



**CENTRE FOR
SOCIAL DATA ANALYTICS**

Predictive Analytics in Child Welfare:

Implementation of a predictive model to support child maltreatment hotline screening decisions: methods, challenges and lessons learned

Diana Benavides Prado,
Data Scientist

June 2nd, 2017

Context

- Child maltreatment is an international public health problem
- ~40 million children worldwide are subject to some kind of **maltreatment**, per year
- There are over **3 million calls** per year made to **child welfare agencies** concerning abuse or neglect in the US alone.

Children in the Public Benefit System at Risk of Maltreatment Identification Via Predictive Modeling

Rhema Vaithianathan, PhD, Tim Maloney, PhD, Emily Putnam-Hornstein, PhD, Nan Jiang, PhD

Abstract: A growing body of research links child abuse and neglect to a range of negative short- and long-term health outcomes. Determining a child's risk of maltreatment at or shortly after birth provides an opportunity for the delivery of targeted prevention services. This study presents findings from a predictive risk model (PRM) developed to estimate the likelihood of substantiated maltreatment among children enrolled in New Zealand's public benefit system. The objective was to explore the potential use of administrative data for targeting prevention and early intervention services to children and families.

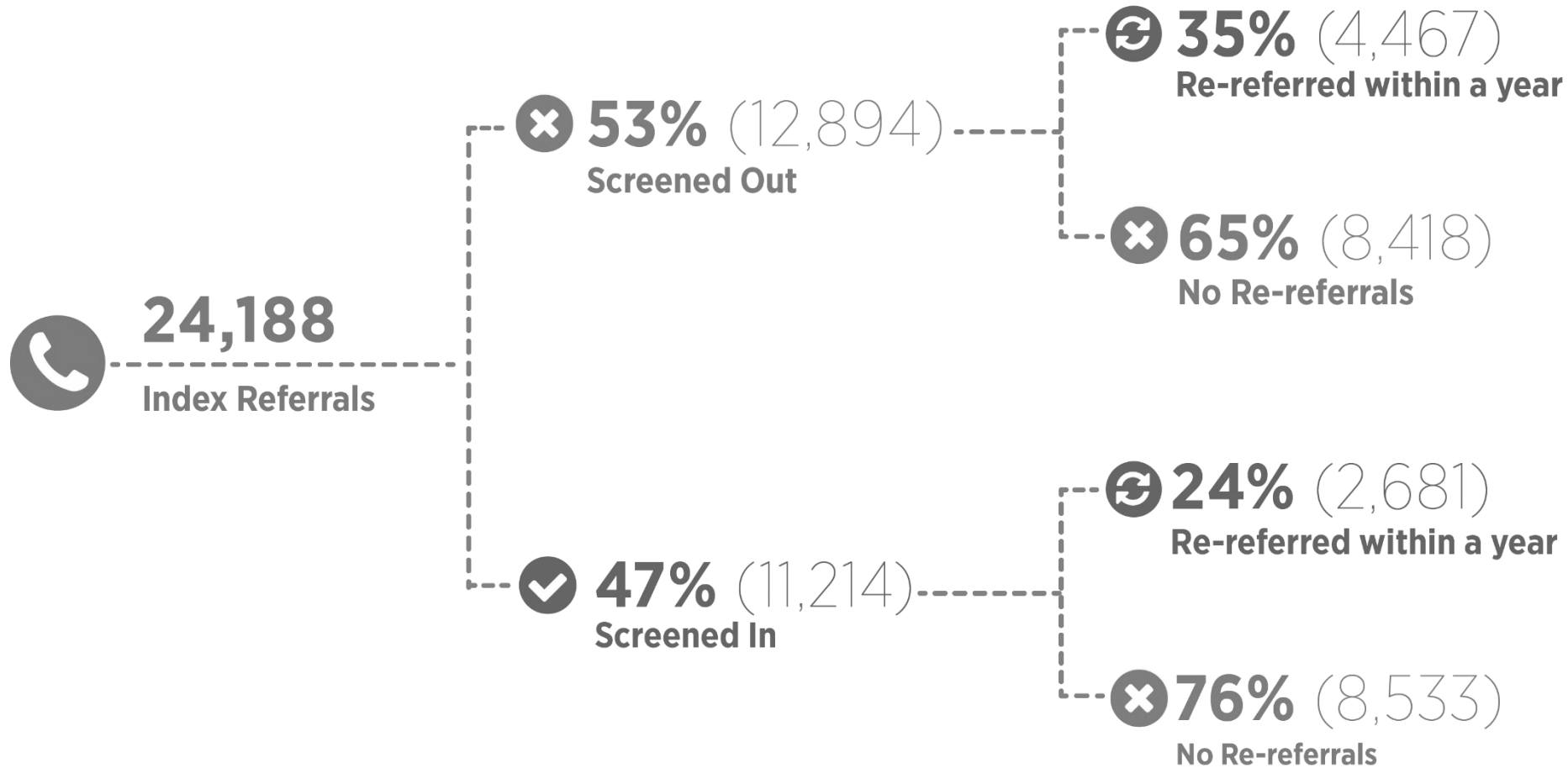
A data set of integrated public benefit and child protection records for children born in New Zealand between January 1, 2003, and June 1, 2006, was used to develop a risk algorithm using stepwise probit modeling. Data were analyzed in 2012. The final model included 132 variables and produced an area under the receiver operating characteristic curve of 76%. Among children in the top decile of risk, 47.8% had been substantiated for maltreatment by age 5 years. Of all children substantiated for maltreatment by age 5 years, 83% had been enrolled in the public benefit system before age 2 years. This analysis demonstrates that PRMs can be used to generate risk scores for substantiated maltreatment. Although a PRM cannot replace more-comprehensive clinical assessments of abuse and neglect risk, this approach provides a simple and cost-effective method of targeting early prevention services.

(Am J Prev Med 2013;45(3):354–359) © 2013 American Journal of Preventive Medicine

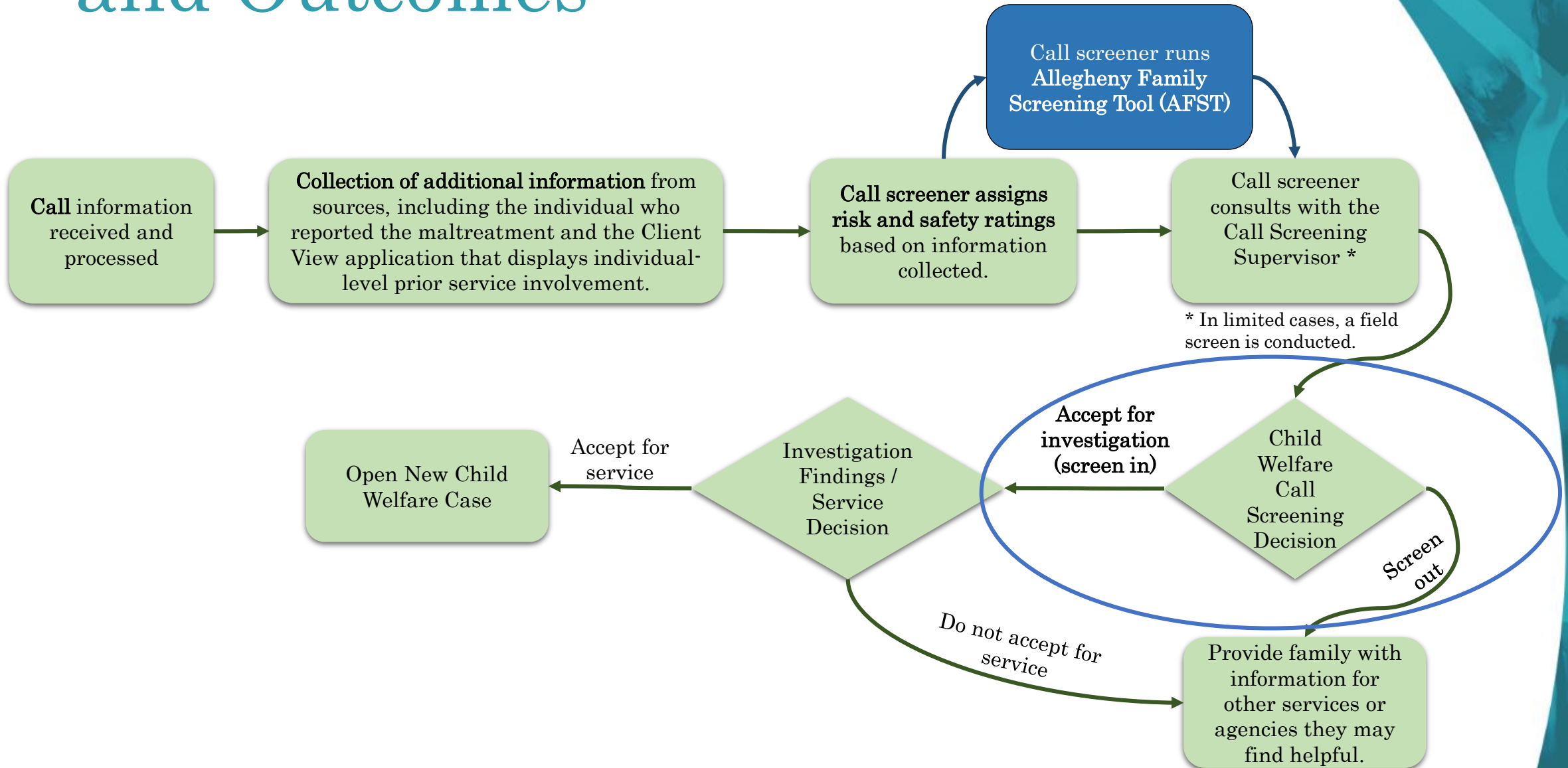
Context: Allegheny County (Pittsburgh, PA, US)



Problem: Screening Decisions and Outcomes

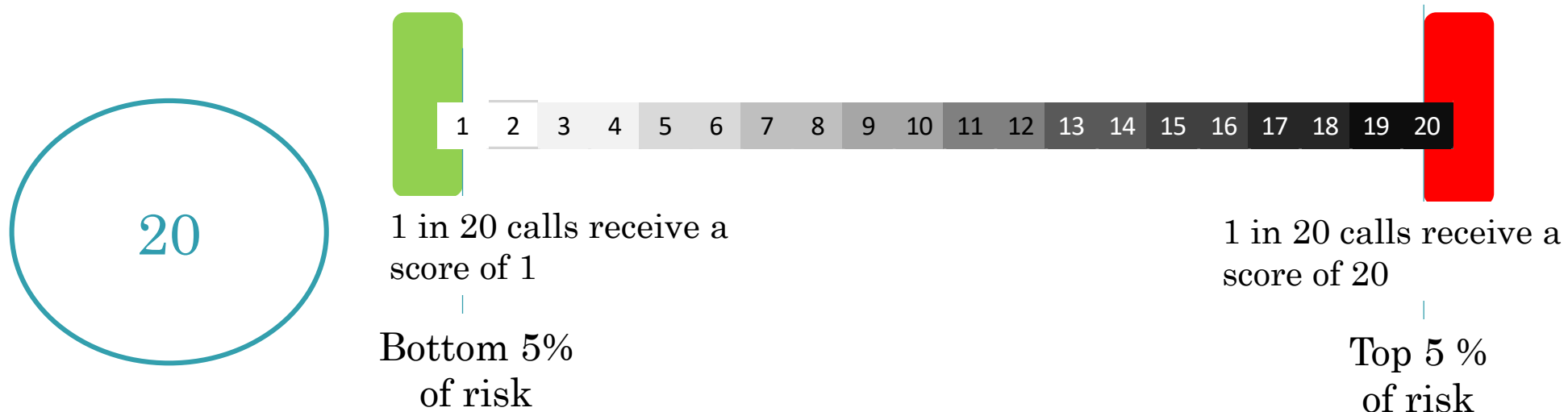


Problem: Screening Decisions and Outcomes

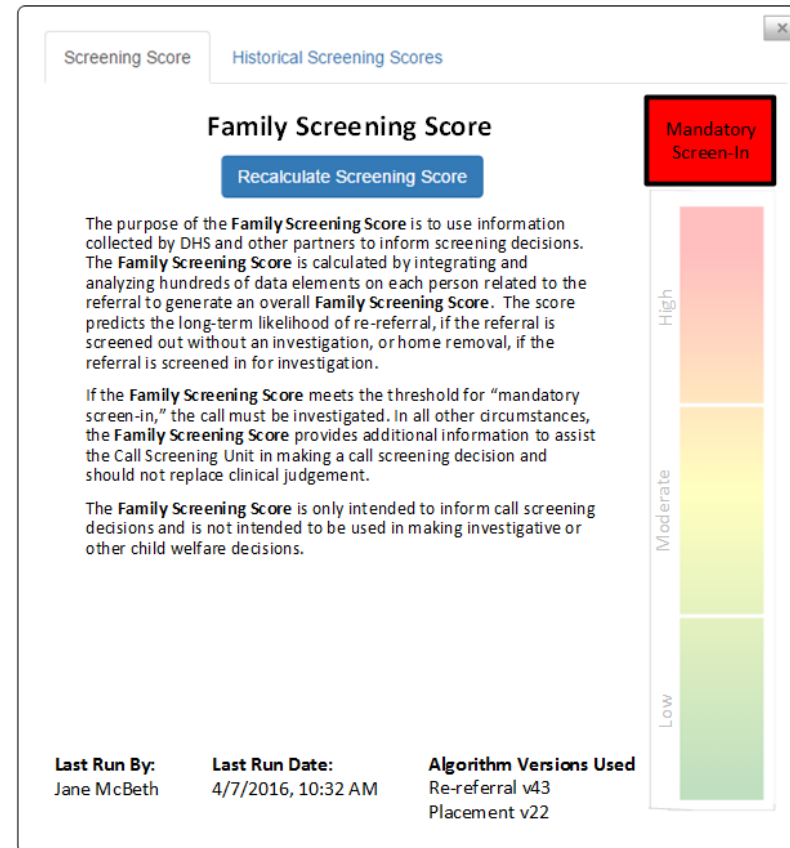
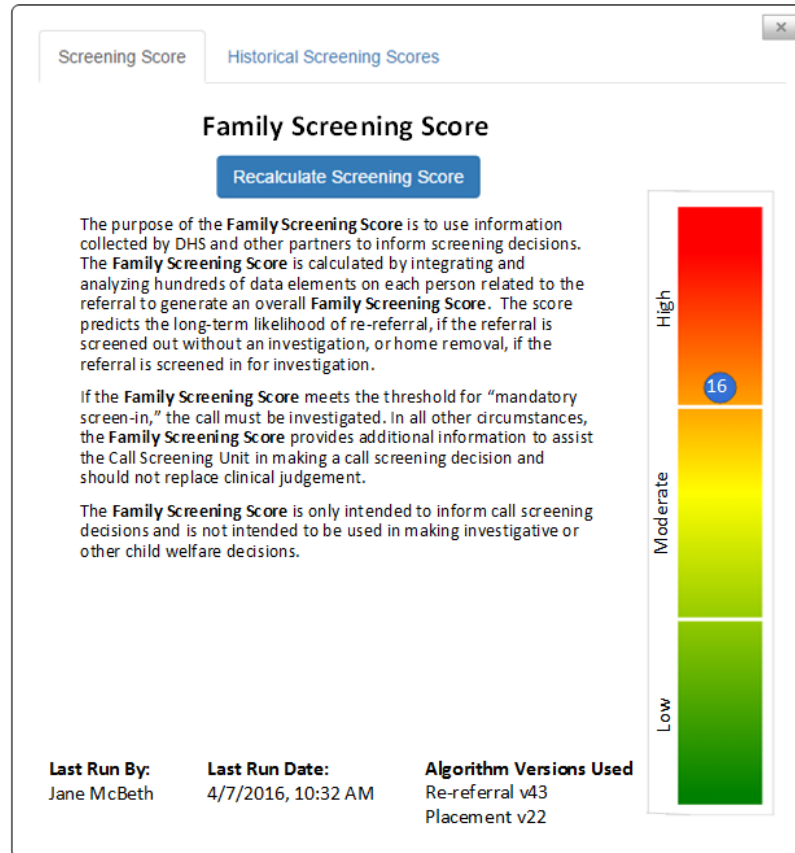


Problem: Screening Decisions and Outcomes

- What's the right question that a predictive model should help to answer?
 - What's the risk that a child will be **re-referred conditional on being screened-out**.
 - What's the risk that a child will be **placed in foster care conditional on being screened-in**.

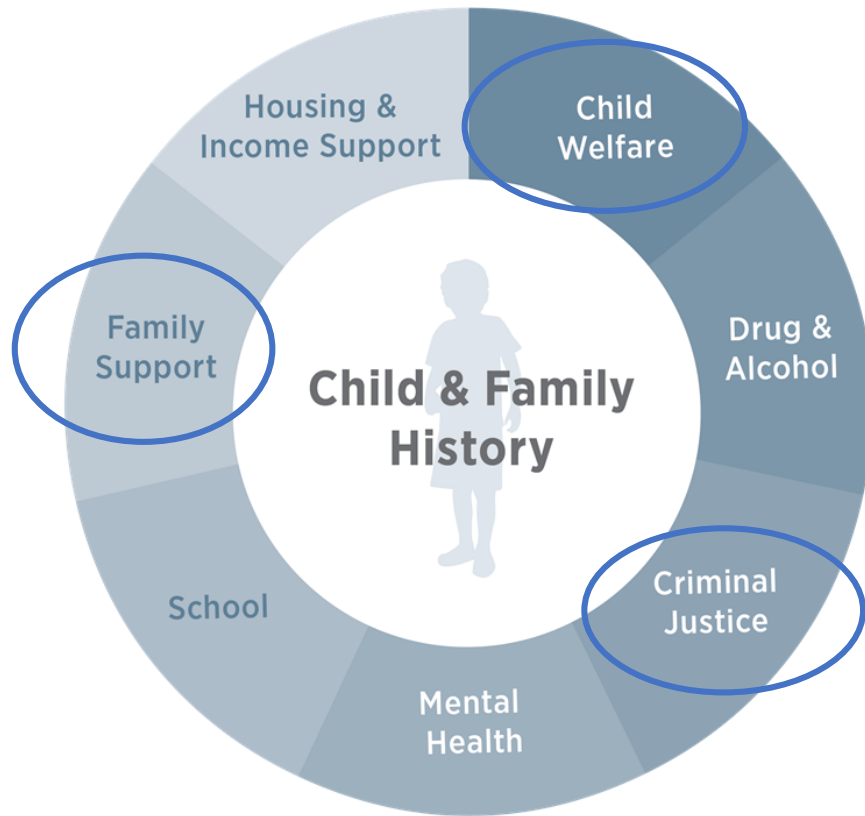


Problem: Screening Decisions and Outcomes



A subset of high scores that are particularly correlated with possible critical incidents are flagged for **mandatory screen-in**, with an override option for unique circumstances

Data



- Several administrative data sources
- More than 800 variables/predictors
- De-identified data for research/analysis purposes
- What's the ideal “shape” of the data for analysis?

Referral ID	Individual ID	Role	Individual Info...	Referral Info...

Modelling v1.0

Extract:

The “right” patterns



Identify correctly the high risk population



Initial criteria:

“Good” balance between true positives
and false negatives

Modelling v1.0

- Logistic Regression

$$\hat{p} = \frac{\exp(b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p)}{1 + \exp(b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p)}$$

59 predictors for the re-referral
model

71 predictors for the placement
model

Modelling v2.0

- Random Forests

(Breiman, 2001)

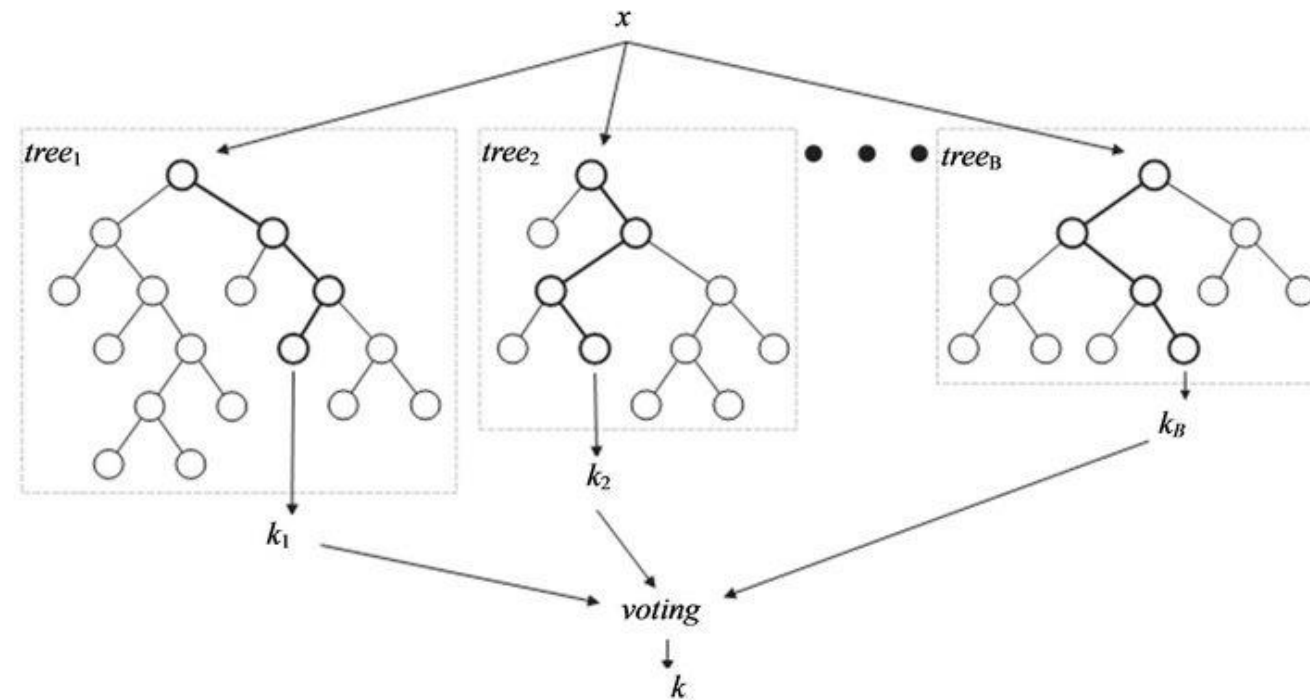


Figure:

http://file.scirp.org/Html/6-9101686_31887.htm

Modelling v2.0

- XGBoost

(Chen, Guestrin, 2016)

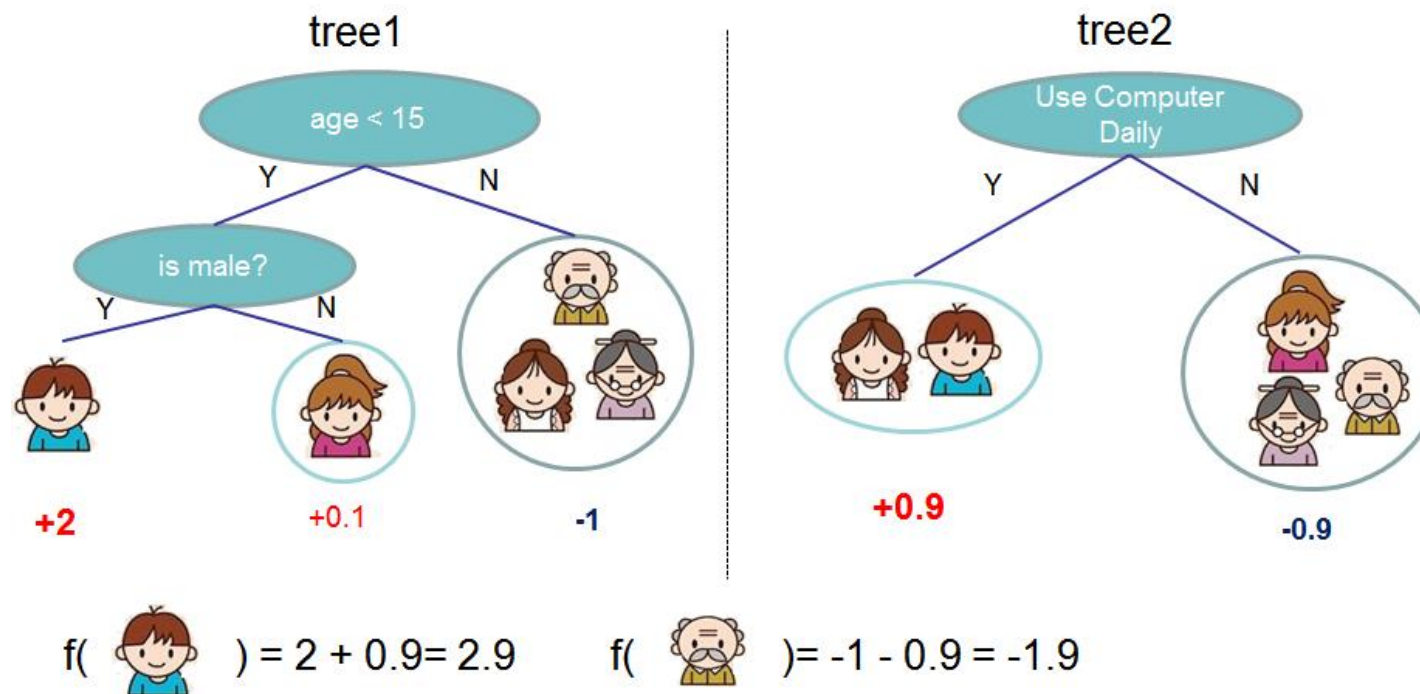


Figure:
<http://xgboost.readthedocs.io/en/latest/model.html>

Modelling v2.0

- SVM

(Vapnik, Chervonenkis, 1963)

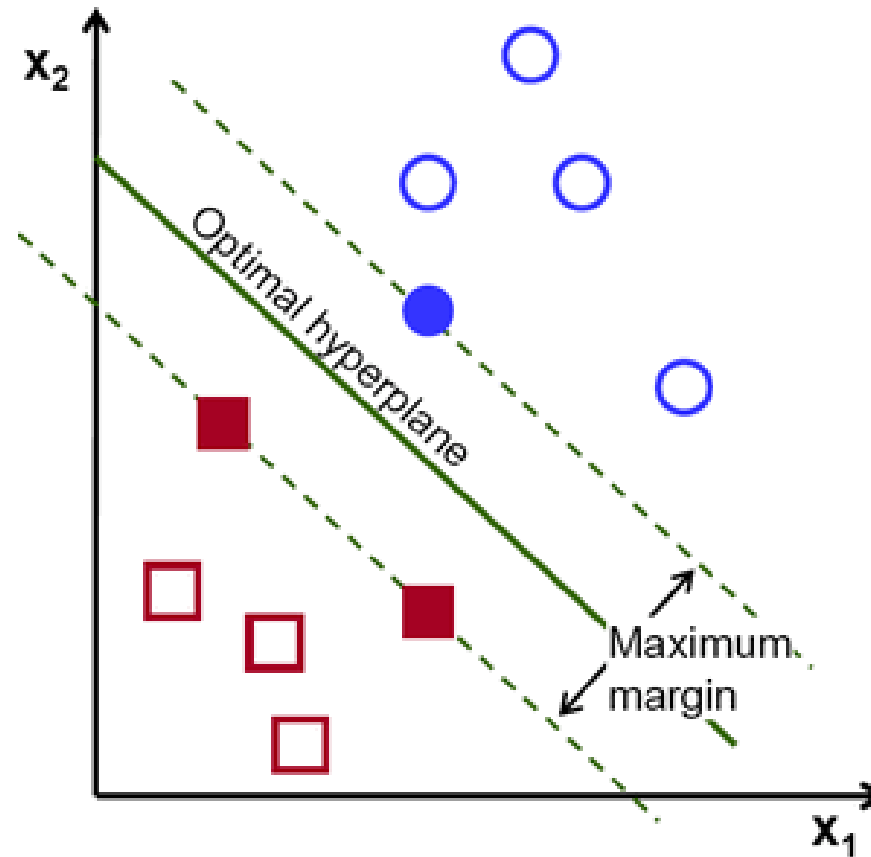


Figure:
http://docs.opencv.org/2.4/doc/tutorials/ml/introduction_to_svm/introduction_to_svm.html

Modelling v2.0

- Logistic Regression

- More familiar to a wider audience
- Not so accurate

- Random Forests

- Human-readable, to an extent (rules)
- Usually accurate for high dimensional data

- XGBoost

- Super-extreme and scalable
- Often complex, though human-readable (to an extent)

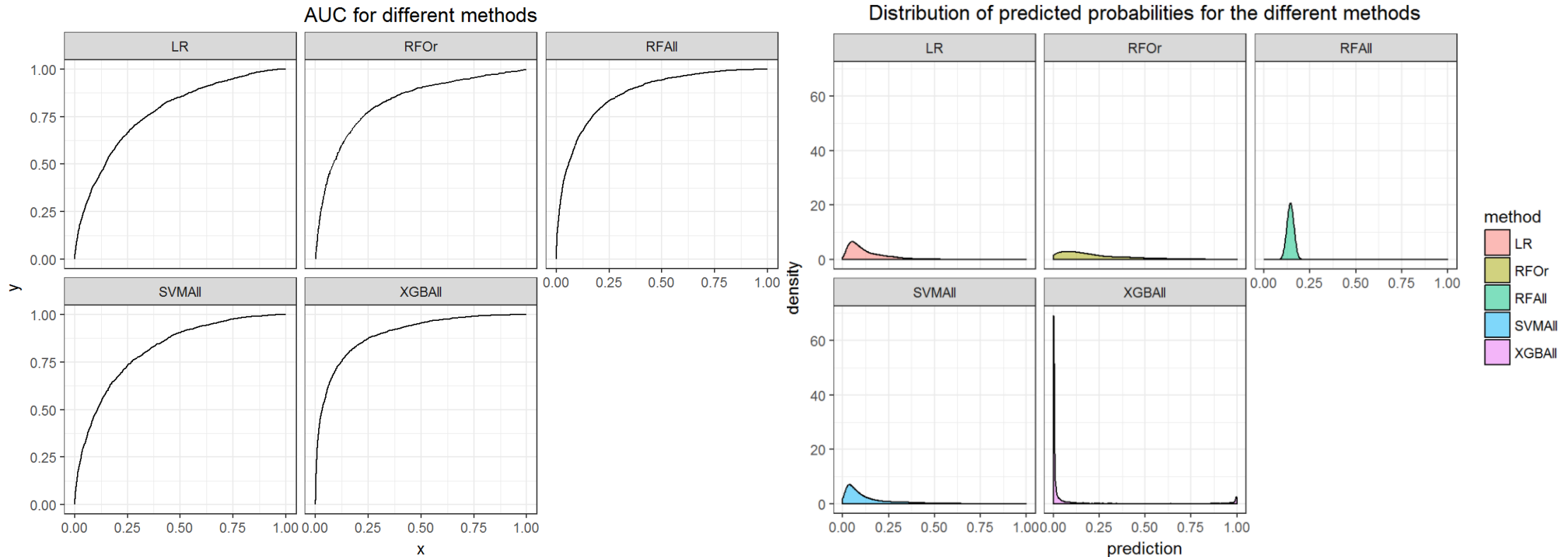
There's
NFL...

- SVM

- Less explanatory
- Usually accurate for high dimensional numeric data

Evaluating what's been learned...

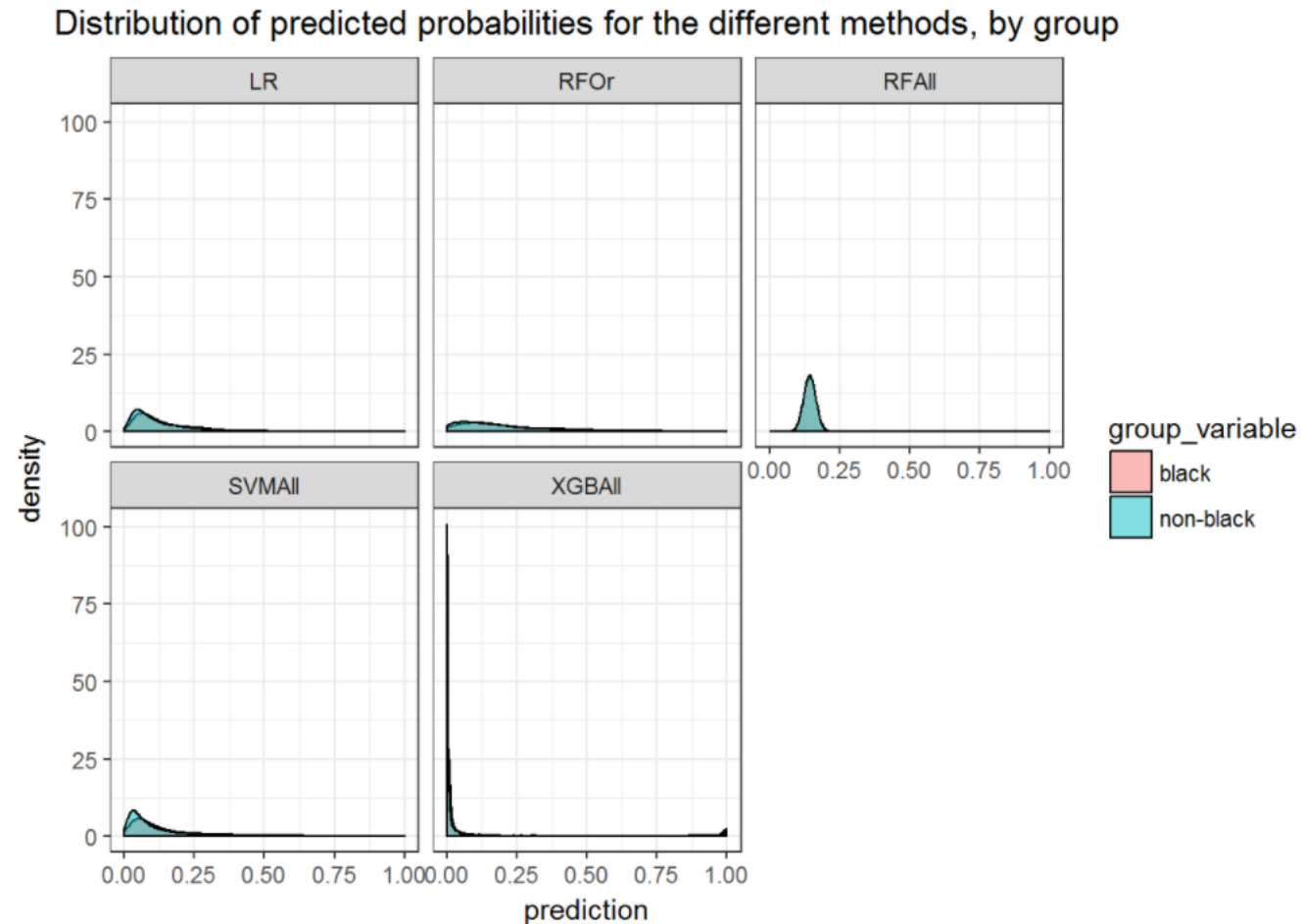
- What are the right metrics to compare – what do we want to achieve?



Evaluating what's been learned...

- What's the impact of what we are learning from data?

Disparities
and
Fairness...



Evaluation and Implementation v2.0

Extract:

The “right” patterns



Identify correctly the high risk population

Subject to:

A “good” explanation



Additional criteria:

Interpretability/Transparency... Fairness

The right to an explanation...

```
-- Model Information ---
Filename: rf_ordinal_model
Schema: waka.classifiers.meta.CPParametersSelection -K 10 -B 1 -M waka.classifiers.trees.RandomForest --P 100 -O -point -attribute-importance=1 100 -m-aleuts 1 -X 0 -M 1.0 -V 0.001 -B 1
Relation: historical_referrals_labelled_allvers_placement_hactVP_apprfist_subpublic_testing_resampled_both-waka.filters.unsupervised.attribute.HumanisticSubclass-Meat-waka.filters.unsupervised.attribute.Ramond-Wid-B-waka.filters.[List of Attributes omitted]
Attribute: 840
-- Classifier model ---
Cross-validated Parameter selection.
Classifier: waka.classifiers.trees.RandomForest
Classifier Options: -P 100 -O -point -attribute-importance=1 100 -m-aleuts 1 -X 0 -M 1.0 -V 0.001 -B 1
RandomForest
Bagging with 100 iterations and base learner
waka.classifiers.trees.RandomTree -K 0 -M 1.0 -V 0.001 -B 1 -do-not-check-splittability=all the base classifiers;
```

[illegible]

Complex formulas (SVM)

The (near) future...

- Challenges at the State level (CA, US) - more data, possibly more concerns about disparities
- ‘Hot’ topics:
 - Ethics: disparities, fairness, stigmatization
 - Interpretability/Transparency, for a wide variety of audiences
 - Accountability
 - Engagement with, and high-level knowledge of, ML from the involved parties

Machine learning: the power and promise of computers that learn by example, Royal Society, April 2017.

Thanks.

Allegheny's Methodology Paper:

<http://www.alleghenycountyanalytics.us/index.php/2017/04/17/developing-predictive-risk-models-support-child-maltreatment-hotline-screening-decisions/>

Visit/follow CSDA:

<https://csda.aut.ac.nz/>

@AUTCSDA

@rvaithianathan