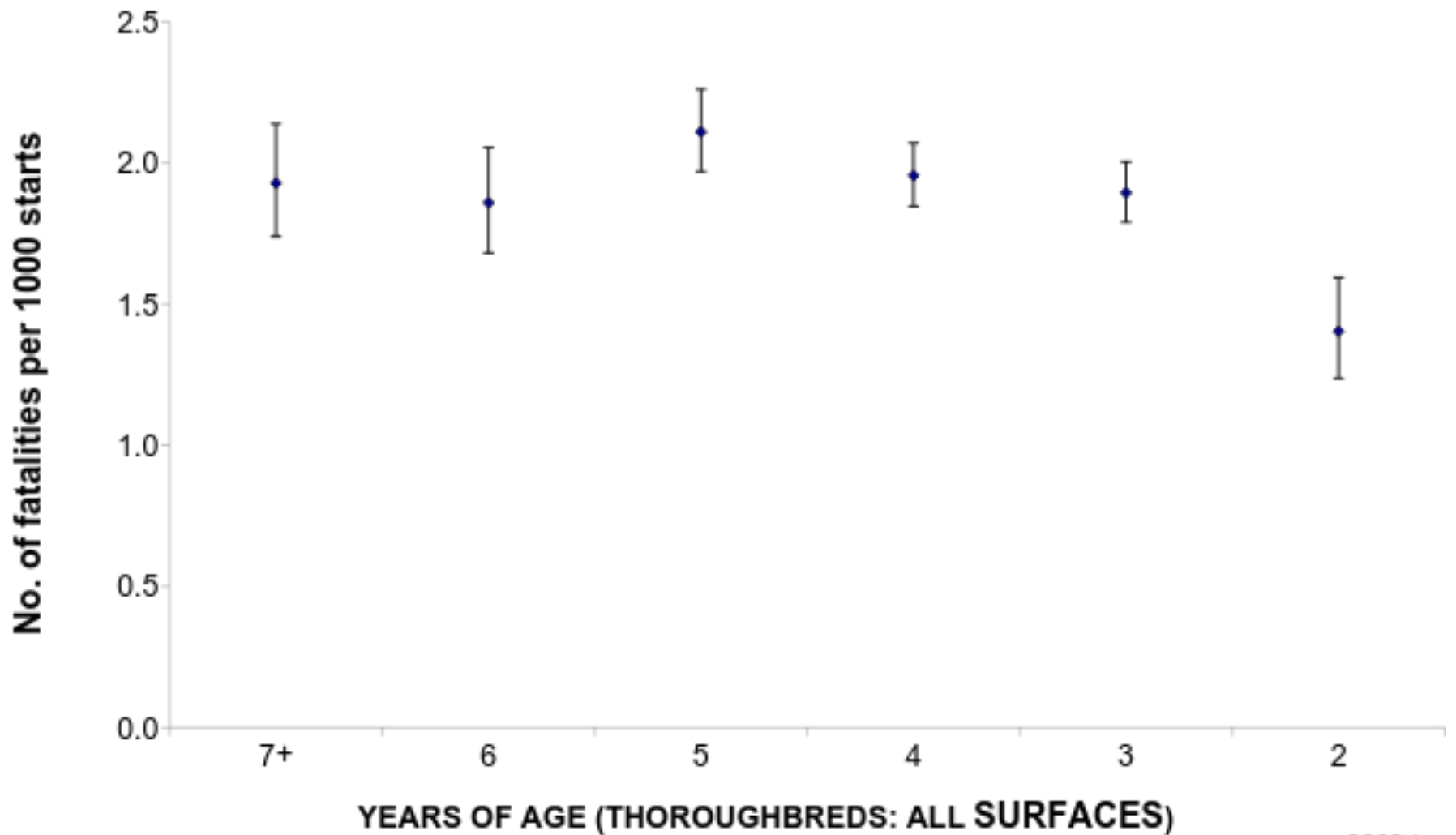


TRACK SURFACE (THOROUGHBREDS)

2009 to 2014

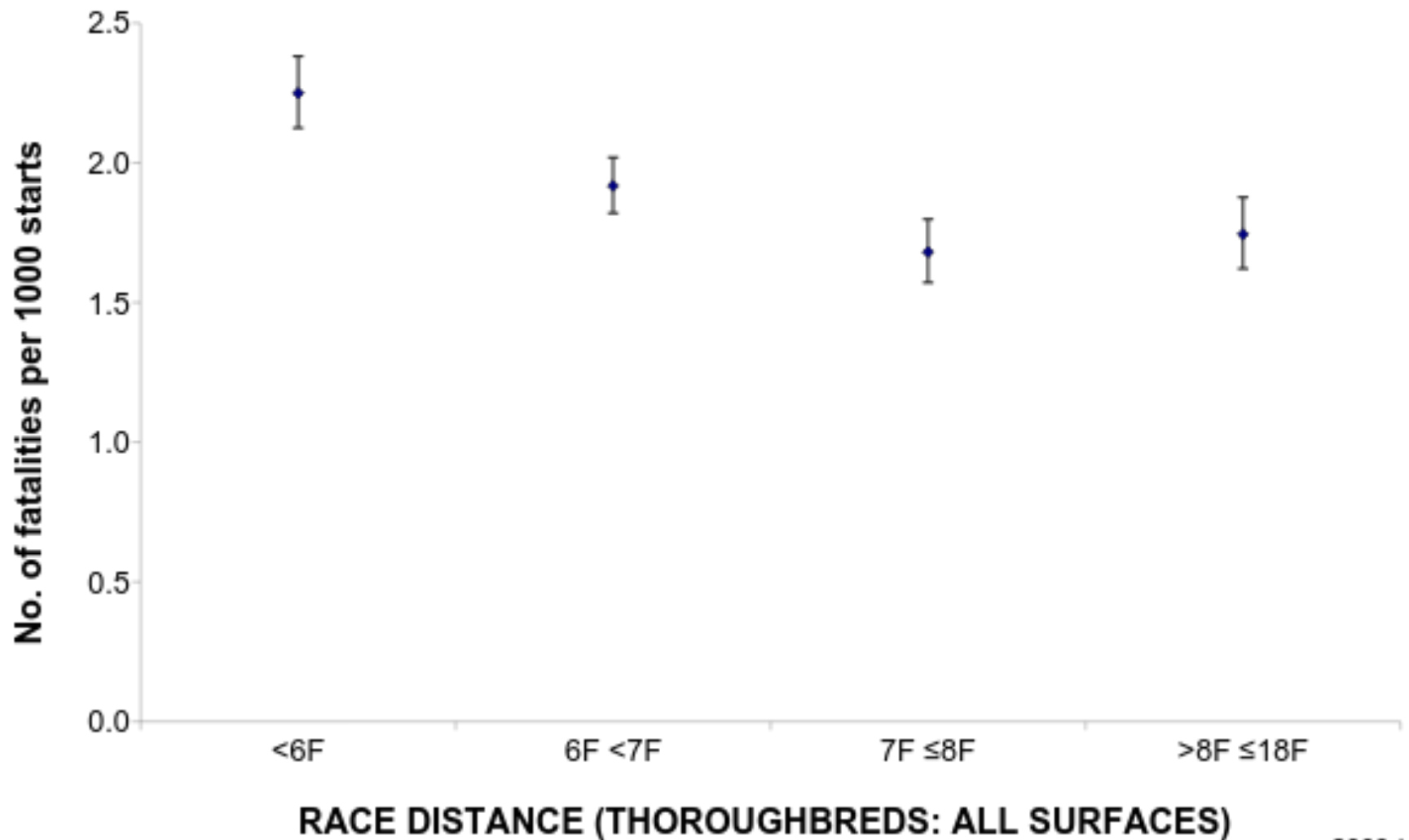


## By age





## By race distance



## Models

- Account for effect of risk factor upon each other and the risk of fatal injury
- National and Track Specific models
- National models built using 6-years of data
  - All races and claiming races only
- Track-specific models for 8 tracks
  - Dependent on sufficient number of starts at these tracks to provide adequate statistical power

## National and track-specific models

- 2.2 million starts
- 150,000 horses
- 94% of all starts in North America (2009 to 2014)
- A selection of important risk factors:
  - Previous EID injuries
  - Appearance on a vet list
  - Time with same trainer
  - Race distance
  - Surface
  - Previous race history
  - Drop in claim price since previous race
  - Age at first race

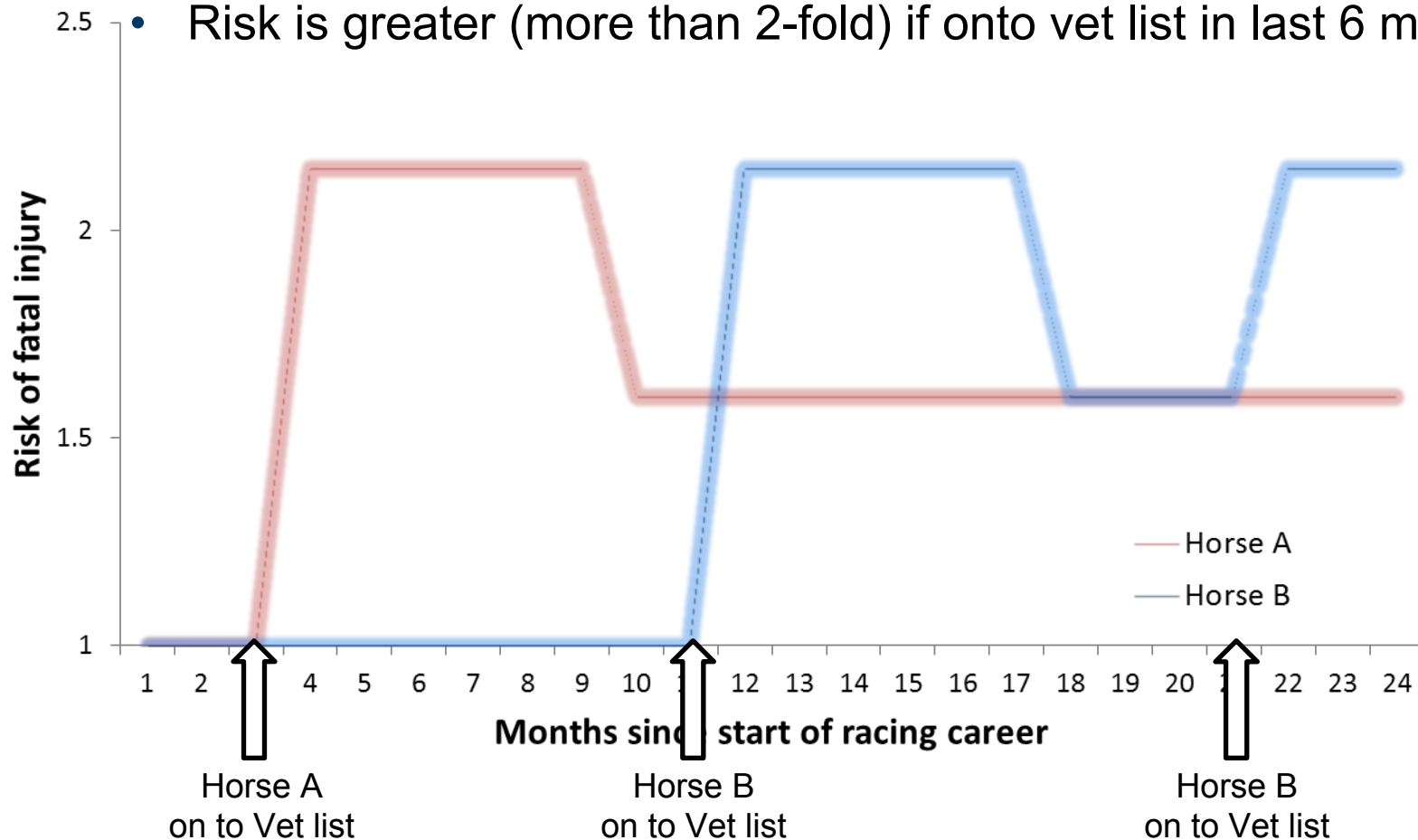
## Previous injuries

- Note: Only EID reported injuries
  - Actual relationship could be much bigger
- For every extra previous injury the risk of fatal injury during racing increases by 30%
  - Compared with a horse with no previous EID injury:
    - 1 previous injury – 30% greater risk (about 2% of starts)
    - 2 previous injuries – 70% greater risk (0.1% of starts)
    - 3 previous injuries – 110% greater risk (0.01% of starts)
- Could be much more valuable IF we could include injuries that are not recorded on EID



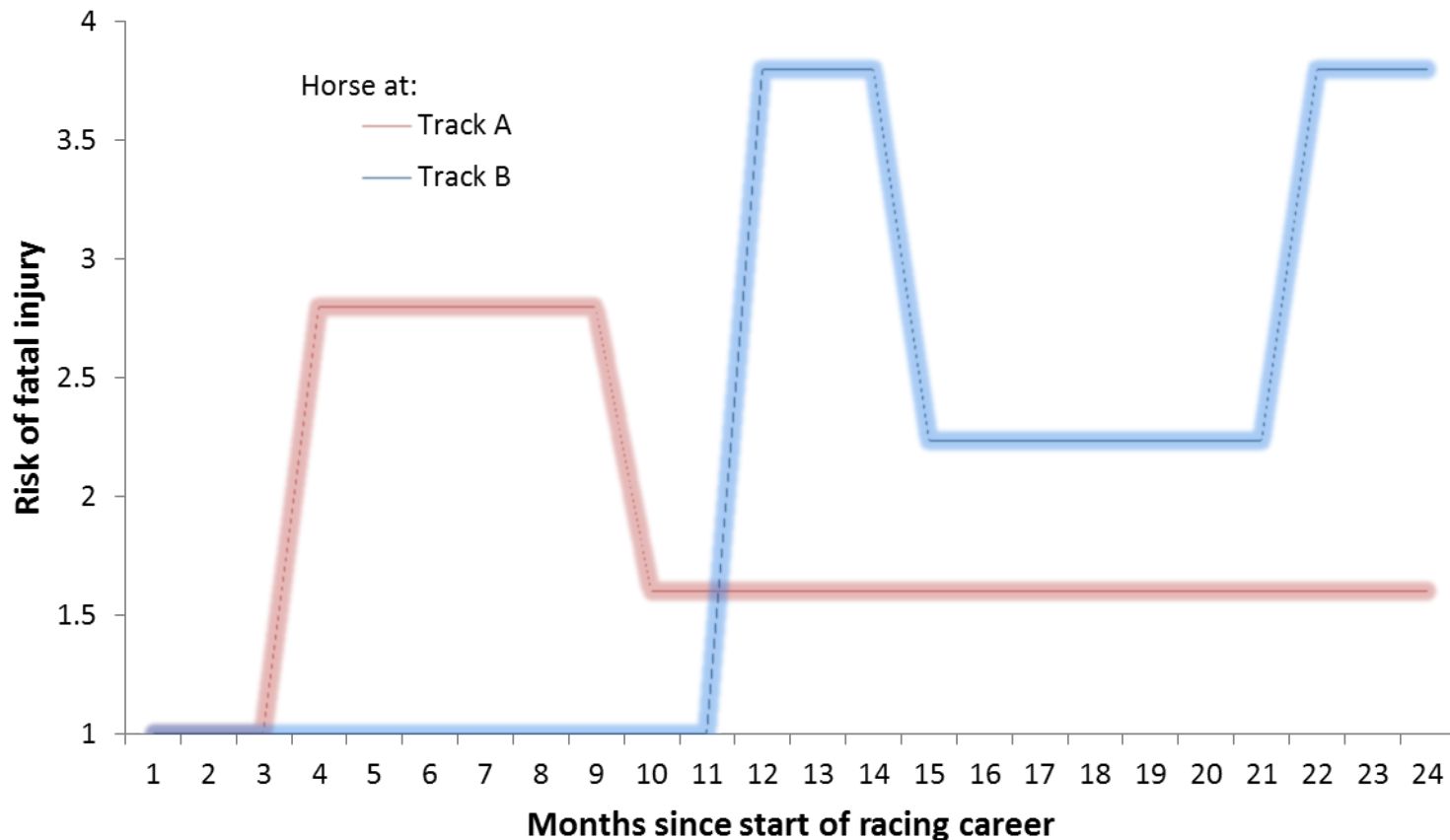
## Vet list

- No difference if include when come off the vet list
- Risk does not return to 'base line' once been on the vet list
- Risk is greater (more than 2-fold) if onto vet list in last 6 months

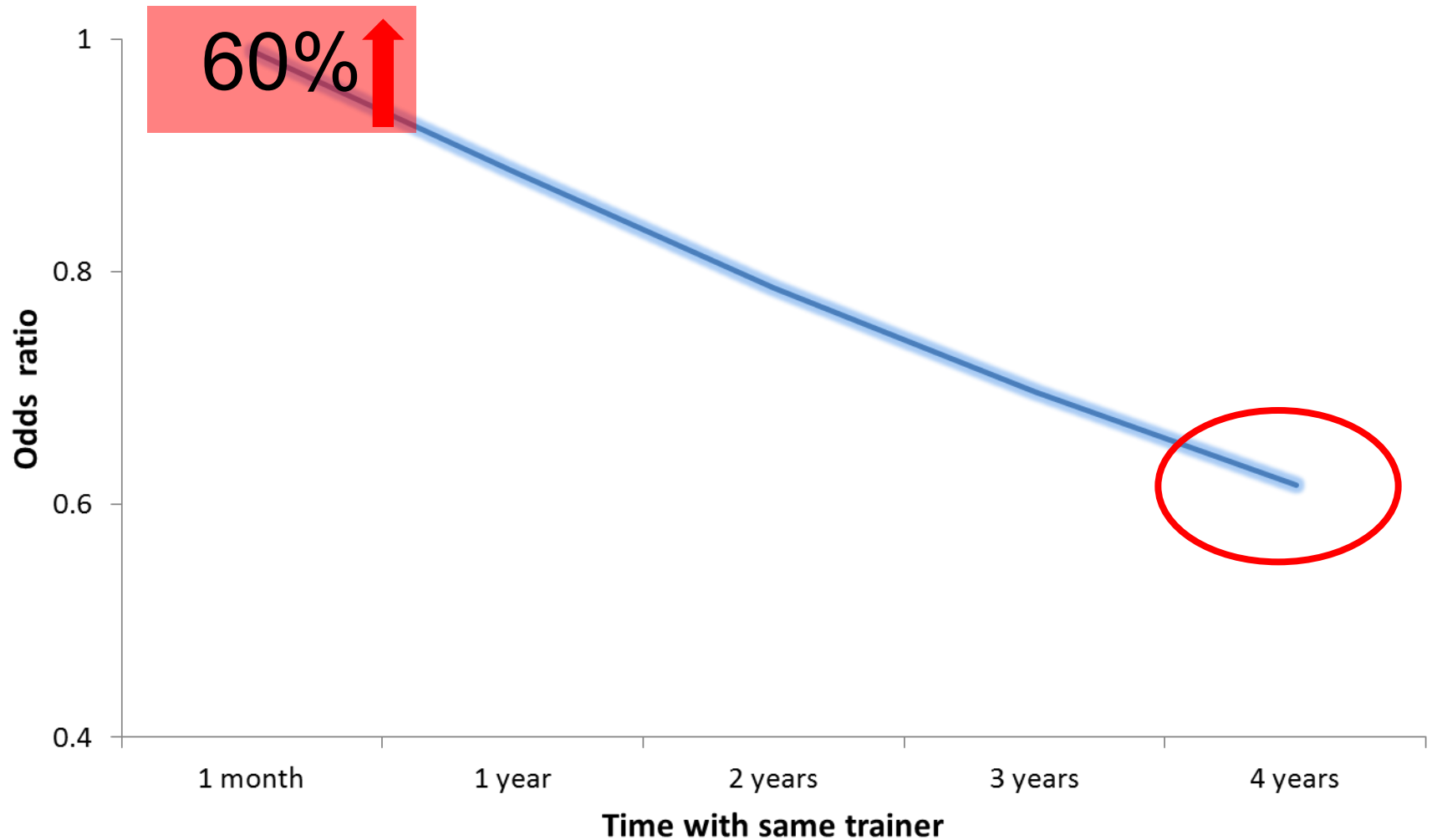


## Vet list

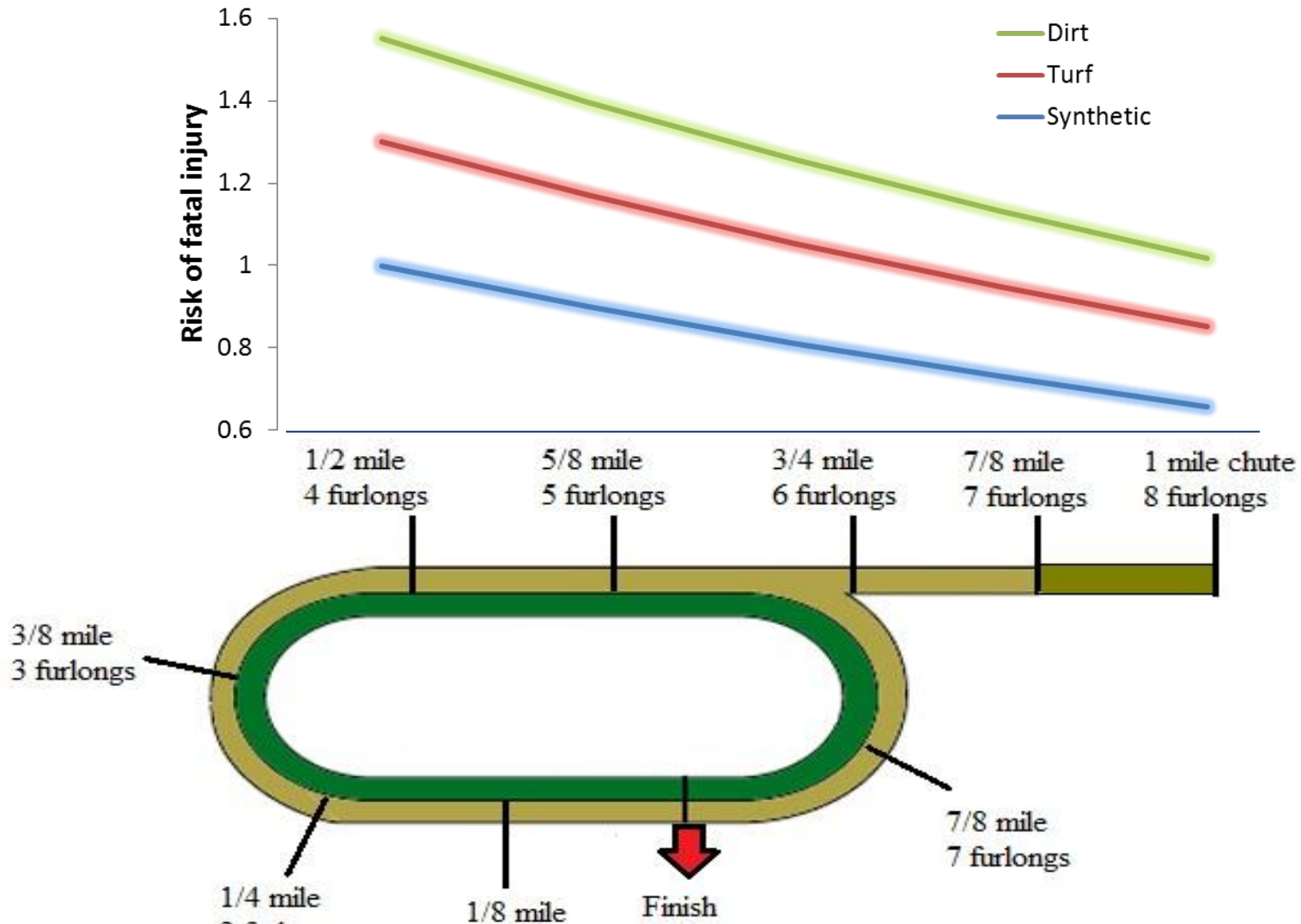
- Each track is different
  - Amount of time after onto vet list that risk is increased
  - After onto vet list 'baseline' risk



## Time with same trainer



## Surface and race distance



## Previous race history

### Horse A

2014

January							February							March								
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S		
1	2	3	4	5	6	7				1	2	3	4					1	2	3		
8	9	10	11	12	13	14		5	6	7	8	9	10	11		4	5	6	7	8	9	10
15	16	17	18	19	20	21		12	13	14	15	16	17	18		11	12	13	14	15	16	17
22	23	24	25	26	27	28		19	20	21	22	23	24	25		18	19	20	21	22	23	24
29	30	31						26	27	28	29					25	26	27	28	29	30	31

April							May							June								
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S		
1	2	3	4	5	6	7				1	2	3	4	5					1	2		
8	9	10	11	12	13	14		6	7	8	9	10	11	12		3	4	5	6	7	8	9
15	16	17	18	19	20	21		13	14	15	16	17	18	19		10	11	12	13	14	15	16
22	23	24	25	26	27	28		20	21	22	23	24	25	26		17	18	19	20	21	22	23
29	30							27	28	29	30	31				24	25	26	27	28	29	30
July							August							September								
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S		
1	2	3	4	5	6	7				1	2	3	4							1		
8	9	10	11	12	13	14		5	6	7	8	9	10	11		2	3	4	5	6	7	8
15	16	17	18	19	20	21		12	13	14	15	16	17	18		9	10	11	12	13	14	15
22	23	24	25	26	27	28		19	20	21	22	23	24	25		16	17	18	19	20	21	22
29	30	31						26	27	28	29	30	31		23	24	25	26	27	28	29	30
October							November							December								
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S		
	1	2	3	4	5	6					1	2	3							1		
7	8	9	10	11	12	13		4	5	6	7	8	9	10		2	3	4	5	6	7	8
14	15	16	17	18	19	20		11	12	13	14	15	16	17		9	10	11	12	13	14	15
21	22	23	24	25	26	27		18	19	20	21	22	23	24		16	17	18	19	20	21	22
28	29	30	31					25	26	27	28	29	30		23	24	25	26	27	28	29	30

### Horse B

2014

January							February							March						
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S
1	2	3	4	5	6	7				1	2	3	4					1	2	3
8	9	10	11	12	13	14	5	6	7	8	9	10	11	12	13	14	15	16	17	18
15	16	17	18	19	20	21	12	13	14	15	16	17	18	19	20	21	22	23	24	25
22	23	24	25	26	27	28	19	20	21	22	23	24	25	26	27	28	29	30	31	
29	30	31					26	27	28	29										

April							May							June						
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S
1	2	3	4	5	6	7														
8	9	10	11	12	13	14	6	7	8	9	10	11	12	13	14	15	16	17	18	19
15	16	17	18	19	20	21	13	14	15	16	17	18	19	20	21	22	23	24	25	26
22	23	24	25	26	27	28	20	21	22	23	24	25	26	27	28	29	30	31		
29	30						27	28	29	30	31									

July							August							September						
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S
1	2	3	4	5	6	7														
8	9	10	11	12	13	14	6	7	8	9	10	11	12	13	14	15	16	17	18	19
15	16	17	18	19	20	21	13	14	15	16	17	18	19	20	21	22	23	24	25	26
22	23	24	25	26	27	28	20	21	22	23	24	25	26	27	28	29	30	31		
29	30	31					27	28	29	30	31									

October							November							December						
S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S
	1	2	3	4	5	6														
7	8	9	10	11	12	13	4	5	6	7	8	9	10	11	12	13	14	15	16	17
14	15	16	17	18	19	20	11	12	13	14	15	16	17	18	19	20	21	22	23	24
21	22	23	24	25	26	27	18	19	20	21	22	23	24	25	26	27	28	29	30	31
28	29	30	31				25	26	27	28	29	30								

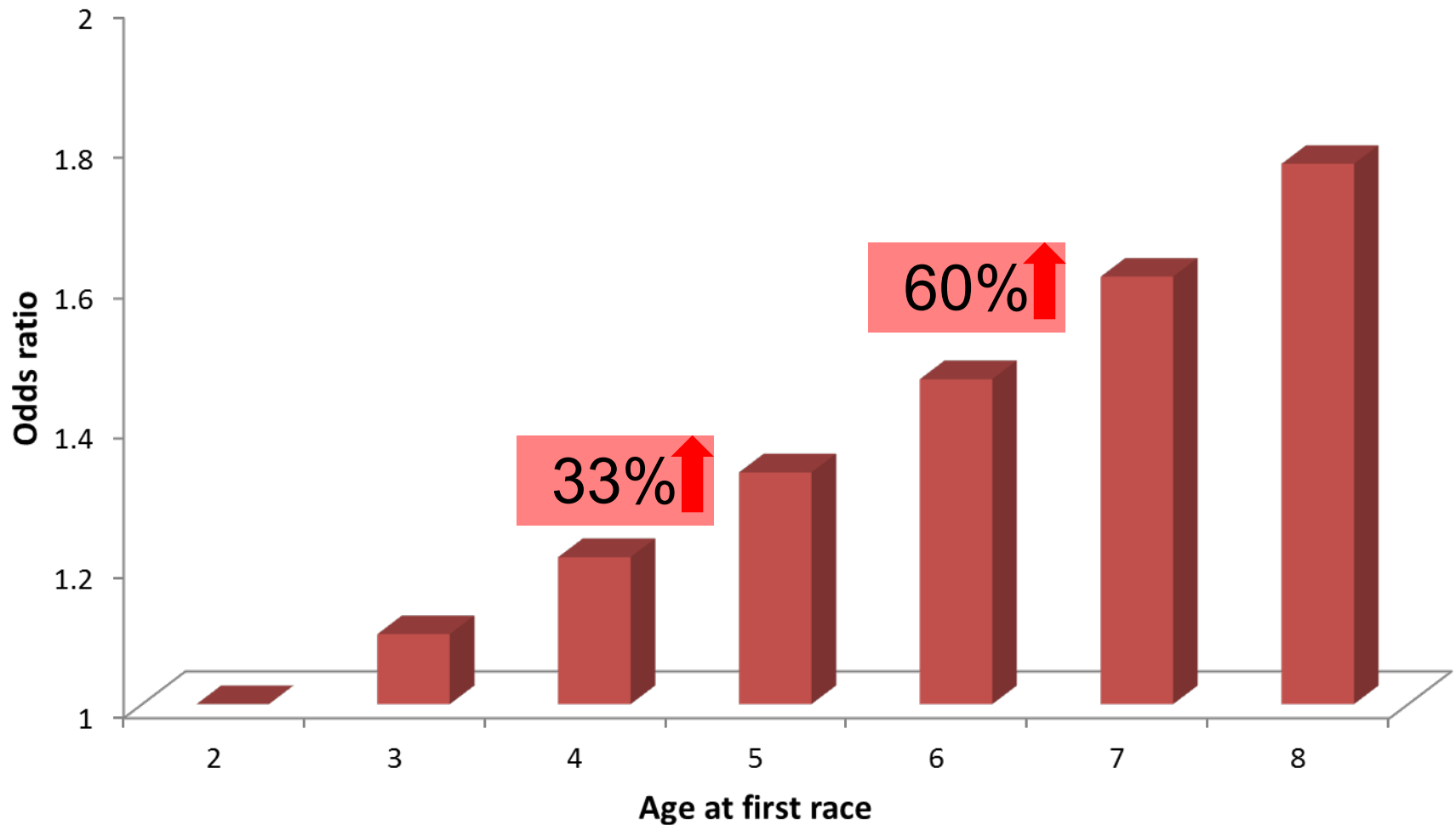
## Drop in claiming price since last race



Little change since last race (+/- \$500)	Drop of between \$500 and \$10,000	Drop of more than \$10,000
Reference	14%	16%

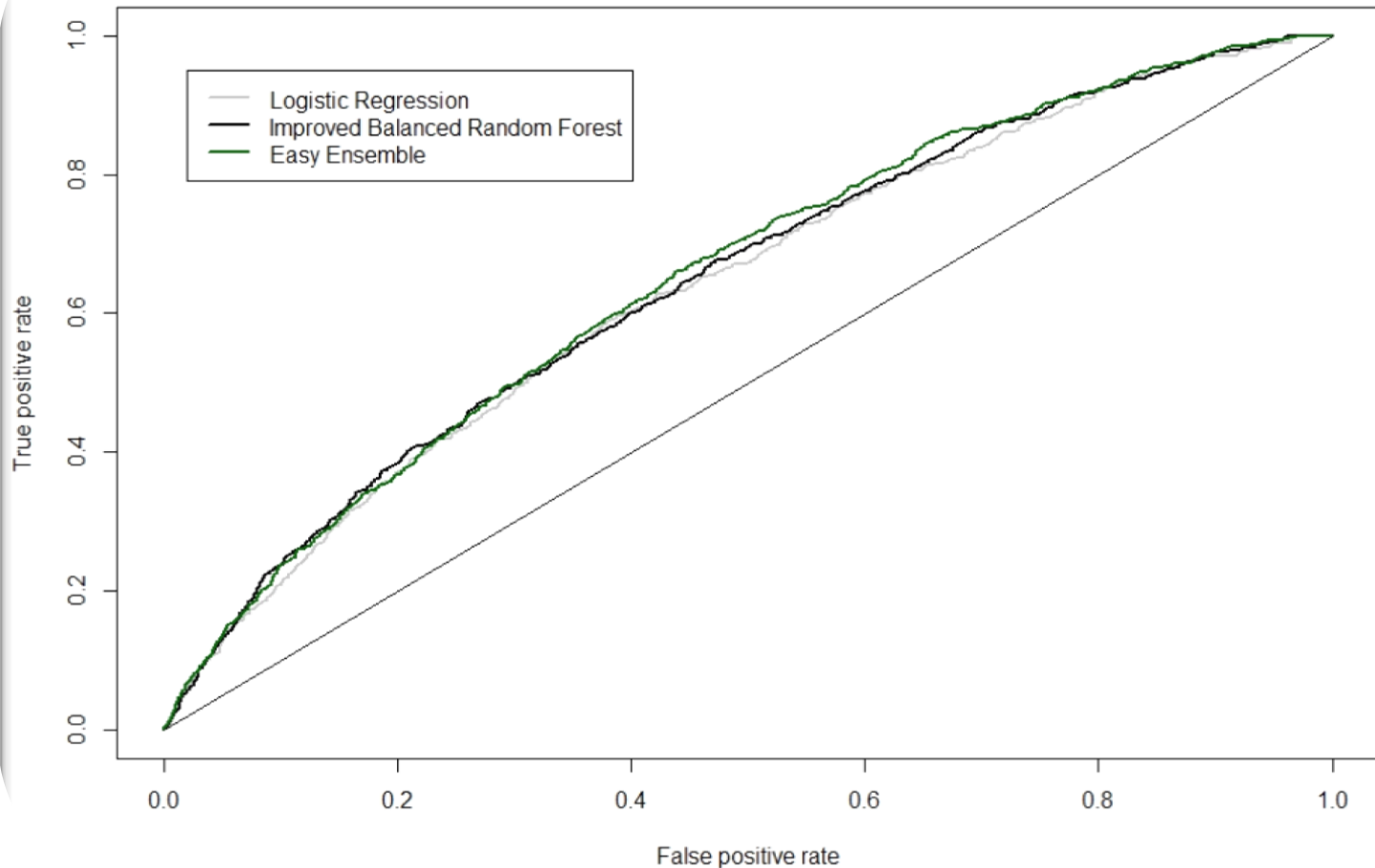


## Age at first race



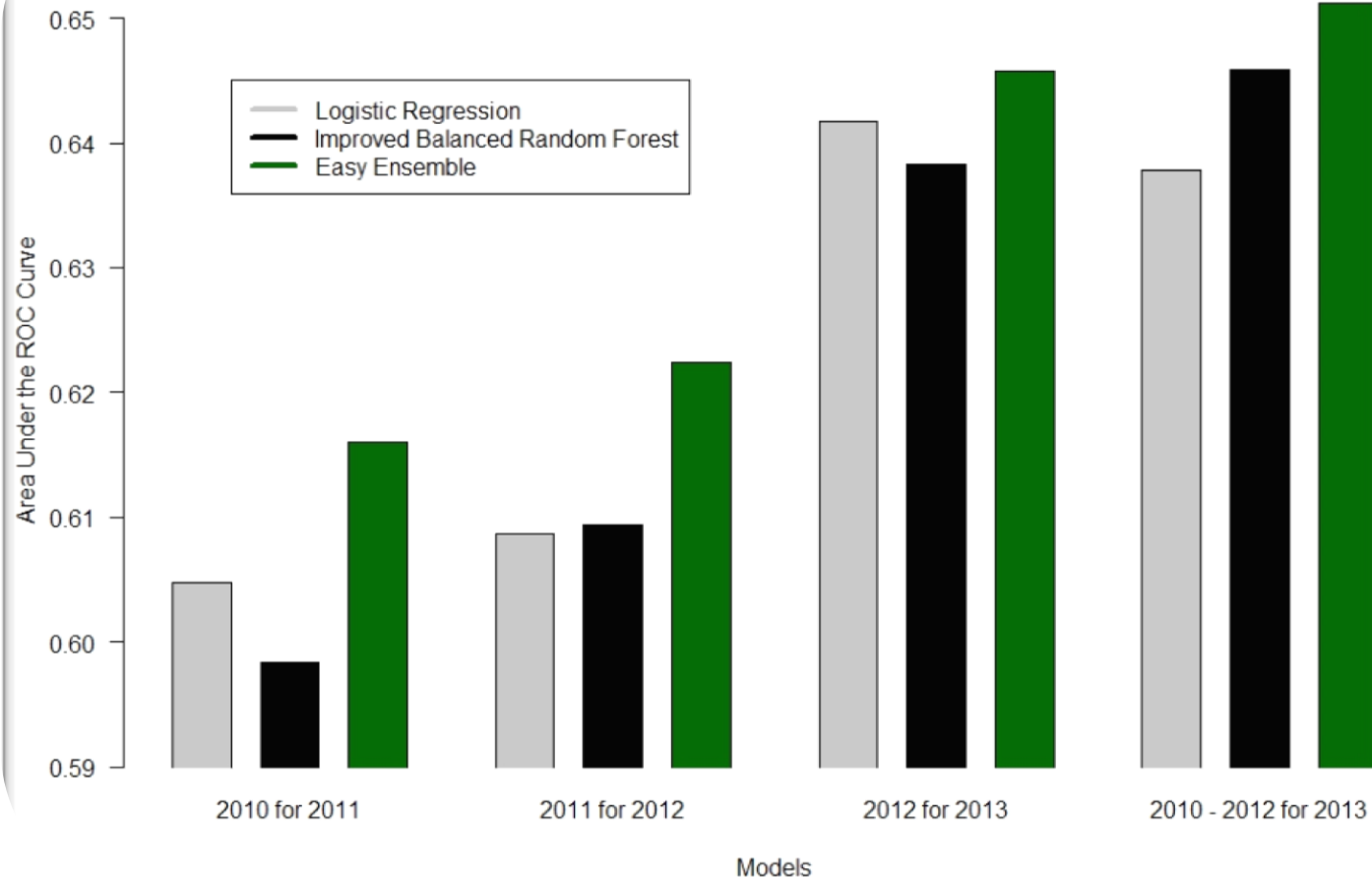
## Predictive ability of models ~ 65%

### Area under the ROC curve



## Predictive ability

### Area under the ROC curve



## Variable predictive ability at different tracks

- AUC at different tracks:
- Range from 53% to 68%
  - Most individual track models are slightly less predictive
- A lot of ‘local’ factors that are simply missed in EID or not recorded at all
- Importance of ‘local’ knowledge and working with those on the ground at different tracks

## Overall 3-fold greater risk for top 5%

Multivariable Logistic Regression									
Quantiles of Score		0-5%	5-20%	20-35%	35-50%	50-65%	65-80%	80-95%	95-100%
Relative Horse Risk		0.46	0.47	0.66	0.87	0.94	1.09	1.60	2.71
Improved Balanced Random Forest									
Quantiles of Score		0-5%	5-20%	20-35%	35-50%	50-65%	65-80%	80-95%	95-100%
Relative Horse Risk		0.49	0.45	0.62	0.72	1.05	1.14	1.46	3.26
Easy Ensemble									
Quantiles of Score		0-5%	5-20%	20-35%	35-50%	50-65%	65-80%	80-95%	95-100%
Relative Horse Risk		0.43	0.47	0.67	0.74	0.94	1.20	1.48	3.10

## Predictive ability of the models

- How close are we at being able to more accurately find horses of interest BEFORE they race?
- Topping out at 65% on predictive models?
  - Maybe best possible
  - Unmeasured variables
  - Inherent variability i.e. unmeasurable variables
- Risk factors & predictive models for injuries/triage 2+
- Keep with analysis from all tracks
- Focus in on tracks with available training data
- Availability of medical/treatment records?
  - Importance of being on the vet list/previous injuries and from work we have done with BHA



## Further analyses

### Variables

- Number of times on vet list
- Work to get off vs. automatically off vet list
- Type of previous injury (fetlock)
- Vet scratches vs. trainer scratches
- Length of meet

### Fast work data models

### Use of “National” model

- Examine predictive ability of National model for each track

## What to do with this information?

- Is a three-fold difference in risk important for you to be aware of?
  - 3-fold difference in risk between ‘average’ horse and horse in ‘top 5%’
- Which outcome would be best to try to embed within automatic risk profiling for each start?
  - Fatality – clearly important but rare
  - Injury/triage 2+ – important and more common, but case definition will include a lot of variation
  - Fracture of distal limb (fatal and non-fatal)

## Acknowledgements

- US Jockey Club
  - Matt Iuliano
  - Kristin Werner Leshney
  - Jamie Haydon
- University of Glasgow
  - Stamatis Georgopoulos (who did all the work!)