# Overcoming the Barriers to Production-Ready Machine Learning Workflows

Josh Bloom     Henrik Brink

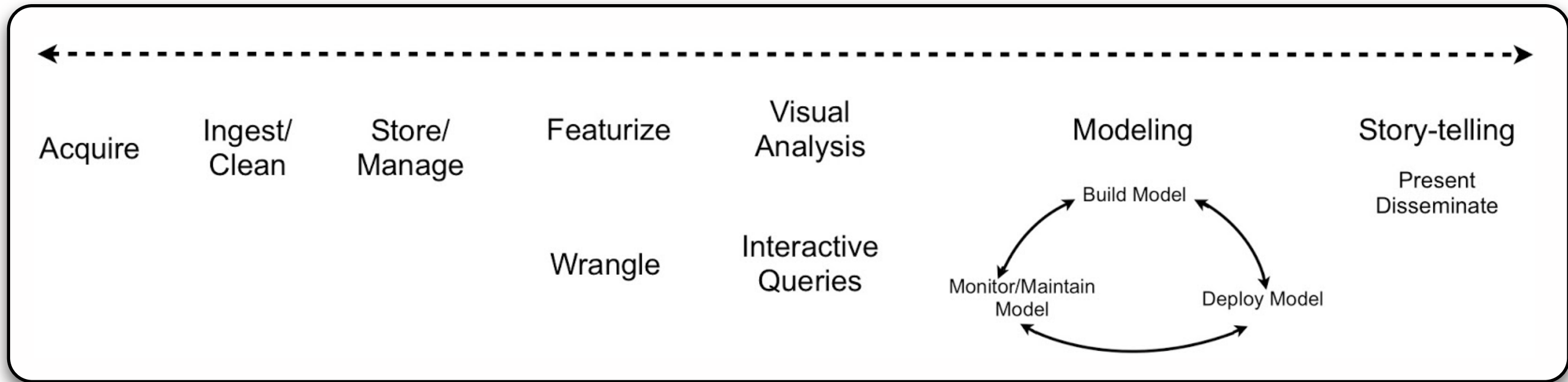@profjsb   @brinkar   @wiseio
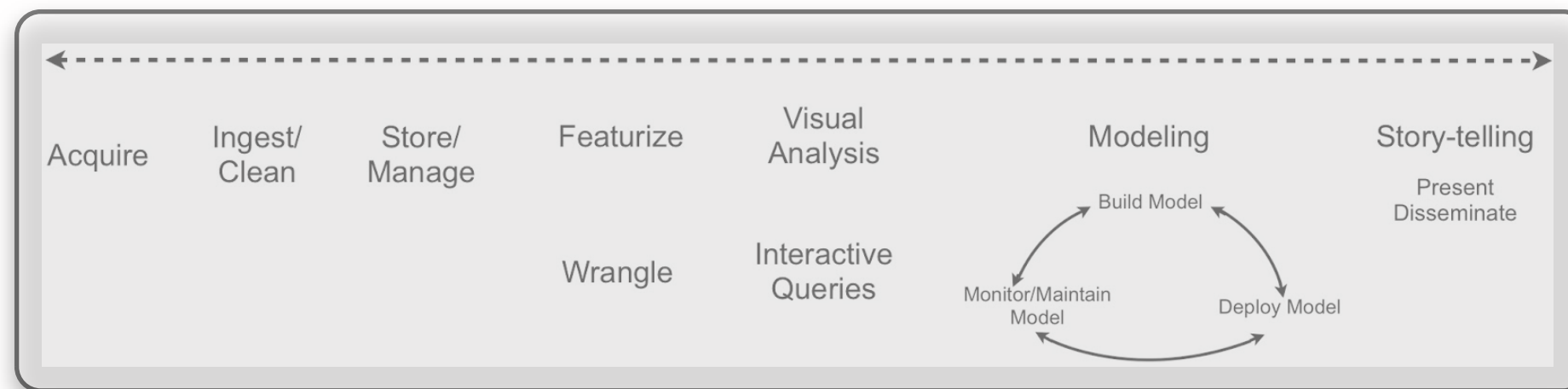
wise.io

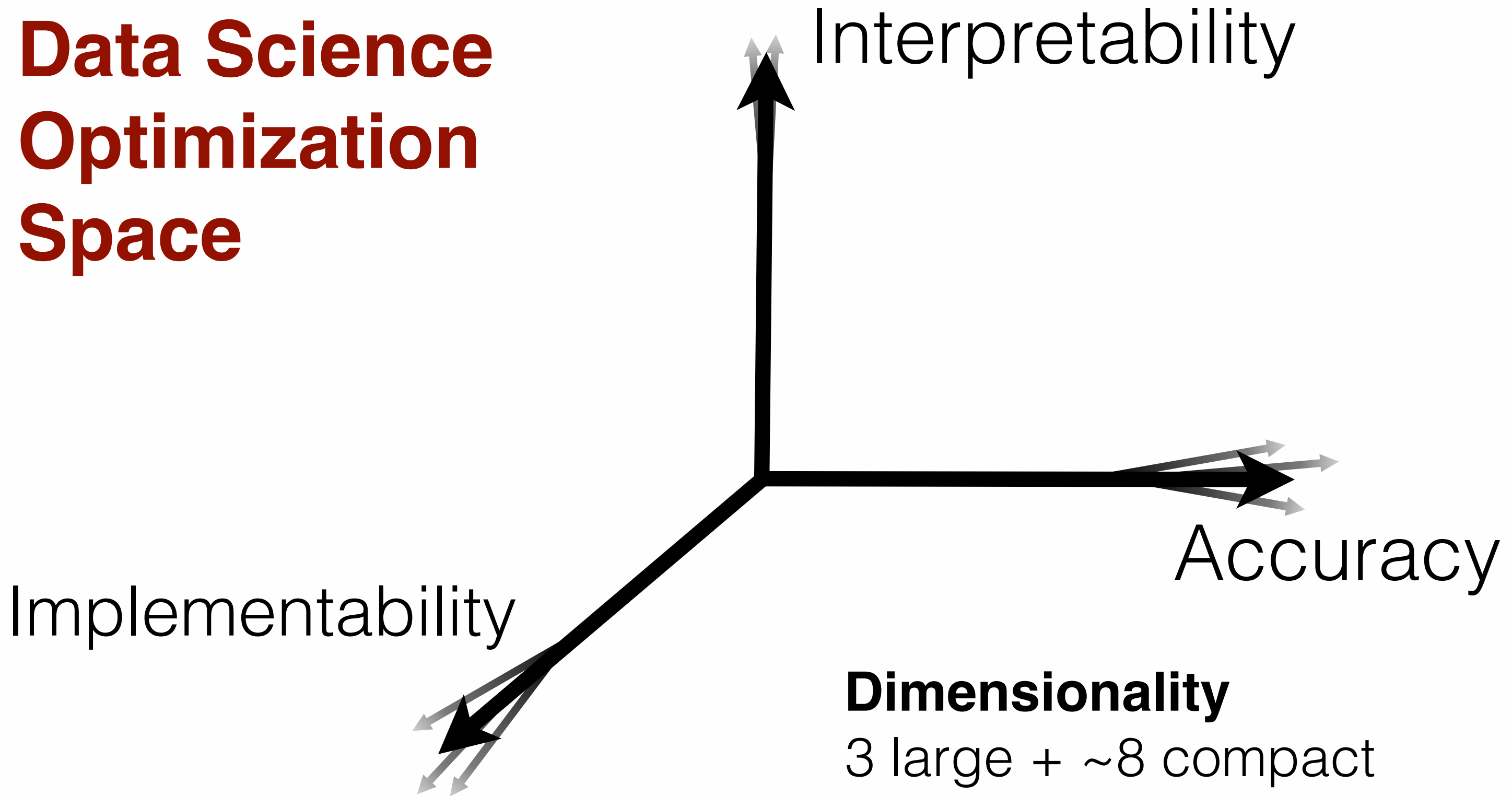University of California, Berkeley

Lorica's "Data Science Workflow"

# Real-World Data Science =
# *Optimization* over
# this *full* Workflow



Lorica's "Data Science Workflow"

# Our Background …
## "Data-Driven Scientists"

▸ Built & Deployed Real-time ML framework, discovering >10,000 events in > 10 TB of imaging
→ 50+ journal articles

▸ Built Probabilistic Event classification catalogs with innovative active learning

▸ Collective over 350 refereed journal articles including ML & timeseries analysis



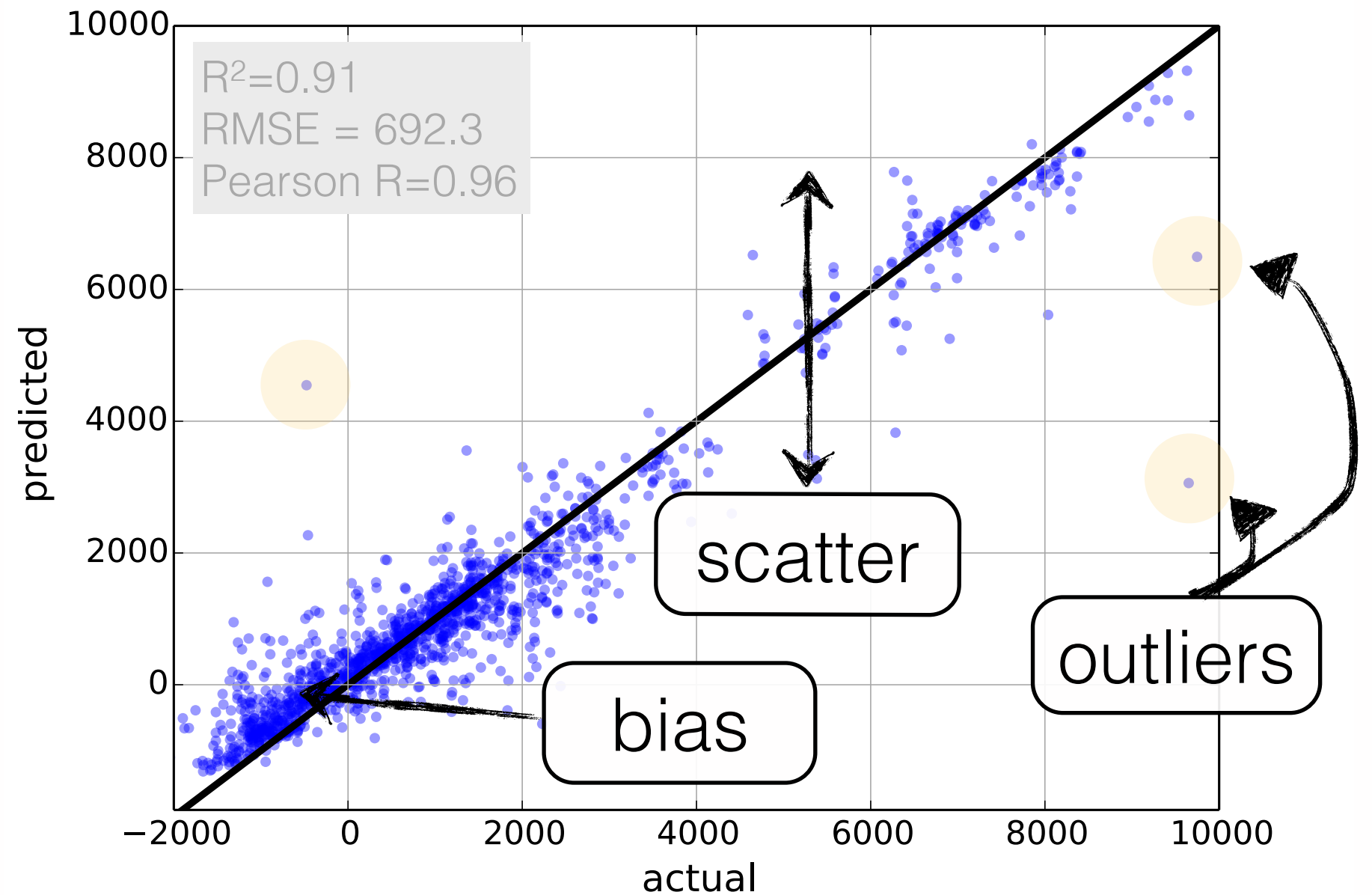*Our ML framework found the Nearest Supernova in 3 Decades ..*

SN 2011fe

# **Accuracy**

## Evaluation Metric: What's the essence of what I care about?

## Scalar proxies

- RMSE
- RMSLE
- [adjusted] $R^2$
- ...

cf. sklearn.metrics



$R^2$=0.91
RMSE = 692.3
Pearson R=0.96

scatter

bias

outliers

predicted

actual

# **Accuracy**

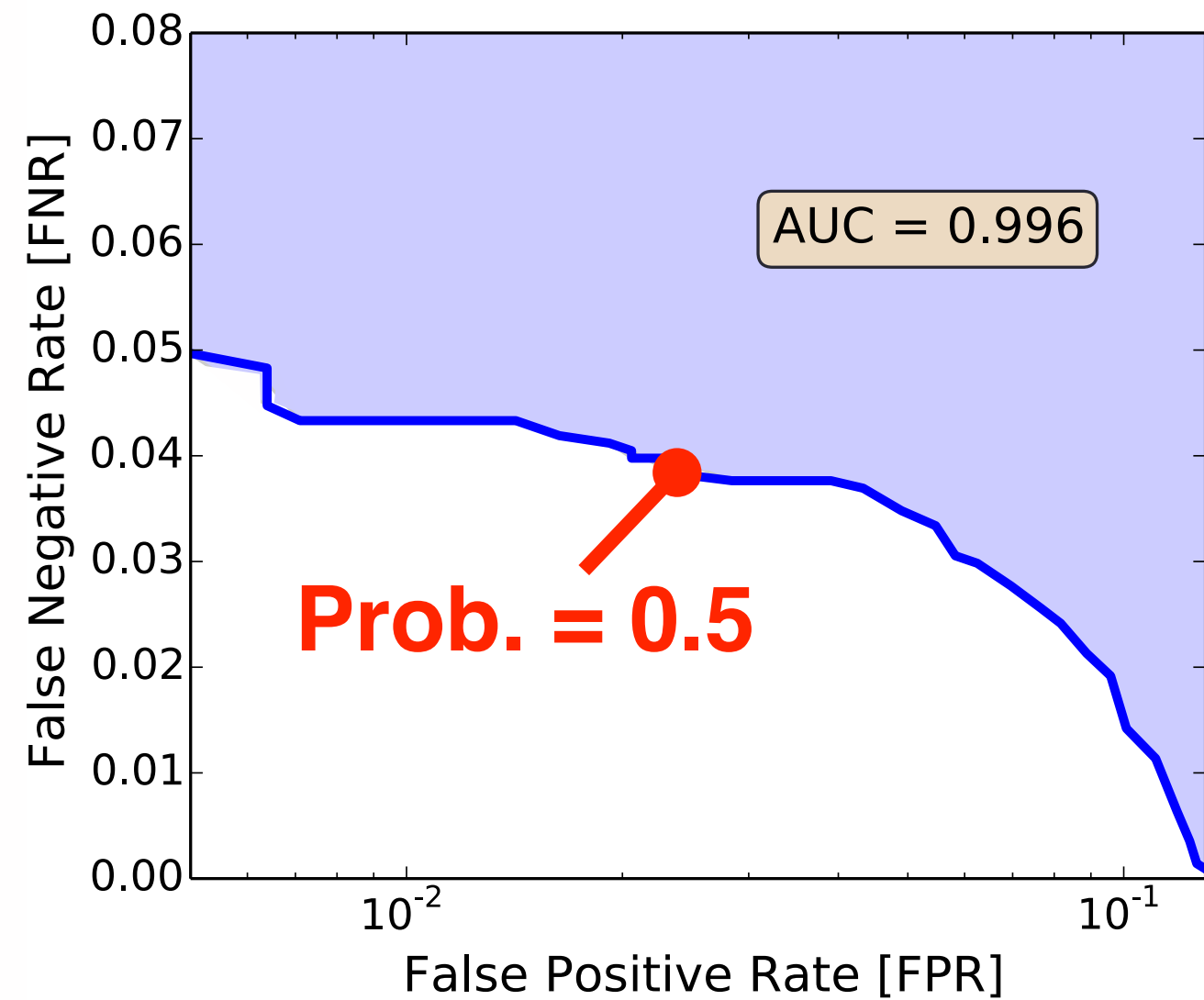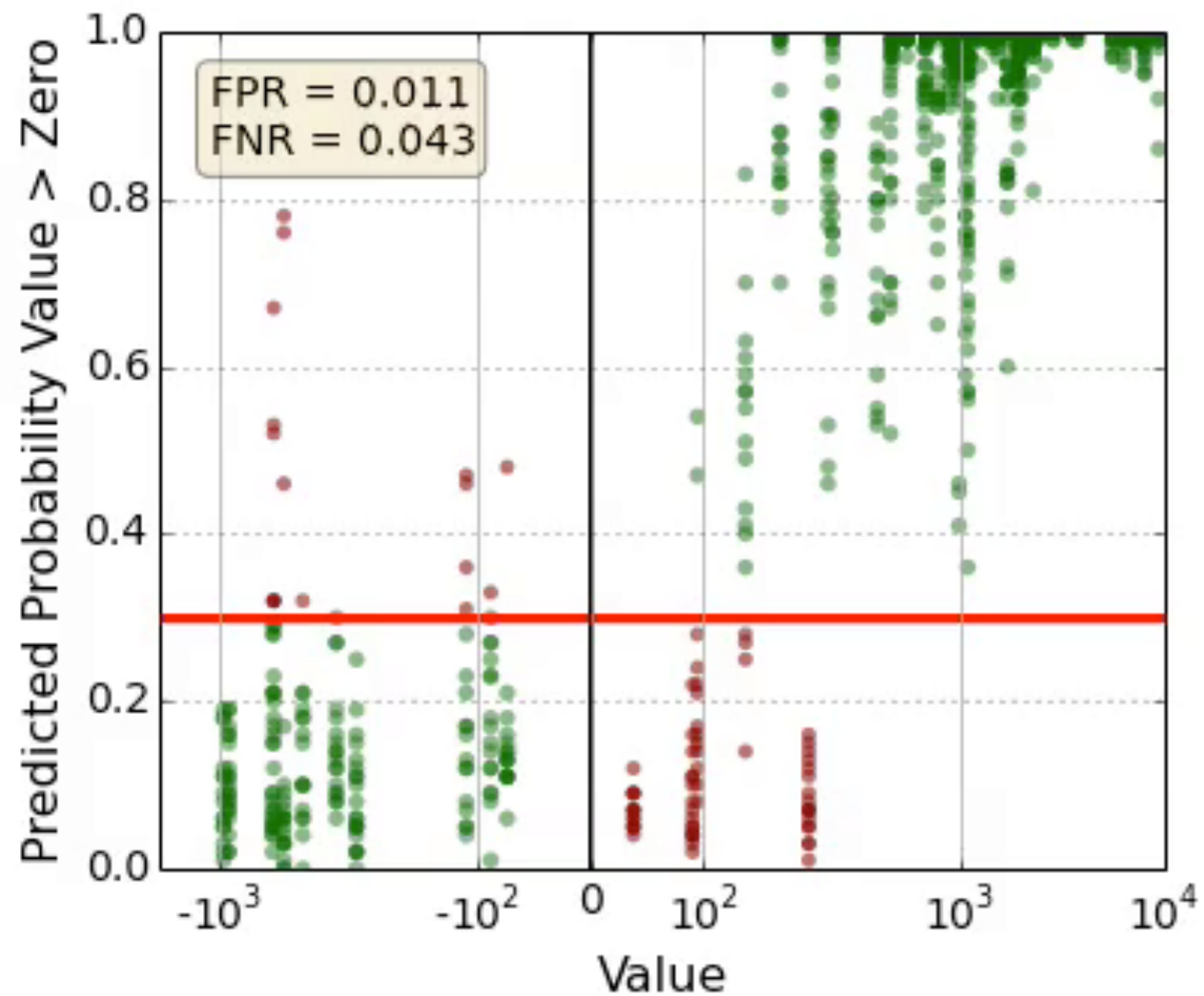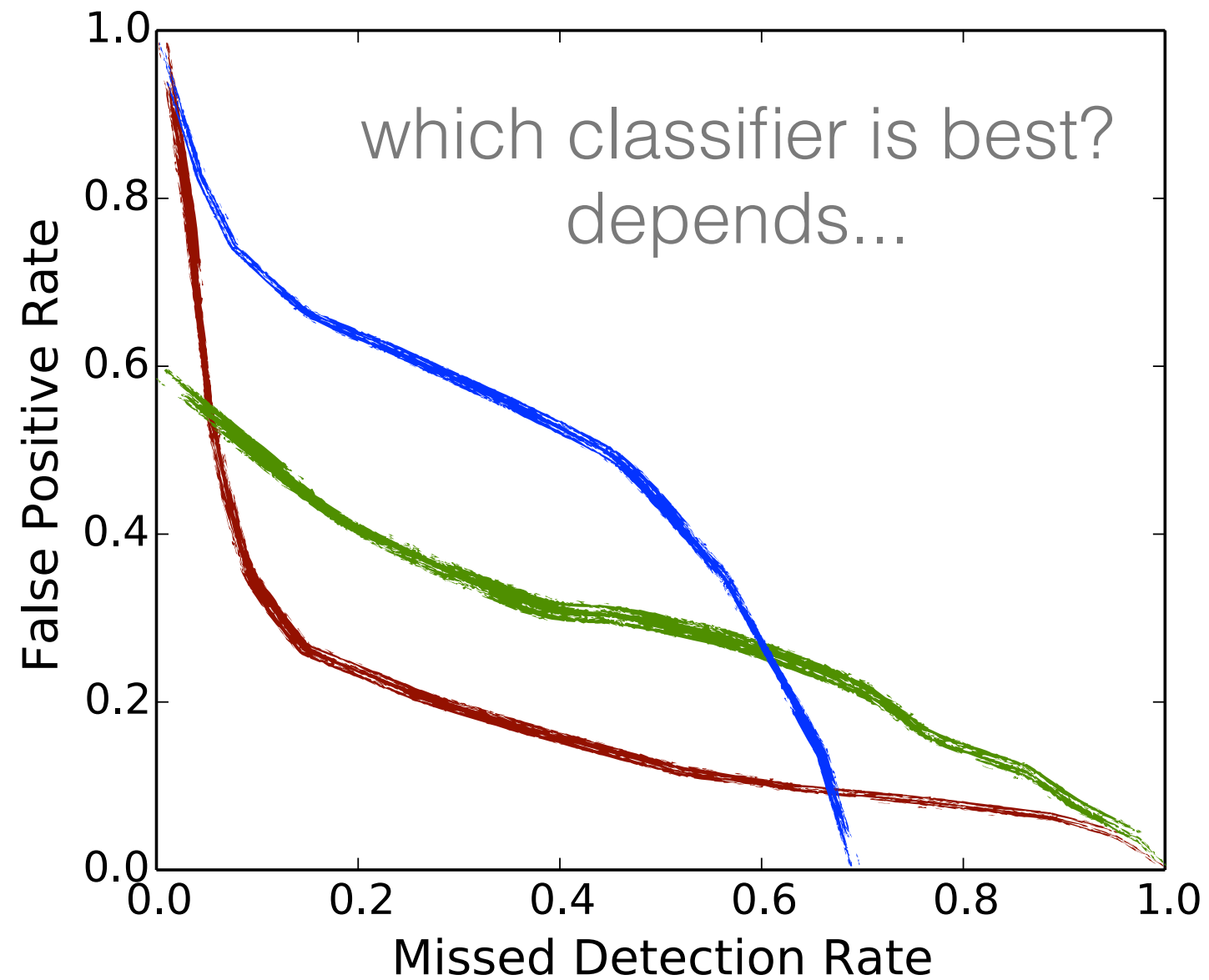**Accuracy**

Evaluation Metric: *What's the essence of what I care about?*

# **Accuracy**

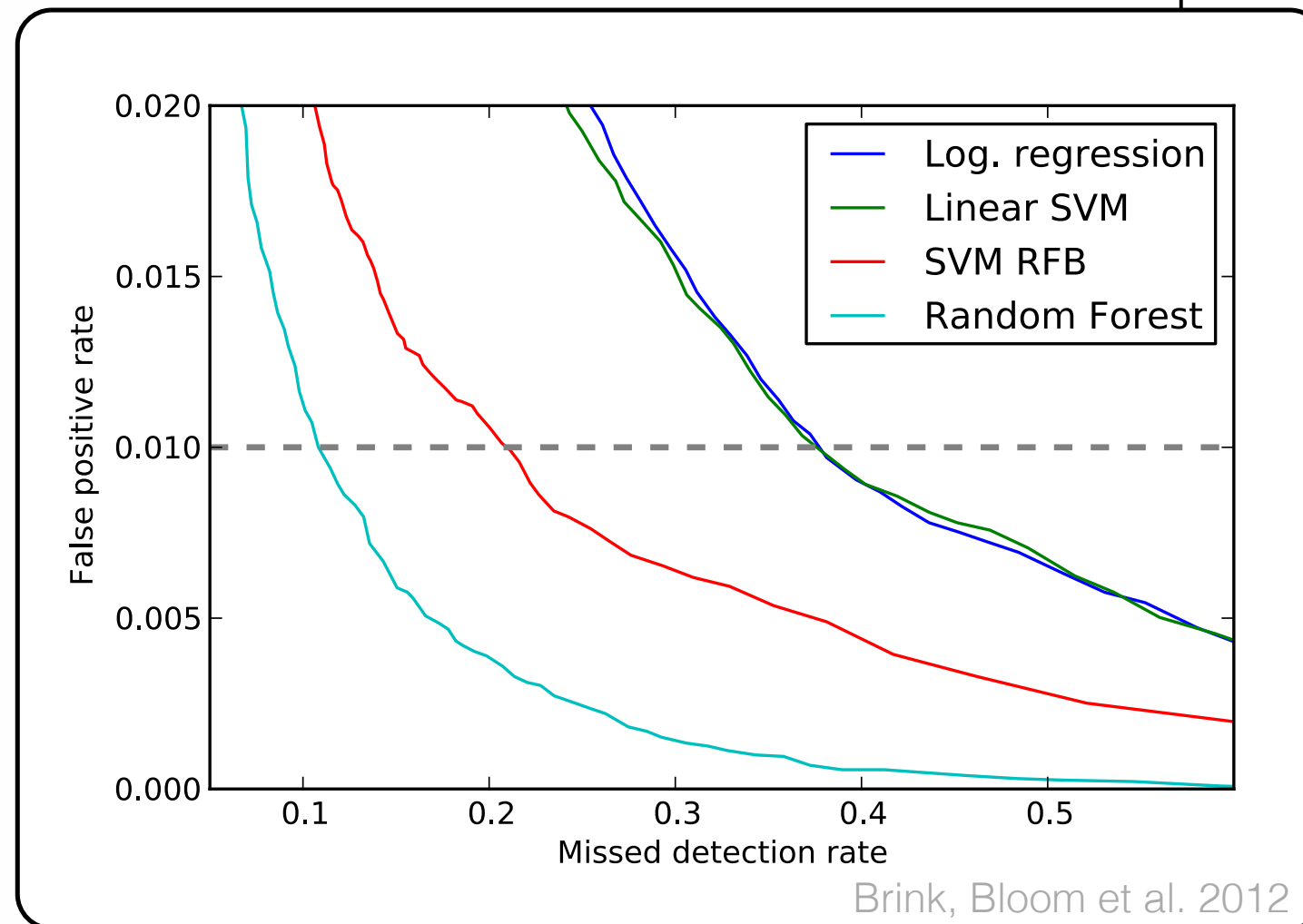# Evaluation Metric: *What's the essence of what I care about?*

42-dimensional feature space



Brink, Bloom et al. 2012

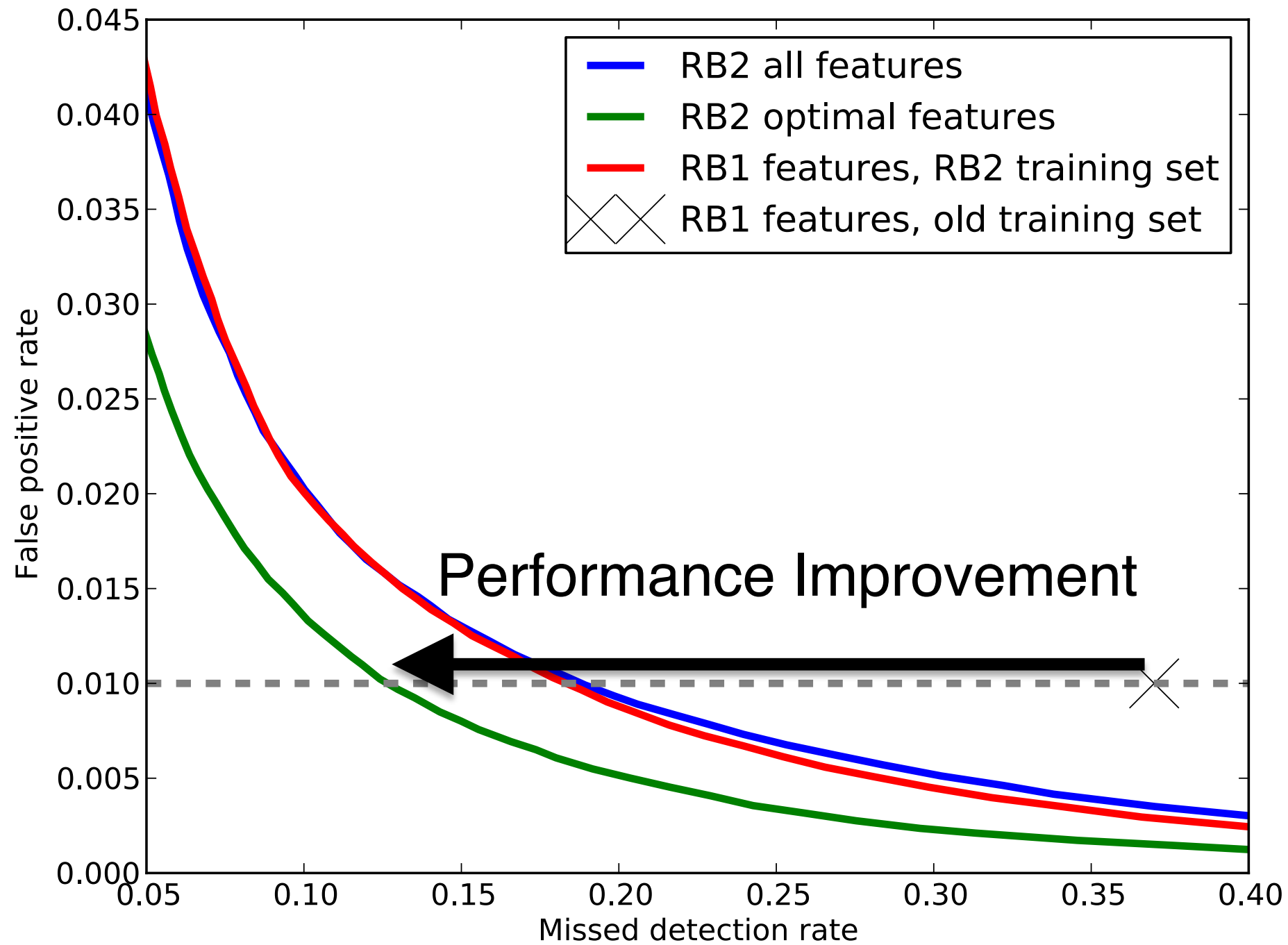Some ML algorithms just do *better*

# **Accuracy**

## More Data (Dimensions) is better, but Protect Against Curse of Dimensionality



Legend:
- RB2 all features
- RB2 optimal features
- RB1 features, RB2 training set
- RB1 features, old training set
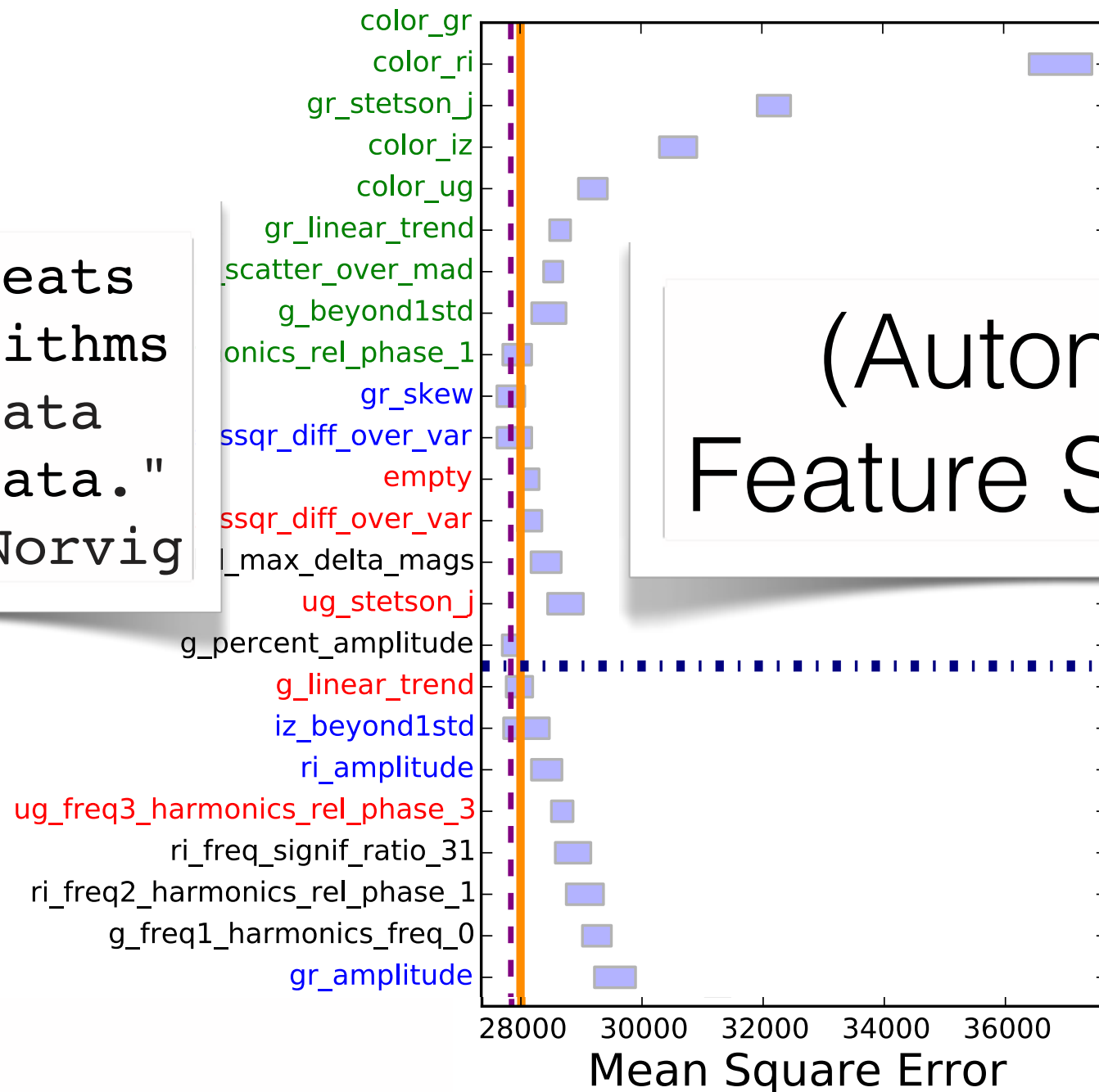
Performance Improvement

**Accuracy**

More Data (Dimensions) is better, but Protect Against Curse of Dimensionality

"More data beats clever algorithms but better data beats more data."
   - Peter Norvig

(Automatic) Feature Selection

**Accuracy**
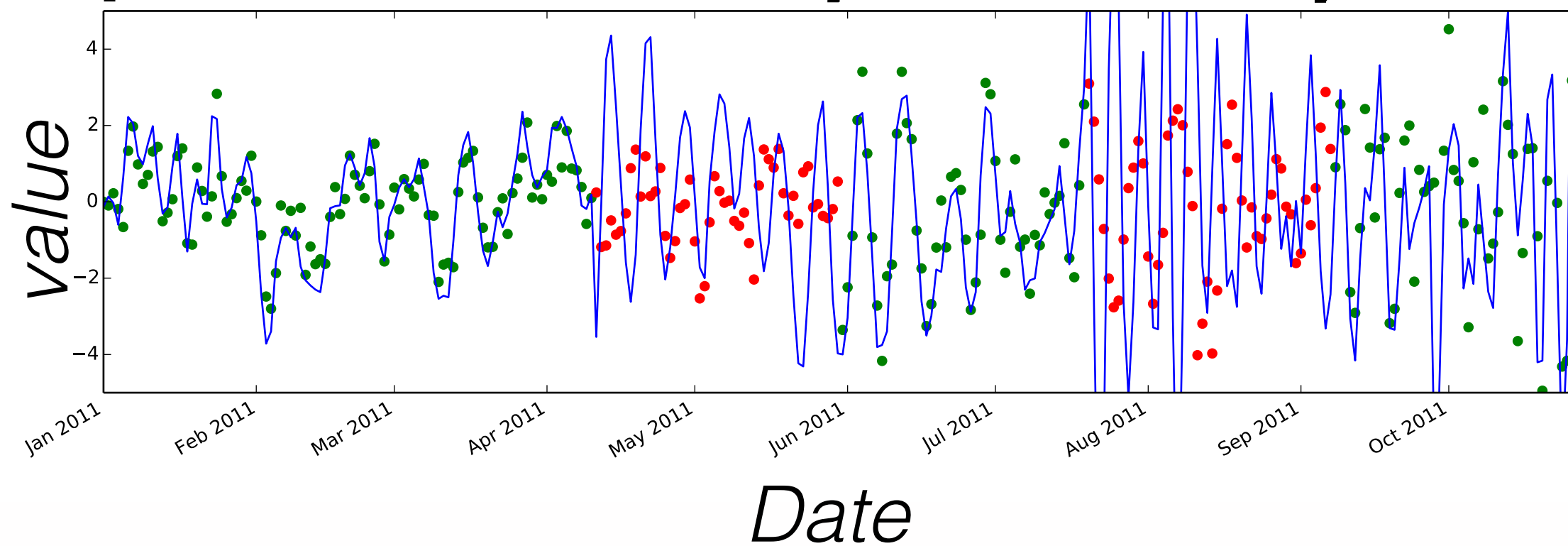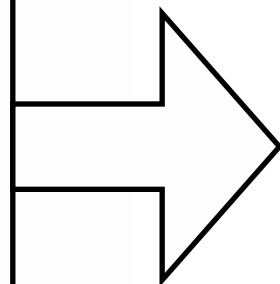
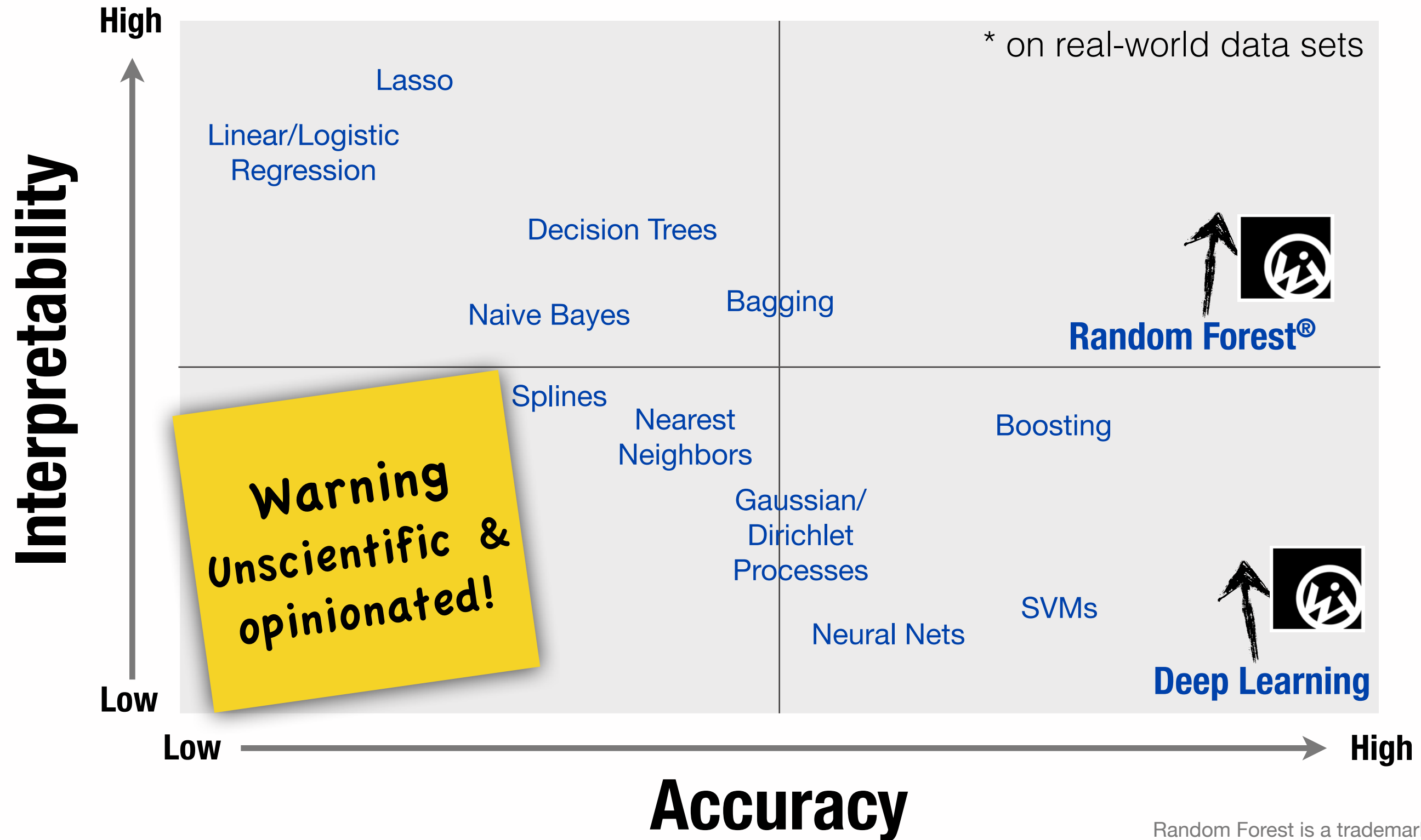Testing Set & Continuous (Streaming) Testing & Model Updates

Model # in production

1     2     3

model 1 building + validation on historical data

*value*

*Date*

— actual value     ● good prediction     ● "bad" prediction

# Interpretability

**Interpretability**

# How does the model work?

Consider a nonlinear system of equations:

$$\begin{cases} 3x_1 - \cos(x_2 x_3) - \frac{3}{2} = 0 \\ 4x_1^2 - 625x_2^2 + 2x_2 - 1 = 0 \\ \exp(-x_1 x_2) + 20x_3 + \frac{10\pi - 3}{3} \end{cases}$$

suppose we have the function

$$G(\mathbf{x}) = \begin{bmatrix} 3x_1 - \cos(x_2 x_3) \\ 4x_1^2 - 625x_2^2 + 2 \\ \exp(-x_1 x_2) + 20x_3 \end{bmatrix}$$

where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

and the objective function

$$F(\mathbf{x}) = \frac{1}{2}G^{\mathrm{T}}(\mathbf{x})G(\mathbf{x})$$
$$= \frac{1}{2}((3x_1 - \cos(x_2 x_3) - \frac{3}{2})^2 +$$

With initial guess

$$\mathbf{x}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

We know that

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \gamma_0 \nabla F(x^{(0)})$$

where

$$\nabla F(\mathbf{x}^{(0)}) = J_G(\mathbf{x}^{(0)})^{\mathrm{T}} G(\mathbf{x}^{(0)})$$

The Jacobian matrix $J_G(\mathbf{x}^{(0)})$

$$J_G = \begin{bmatrix} 3 & \sin(x_2 x_3)x_3 & \sin(x_2 x_3)x_2 \\ 8x_1 & -1250x_2 + 2 & 0 \\ -x_2 \exp(-x_1 x_2) & -x_1 \exp(-x_1 x_2) & 20 \end{bmatrix}$$

Then evaluating these terms at $\mathbf{x}^{(0)}$

$$J_G(\mathbf{x}^{(0)}) = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 20 \end{bmatrix}$$

and

$$G(\mathbf{x}^{(0)}) = \begin{bmatrix} -2.5 \\ -1 \\ 10.472 \end{bmatrix}$$

So that

$$\mathbf{x}^{(1)} = 0 - \gamma_0 \begin{bmatrix} -7.5 \\ -2 \\ 209.44 \end{bmatrix}.$$

and

$$F(\mathbf{x}^{(0)}) = 0.5((-2.5)^2 + (-1)^2 + (10.472)^2) = 58.456$$

Now a suitable $\gamma_0$ must be found such that $F(\mathbf{x}^{(1)}) \leq F(\mathbf{x}^{(0)})$. This can be done with algorithms. One might also simply guess $\gamma_0 = 0.001$ which gives

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0.0075 \\ 0.002 \\ -0.20944 \end{bmatrix}$$

evaluating at this value,

$$F(\mathbf{x}^{(1)}) = 0.5((-2.48)^2 + (-1.00)^2 + (6.28)^2) = 23.306$$

**Interpretability**

How does the model work?

Consider a nonlinear system of equations:

$$\begin{cases} 3x_1 - \cos(x_2x_3) - \frac{3}{2} = 0 \\ 4x_1^2 - 625x_2^2 + 2x_2 - 1 = 0 \\ \exp(-x_1x_2) + 20x_3 + \frac{10\pi-3}{3} \end{cases}$$

suppose we have the function

$$G(\mathbf{x}) = \begin{bmatrix} 3x_1 - \cos(x_2x_3) \\ 4x_1^2 - 625x_2^2 + 2x \\ \exp(-x_1x_2) + 20x_3 \end{bmatrix}$$

where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

and the objective function

$$F(\mathbf{x}) = \frac{1}{2}G^{\mathrm{T}}(\mathbf{x})G(\mathbf{x})$$
$$= \frac{1}{2}((3x_1 - \cos(x_2x_3) - \frac{3}{2})^2 + $$

With initial guess

$$\mathbf{x}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

We know that

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \gamma_0 \nabla F(x^{(0)})$$

where

$$\nabla F(\mathbf{x}^{(0)}) = J_G(\mathbf{x}^{(0)})^{\mathrm{T}} G(\mathbf{x}^{(0)})$$

The Jacobian matrix $J_G(\mathbf{x}^{(0)})$

$$J_G = \begin{bmatrix} 3 & \sin(x_2x_3)x_3 & \sin(x_2x_3)x_2 \\ 8x_1 & -1250x_2 + 2 & 0 \\ -x_2 \exp(-x_1x_2) & -x_1 \exp(-x_1x_2) & 20 \end{bmatrix}$$

Then evaluating these terms at $\mathbf{x}^{(0)}$

$$J_G(\mathbf{x}^{(0)}) = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 20 \end{bmatrix}$$

and

$$G(\mathbf{x}^{(0)}) = \begin{bmatrix} -2.5 \\ -1 \\ 10.472 \end{bmatrix}$$

So that

$$\mathbf{x}^{(1)} = 0 - \gamma_0 \begin{bmatrix} -7.5 \\ -2 \\ 209.44 \end{bmatrix}.$$

and

$$F(\mathbf{x}^{(0)}) = 0.5((-2.5)^2 + (-1)^2 + (10.472)^2) = 58.456$$

Now a suitable $\gamma_0$ must be found such that $F(\mathbf{x}^{(1)}) \leq F(\mathbf{x}^{(0)})$. This can be done with

algorithms. One might also simply guess $\gamma_0 = 0.001$ which gives

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0.0075 \\ 0.002 \\ -0.20944 \end{bmatrix}$$

evaluating at this value,

$$F(\mathbf{x}^{(1)}) = 0.5((-2.48)^2 + (-1.00)^2 + (6.28)^2) = 23.306$$

**Interpretability**

# Why do I get these answers?

e.g., Credit score

## Sample FICO® Scoring Model

| Category | Characteristic | Attributes | Points |
|---|---|---|---|
| Payment History | Number of months since the most recent derogatory public record | No public record<br>0 – 5<br>6 – 11<br>12 – 23<br>24+ | 75<br>10<br>15<br>25<br>55 |
| Outstanding Debt | Average balance on revolving trades | No revolving trades<br>0<br>1 – 99<br>100 – 499<br>500 – 749<br>750 – 999<br>1000 or more | 30<br>55<br>65<br>50<br>40<br>25<br>15 |
| Credit History Length | Number of months in file | Below 12<br>12 – 23<br>24 – 47<br>48 or more | 12<br>35<br>60<br>75 |
| Pursuit of New Credit | Number of inquiries in last 6 mos. | 0<br>1<br>2<br>3<br>4+ | 70<br>60<br>45<br>25<br>20 |
| Credit Mix | Number of bankcard trade lines | 0<br>1<br>2<br>3<br>4+ | 15<br>25<br>50<br>60<br>50 |

# **Peering Inside** the Black Box



| Feature | Importance |
|---|---|
| over_draft :'no checking' | |
| over_draft :'<0' | |
| credit_usage | |
| current_balance | |
| cc_age | |
| Average_Credit_Balance :'<100' | |
| credit_history :'critical/other existing credit' | |

Random Forest® model-level **feature importance**

**Interpretability**

# **Peering Inside** the Black Box

Individual-level prediction feature importance

Probability of Default in 1 year:
**76% [deny loan]**

**Driving factors**

☼ Credit history: 10 months | 14%

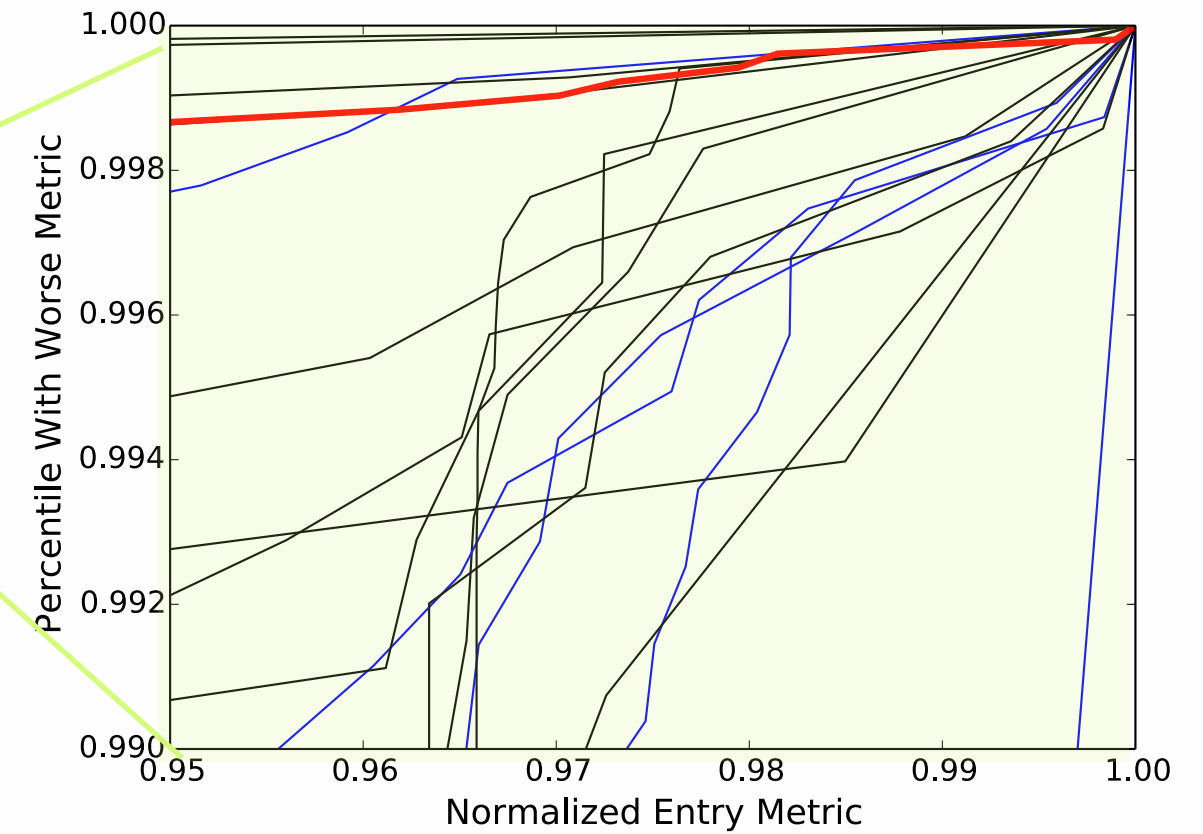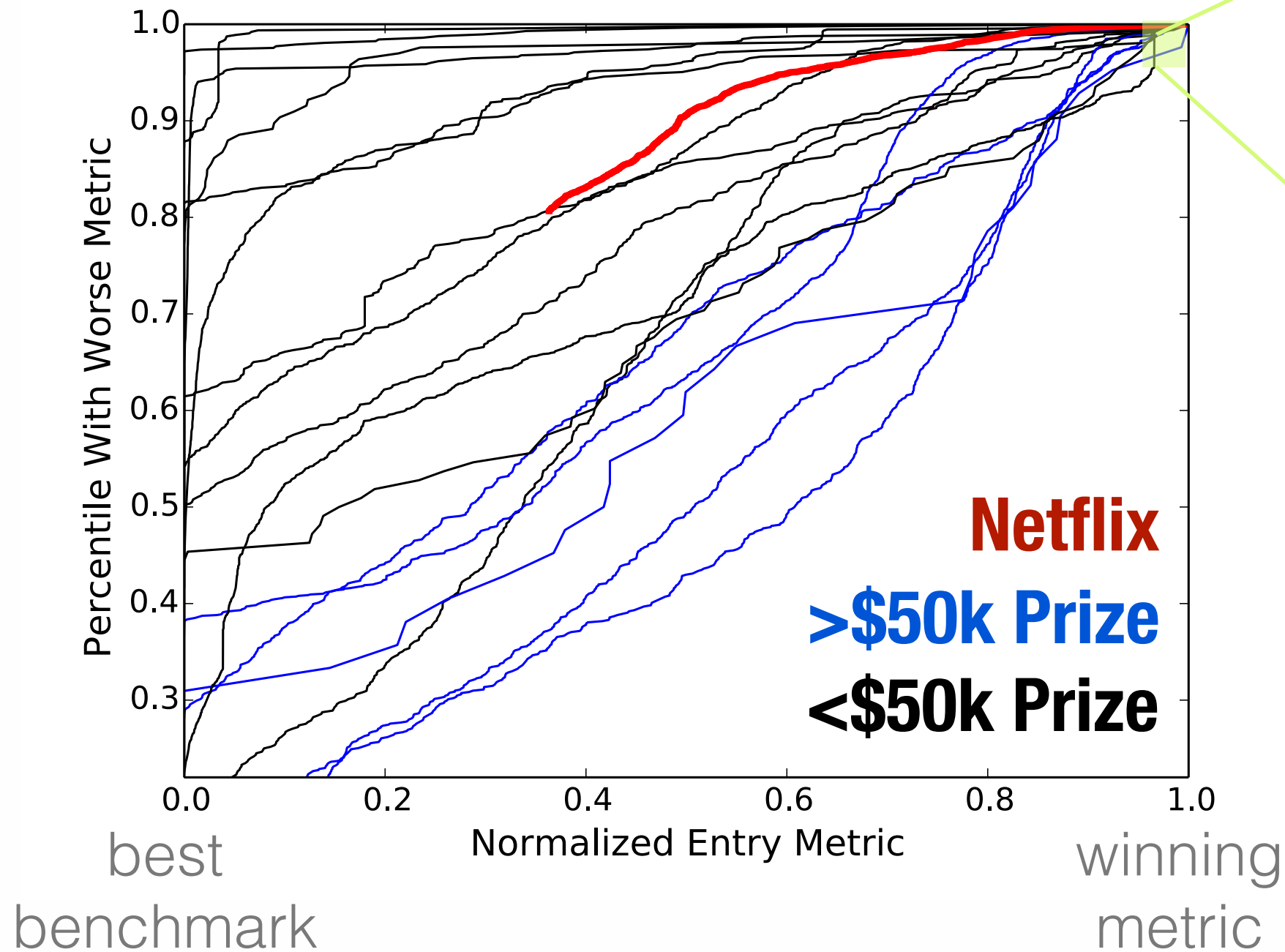☼ Outstanding debt: $1200 | 5%

☼ Inquiries in 6 months: 2 | 1%

e.g. microcredit application scorecard

# Implementability

How long does it take to put the model into production? At what cost?

# Implementability



Netflix

>$50k Prize

<$50k Prize

best benchmark

winning metric

many teams get within ~few % of optimum

**so which is easier to put into production?**

18

# Implementability

On the **NETFLIX** Prize

"We evaluated some of the new methods offline but the **additional accuracy gains** that we measured **did not seem to justify the engineering effort** needed to bring them into a **production environment.**"

*Xavier Amatriain and Justin Basilico (April 2012)*

**Implementability**
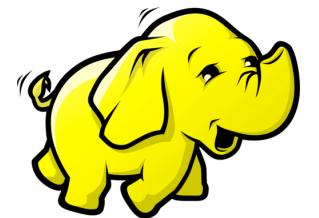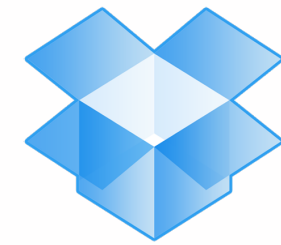
The divide
between
**data science**
& **production**

**Implementability**

Treat Machine Learning Deployment as you would Software

▸ Continuous Deployment
▸ RESTful API
▸ Language bindings
▸ Security
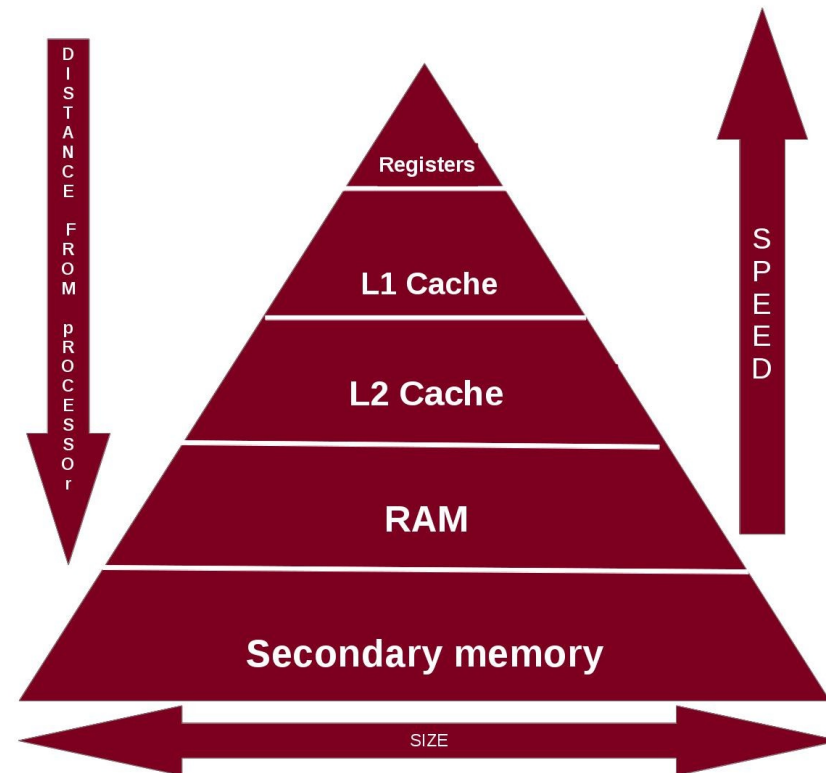▸ SLA

# Integration
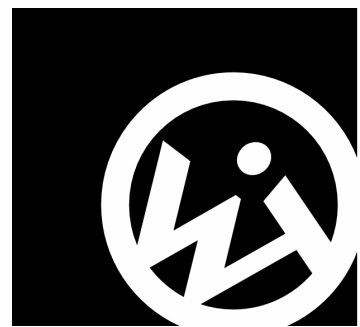
Connect data

Consume predictions

# Scalability & Speed

**Implementability**

**Micro-scaling**

Fast, efficient use of memory hierarchy

**Horizontally** scalable data processing

# We are Hiring!

▸ Full-stack developers

  ▸ Javascript, Python, Spark/Shark

▸ Front end developers

▸ DevOps engineers

▸ C++ engineers

  ▸ C++ template metaprogramming

▸ Data scientists

  ▸ Python, Deep NN, ML expertise

wise.io

jobs@wise.io

http://wise.io/jobs/