# Boosting for Probability Estimation & Cost-Sensitive Learning

Nikos Nikolaou

EPSRC Doctoral Prize Fellow, University of Manchester

# Introduction:
# Supervised Machine Learning Basics
# (Classification)

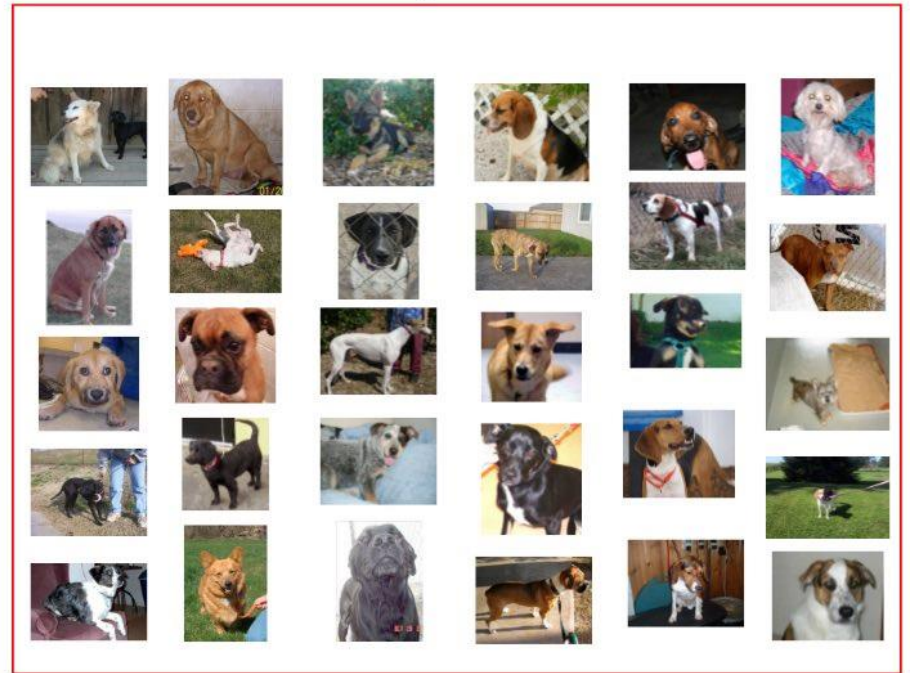# Classification

## Example: Cat vs. Dog classifier

# Machine learning basics: training

**TRAINING DATA**

Features
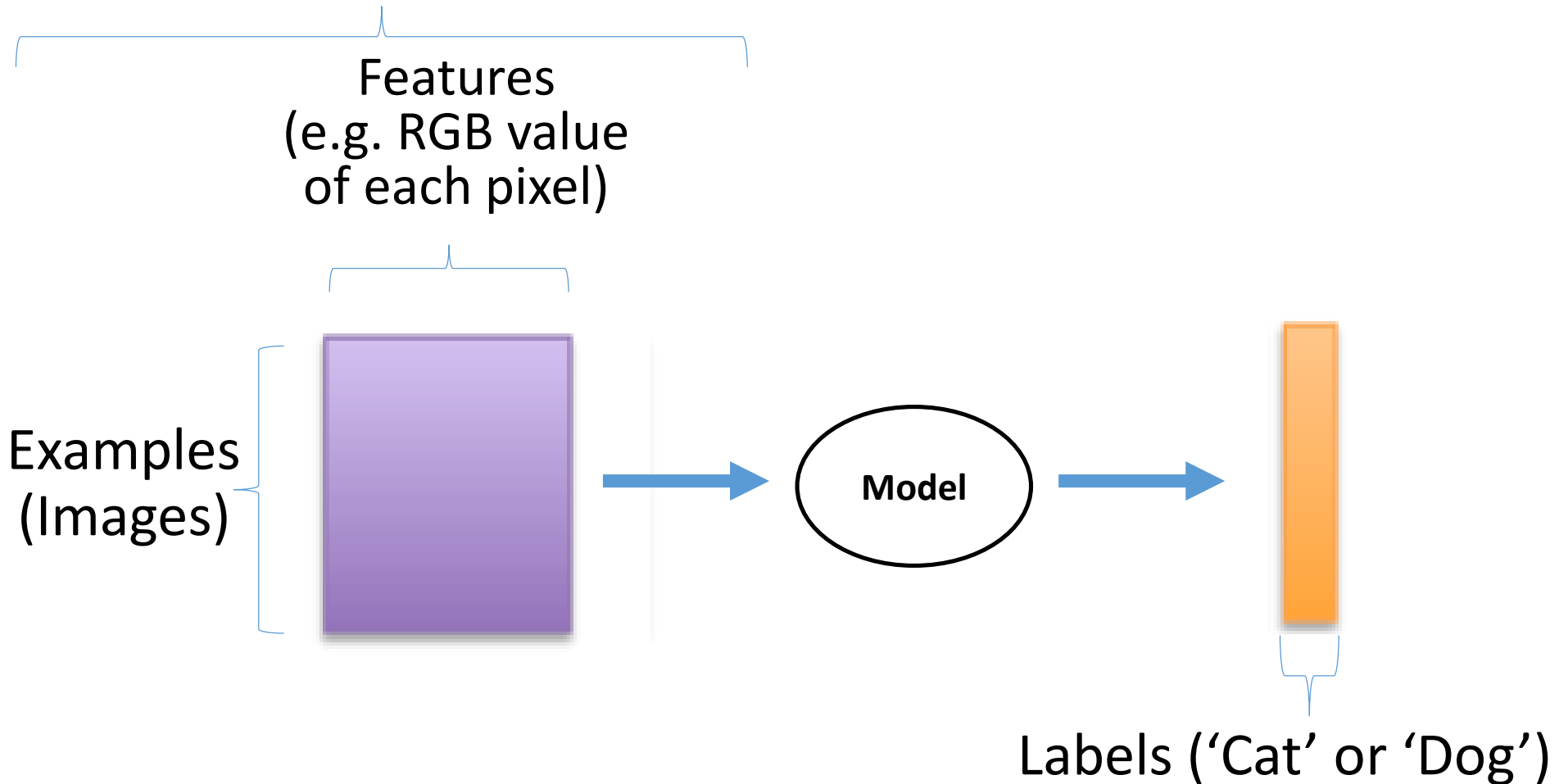(e.g. RGB value
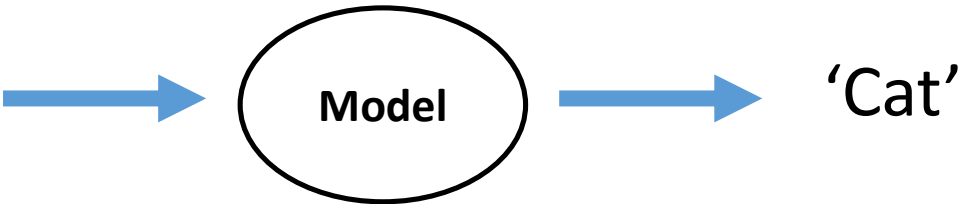of each pixel)

Examples
(Images)



LEARNING
ALGORITHM

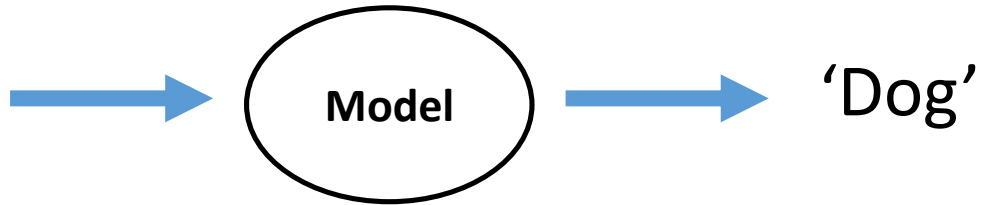Model

Labels ('Cat' or 'Dog')

# Machine learning basics: prediction

**TEST DATA**

Features
(e.g. RGB value
of each pixel)

Examples
(Images)

**Model**
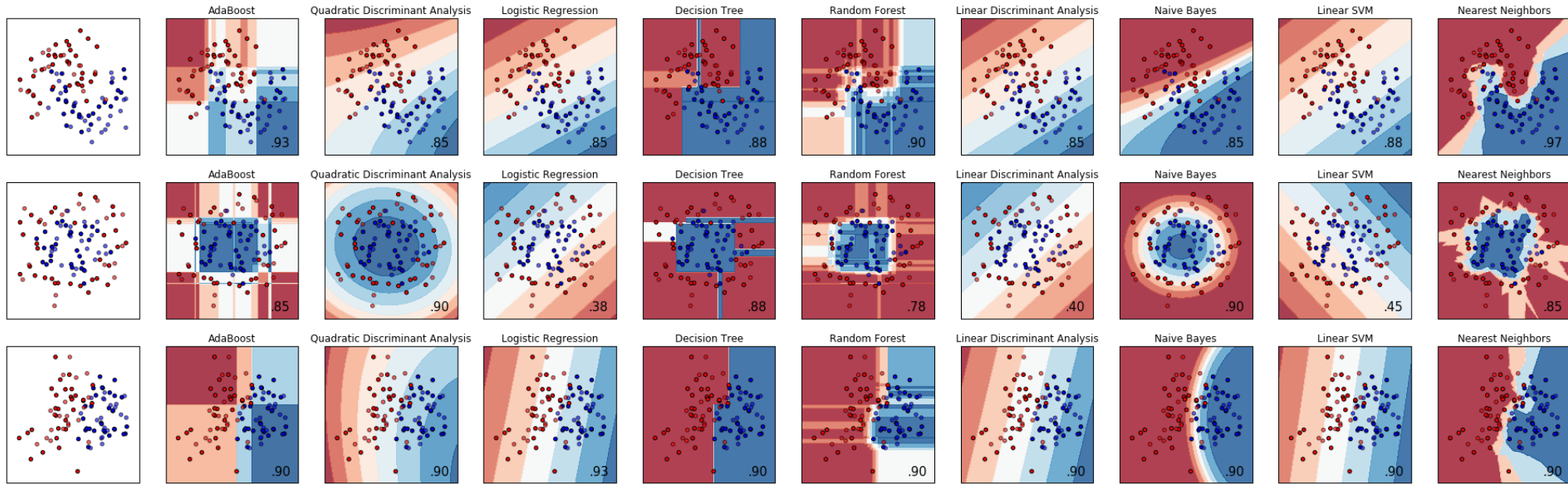
Labels ('Cat' or 'Dog')

# Machine learning basics: prediction

# The learning algorithm's job

Given a set of points in some space belonging to different classes...



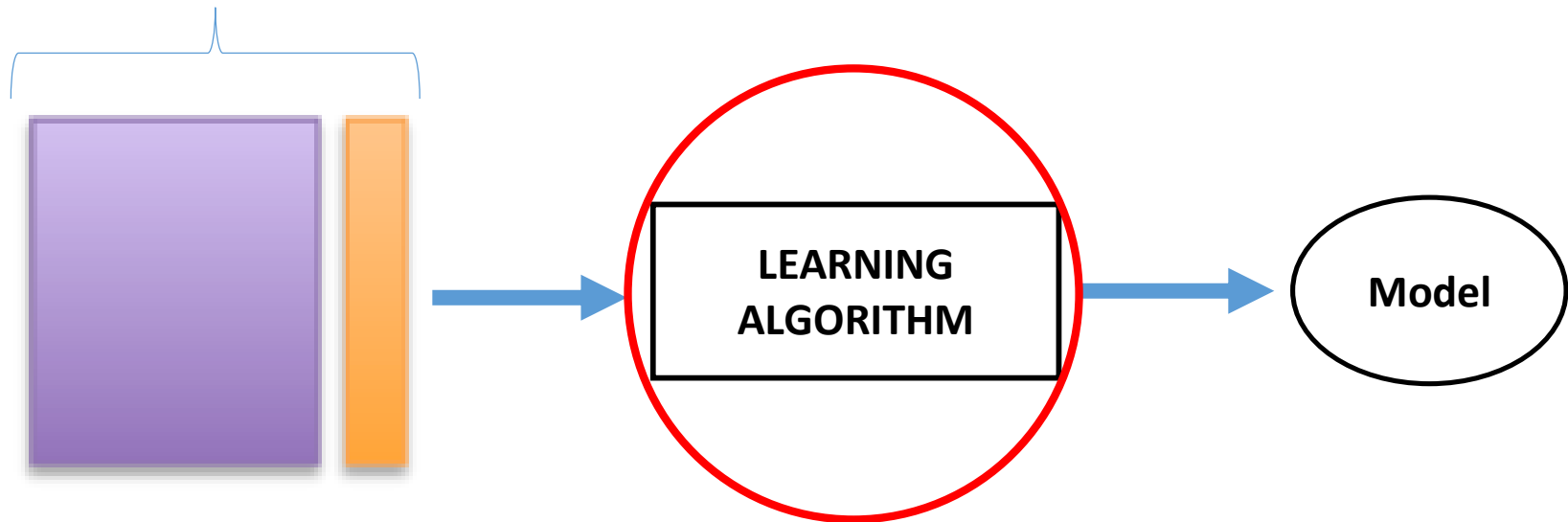Learn a **decision surface** that **'best' separates classes**

Many **learning algorithms** each with its own **assumptions** (statistical, probabilistic, mathematical, geometrical, ...)

# In this talk…

**TRAINING DATA**



We focus on **BOOSTING**, a specific **family** of **learning algorithms**

**Meta-learning algorithms** - can **apply to other learning algorithms** improving their performance

# More specifically...

**BOOSTING** in **cost-sensitive** scenarios



| | | Actual Value | |
|---|---|---|---|
| | | positives | negatives |
| **Predicted Value** | positives | **TP** True Positive | **FP** False Positive |
| | negatives | **FN** False Negative | **TN** True Negative |

cost of a FP ≠ cost of a FN

# Part I:
# What is wrong with cost-sensitive Boosting?

# Boosting

Can we turn a **weak learner** into a **strong learner**? (Kearns, 1988)

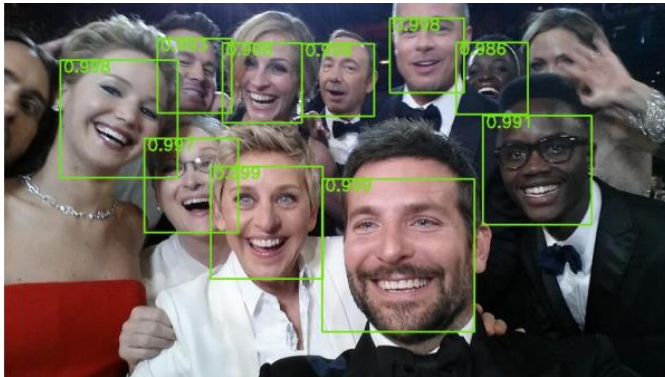| Marginally more accurate than random guessing |
|---|

| Arbitrarily high accuracy |
|---|



**YES!** '**Hypothesis Boosting**' (Schapire, 1990)

**AdaBoost** (Freund & Schapire, 1997)

**Gödel Prize 2003**

**Gradient Boosting** (Friedman, 1999; Mason et al., 1999)

# Boosting

Very **successful** in comparisons, applications & competitions



Rich **theoretical depth**:

PAC learning, VC theory, margin theory, optimization, decision theory, game theory, probabilistic modelling, information theory, dynamical systems, ...

# Adaboost (Freund & Schapire 1997)
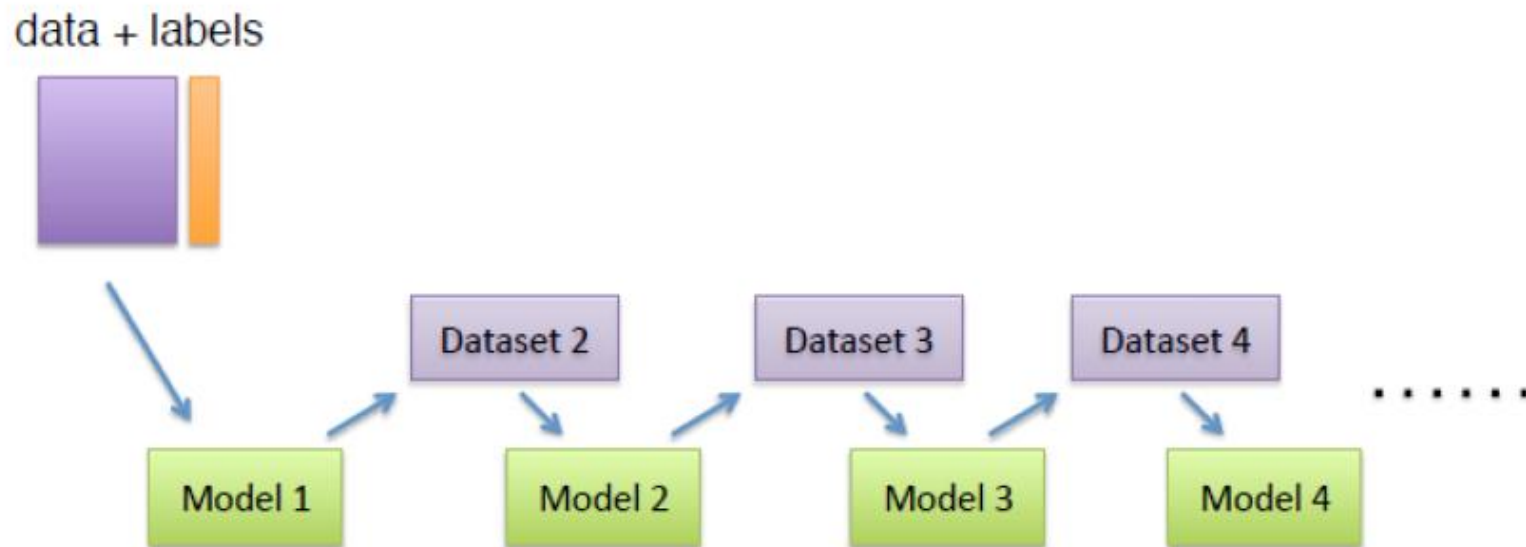
**Ensemble** method.

Train models **sequentially**.

Each model **focuses on examples previously misclassified**.

Combine by **weighted majority vote**.
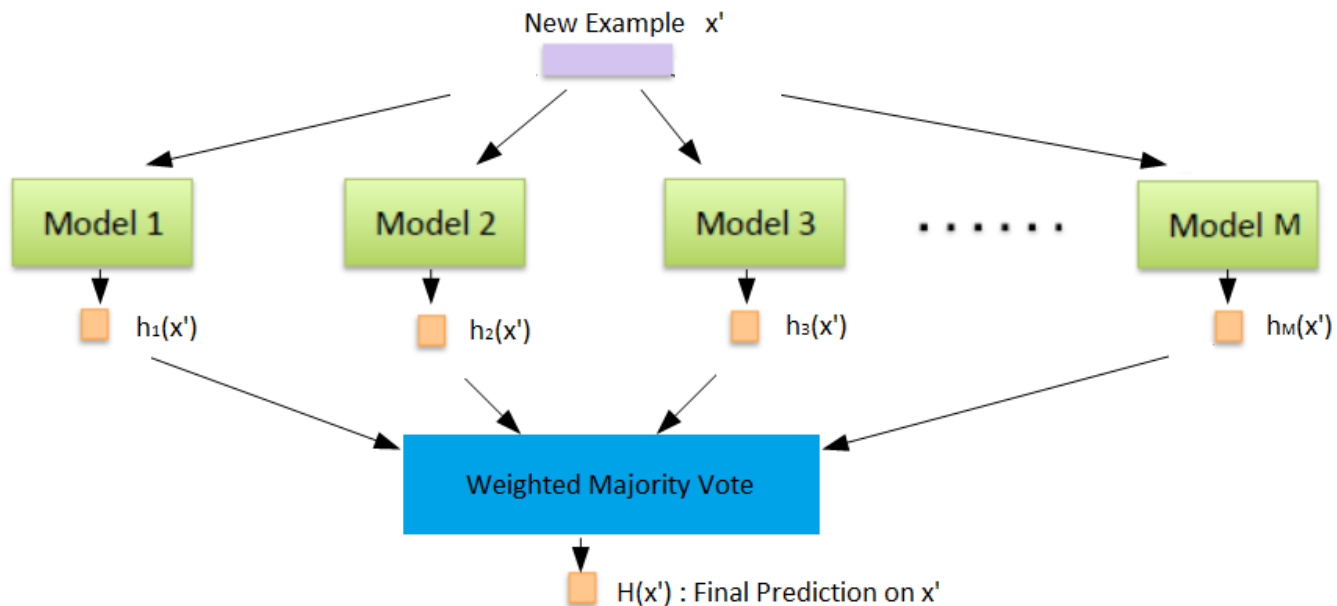
# AdaBoost: training

Construct strong model **sequentially** by combining multiple weak models



Each model reweights/resamples the data, emphasizing on the examples the previous one misclassified – i.e. each model focuses on **correcting the mistakes of the previous one**

# AdaBoost: predictions

Prediction: **weighted majority vote** among M weak learners

# AdaBoost: algorithm

Define a distribution over the training set, $D_1(i) = \frac{1}{N}$, $\forall i$. — **Initial weight distribution**

**for** $t = 1$ to T **do**

    Build a classifier $h_t$ from the training set, using distribution $D_t$.

    Set $\alpha_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$ ——— **Majority voting confidence in classifier t**

    Update $D_{t+1}$ from $D_t$ :

        Set $D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$ ——— **Distribution update**

**end for**

$$H(x') = sign\left( \sum_{t=1}^{T} \alpha_t h_t(x') \right)$$ ——— **Majority vote on test example x'**
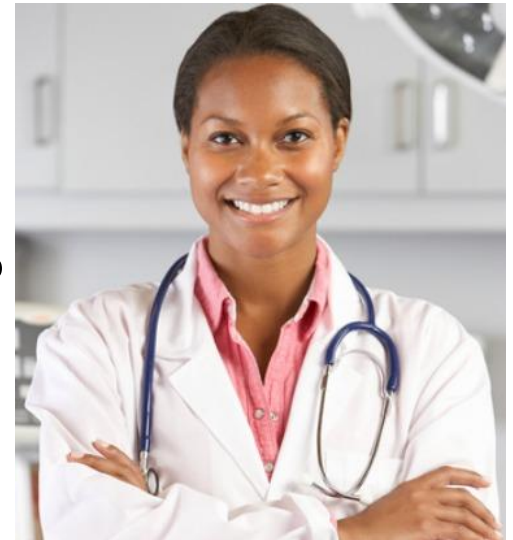
[Pos & Neg class encoded as +1 & -1 respectively for both predictions $h_t(\boldsymbol{x})$ and labels $y$]

# Adaboost

How will it work on cost sensitive* problems?    $\begin{bmatrix} 0 & c_{FN} \\ c_{FP} & 0 \end{bmatrix}$

i.e. with differing cost for a False Positive / False Negative ...

...does it **minimize** the **expected cost** (a.k.a. **risk**)?

*note: cost-sensitive & imbalanced class learning **duality**

# Cost sensitive Adaboost…

AdaBoost (Freund & Schapire 1997)
AdaCost (Fan et al. 1999)
AdaCost($\beta_2$) (Ting 2000)
CSB0 (Ting 1998)
CSB1 (Ting 2000)
CSB2 (Ting 2000)
AdaC1 (Sun et al. 2005, 2007)
AdaC2 (Sun et al. 2005, 2007)
AdaC3 (Sun et al. 2005, 2007)
CSAda (Mashnadi-Shirazi & Vasconselos 2007, 2011)
AdaDB (Landesa-Vázquez & Alba-Castro 2013)
AdaMEC (Ting 2000, Nikolaou & Brown 2015)
CGAda (Landesa-Vázquez & Alba-Castro 2012, 2015)
AsymAda (Viola & Jones 2002)

**15+** boosting variants over **20** years

Some **re-invented** multiple times

Most proposed as **heuristic** modifications to original AdaBoost

Many treat FP/FN costs as **hyperparameters**

# A step back… Why is Adaboost interesting?

*Functional Gradient Descent* (Mason et al., 2000)

*Decision Theory* (Freund & Schapire, 1997)

*Margin Theory* (Schapire et al., 1998)

*Probabilistic Modelling* (Lebanon & Lafferty 2001; Edakunni et al 2011)

Set $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$

Update $D_{t+1}$ from $D_t$ :

Set $D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$
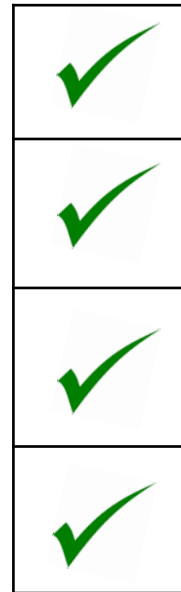
# So for a cost sensitive boosting algorithm...

*Functional Gradient Descent*

*Decision Theory*

*Margin Theory*

*Probabilistic Modelling*

*"Does the algorithm follow from each?"*

Set $\alpha_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$

Update $D_{t+1}$ from $D_t$ :

Set $D_{t+1}(i) = \dfrac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$
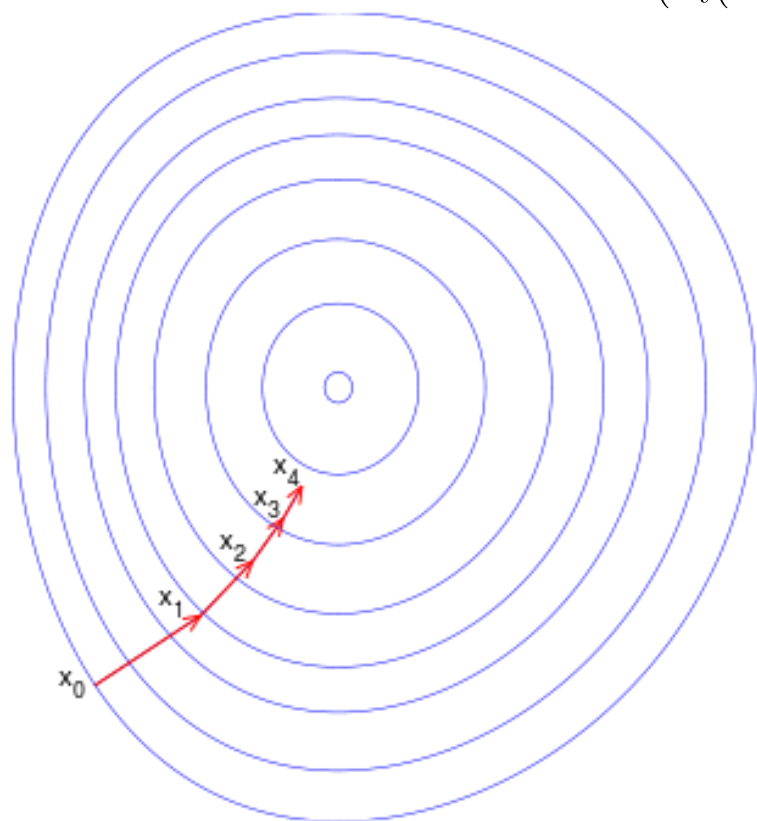
# Functional Gradient Descent



$$J(F_t(\mathbf{x})) = \frac{1}{N} \sum_{i=1}^{N} L(y_i F_t(\mathbf{x}_i)),$$

**Direction in function space**

$$D_i^{t+1} = \frac{\frac{\partial}{\partial y_i F_t(\mathbf{x}_i)} J(F_t(\mathbf{x}))}{\sum_{j=1}^{N} \frac{\partial}{\partial y_j F_t(\mathbf{x}_j)} J(F_t(\mathbf{x}))}$$

**Step size**

$$\alpha_t^* = \arg\min_{\alpha_t} \left[ \frac{1}{N} \sum_{i=1}^{N} L\Big(y_i (F_{t-1}(\mathbf{x}_i) + \alpha_t h_t(\mathbf{x}_i))\Big) \right].$$
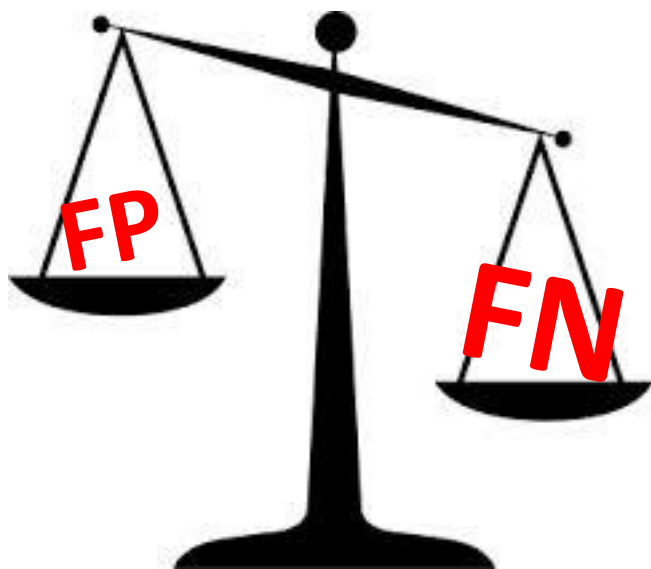
**Property:** FGD-consistency

Are the voting weights and distribution updates consistent with each other?

(i.e. both derivable by FGD on a given loss)

# Decision theory



Ideally: Assign each example to **risk-minimizing** class:

*Predict class $y = 1$ iff*

$$\hat{p}(y = 1|\mathbf{x}) \;>\; \frac{c_{FP}}{c_{FP} + c_{FN}}$$

$$\begin{bmatrix} 0 & c_{FN} \\ c_{FP} & 0 \end{bmatrix}$$

**Property:** Cost-consistency

Does the algorithm use the above
(Bayes Decision Rule)
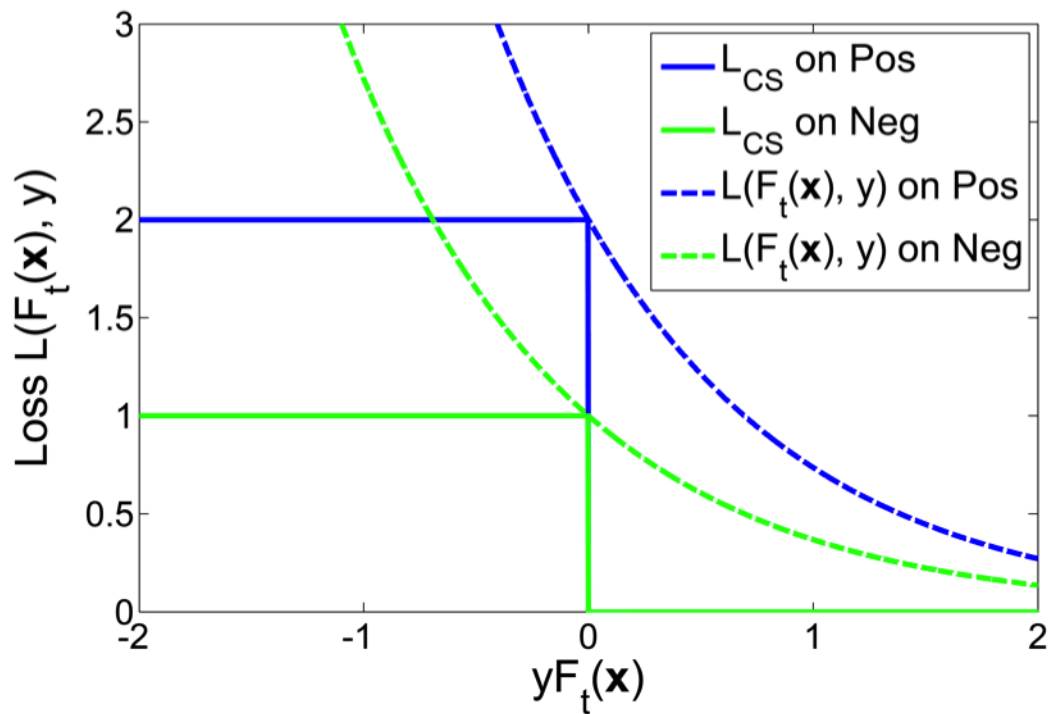to make decisions?

(assuming 'good' probability estimates)

# Margin theory



Sign of margin: encodes correct (>0) or incorrect (<0) classification of (**x**,y)
Magnitude of margin: encodes confidence of boosting ensemble in its prediction

**Large margins** encourage small **generalization error**.
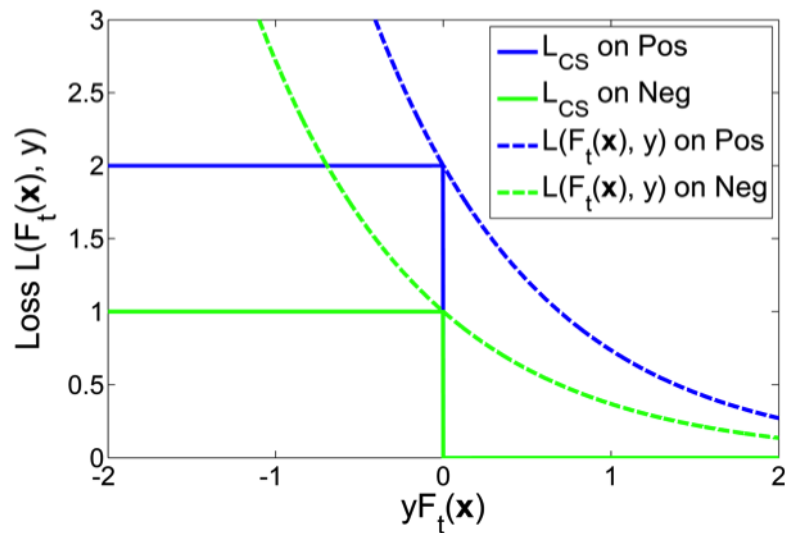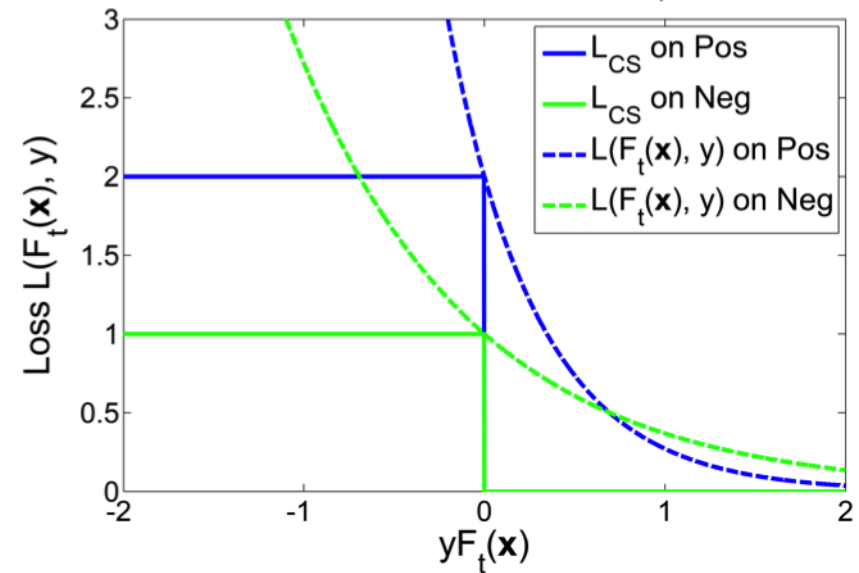Adaboost promotes **large margins**.

# Margin theory – with costs…



**Different surrogate losses for each class.**

# So for a cost sensitive boosting algorithm…

We expect this to be the case.



But some algorithms do this…



**Property:** Asymmetry preservation

Does the loss function preserve the **relative** importance of each class, for all margin values?

# Probabilistic models

'AdaBoost does not produce good probability estimates.'
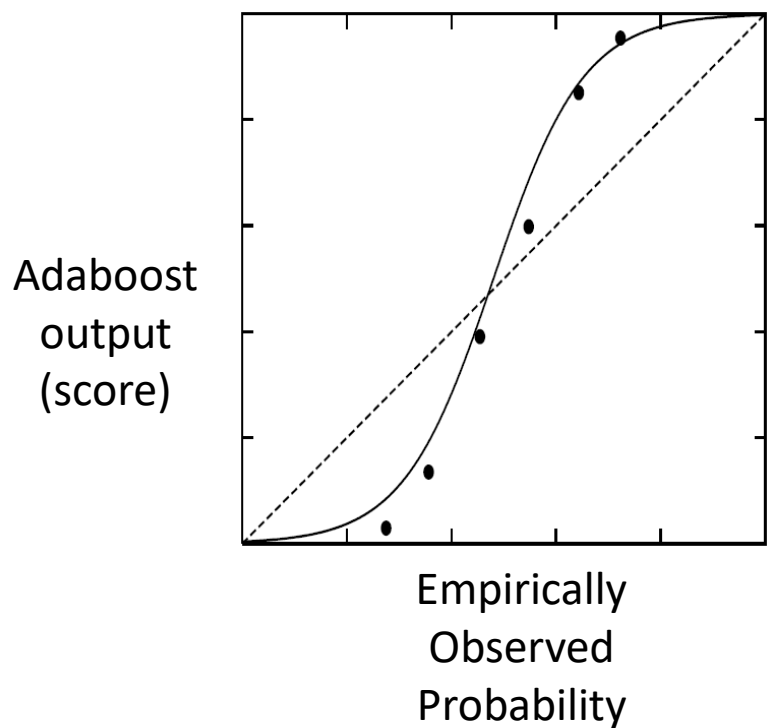
Niculescu-Mizil & Caruana, 2005

'AdaBoost is successful at [..] classification [..] but not class probabilities.'
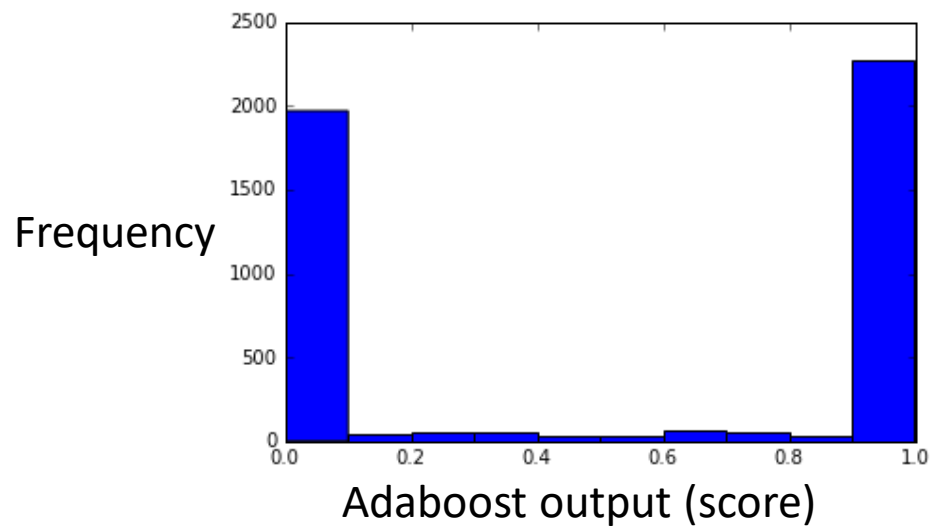
Mease et al., 2007

'This increasing tendency of [the margin] impacts the probability estimates by causing them to quickly diverge to 0 and 1.'

Mease & Wyner, 2008

# Probabilistic models

Adaboost tends to produce probability estimates **close to 0 or 1**.



Adaboost output (score) vs. Empirically Observed Probability



Frequency vs. Adaboost output (score)

# Why this distortion?

Estimates of form:

$$\hat{p}(y=1|\mathbf{x}) = \frac{\sum_{\tau:h_\tau(\mathbf{x})=1}\alpha_\tau}{\sum_{\tau=1}^{t}\alpha_\tau}$$

(Niculescu-Mizil & Caruana, 2005)

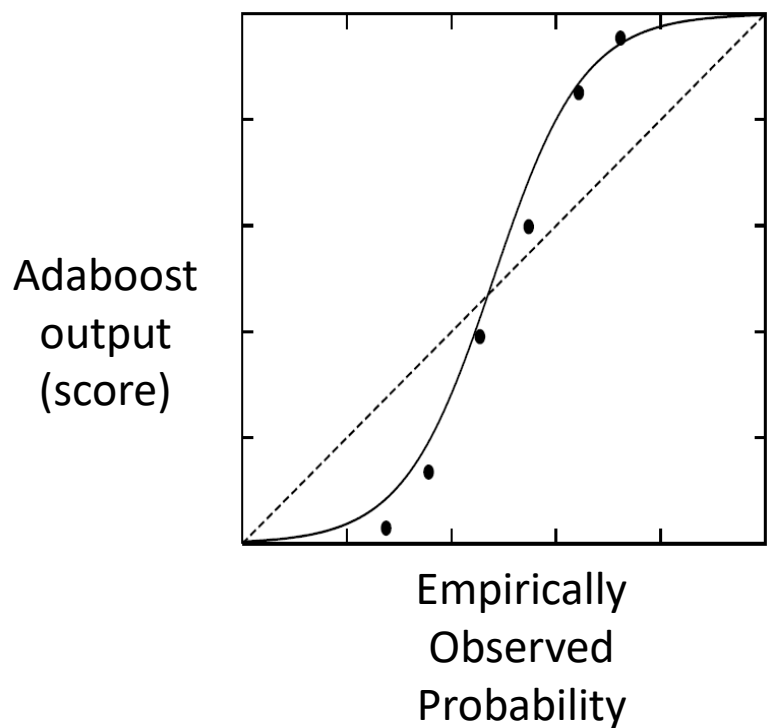As **margin** is **maximized** on training set, scores will tend to 0 or 1.

Estimates of form:

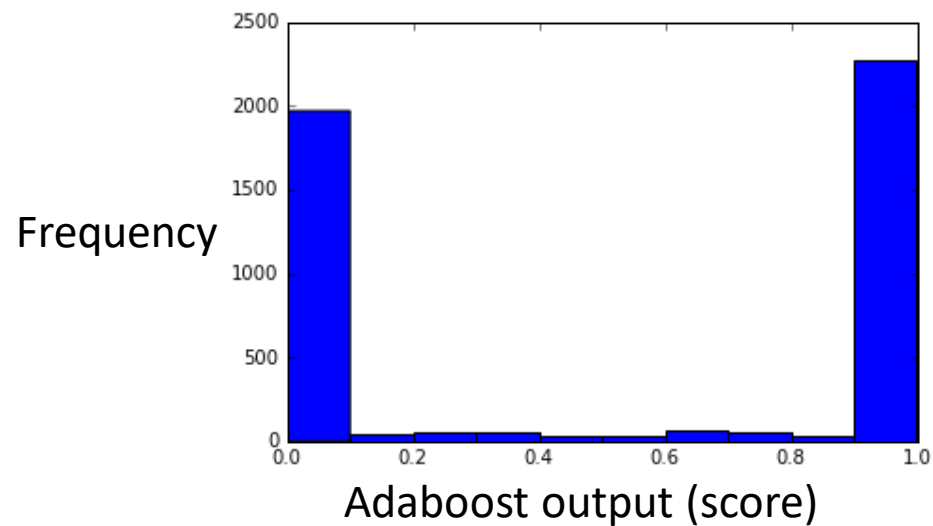$$\hat{p}(y=1|\mathbf{x}) = \frac{1}{1+e^{-2F_t(\mathbf{x})}}$$

(Friedman, Hastie & Tibshirani, 2000)

**Product of Experts**; if one term close to 0 or 1, it dominates.

# Probabilistic Models



Adaboost
output
(score)

Empirically
Observed
Probability

Adaboost tends to produce probability estimates close to 0 or 1.



Frequency

Adaboost output (score)

**Property:** Calibrated estimates

Does the algorithm generate "calibrated" probability estimates?

# Does a given algorithm satisfy…

**Property:** FGD-consistency

Are the **steps consistent** with each other?

(i.e. both voting weights and distribution updates derivable by FGD on same loss)

**Property:** Cost-consistency

Does the algorithm use the **(risk-minimizing) Bayes Decision Rule** to make decisions?

(assuming 'good' probability estimates)

**Property:** Asymmetry preservation

Does the loss function **preserve the relative importance of each class**, for all margin values?

**Property:** Calibrated estimates

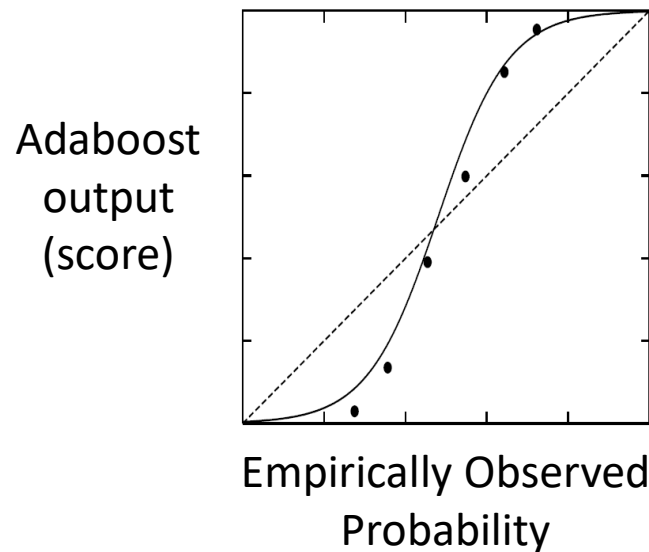Does the algorithm generate **"calibrated" probability estimates**?

# The results are in…

| Method | FGD-consistent | Cost-consistent | Asymmetry-preserving | Calibrated estimates |
|---|---|---|---|---|
| AdaBoost (Freund & Schapire 1997) | ✓ | | ✓ | |
| AdaCost (Fan et al. 1999) | | | | |
| AdaCost($\beta_2$) (Ting 2000) | | | | **All algorithms** produce **uncalibrated** probability estimates! |
| CSB0 (Ting 1998) | | | ✓ | |
| CSB1 (Ting 2000) | | | ✓ | |
| CSB2 (Ting 2000) | | | ✓ | |
| AdaC1 (Sun et al. 2005, 2007) | | ✓ | | |
| AdaC2 (Sun et al. 2005, 2007) | ✓ | | ✓ | |
| AdaC3 (Sun et al. 2005, 2007) | | | | |
| CSAda (Mashnadi-Shirazi & Vasconselos 2007, 2011) | ✓ | ✓ | | |
| AdaDB (Landesa-Vázquez & Alba-Castro 2013) | ✓ | ✓ | | |
| AdaMEC (Ting 2000, Nikolaou & Brown 2015) | ✓ | ✓ | ✓ | |
| CGAda (Landesa-Vázquez & Alba-Castro 2012, 2015) | ✓ | ✓ | ✓ | |
| AsymAda (Viola & Jones 2002) | ✓ | ✓ | ✓ | |

So could we just calibrate these last three?  We use "Platt scaling".
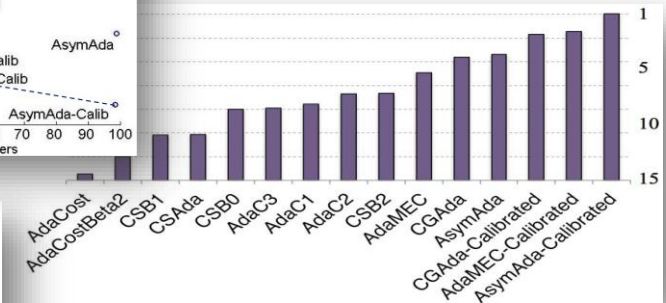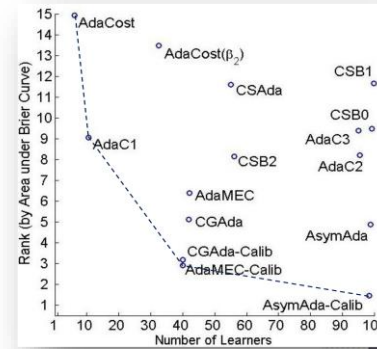
# Platt scaling (logistic calibration)

**Training**: Reserve part of training data (here 50% -more on this later) to fit a sigmoid to correct the distortion:
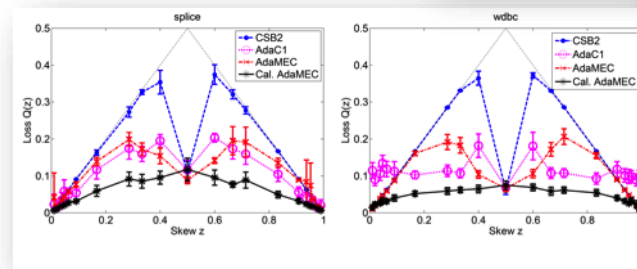


Adaboost output (score)

Empirically Observed Probability

**Prediction**: Apply sigmoid transformation to score (output of ensemble) to get probability estimate

# Experiments

- 15 algorithms.
- 18 datasets.
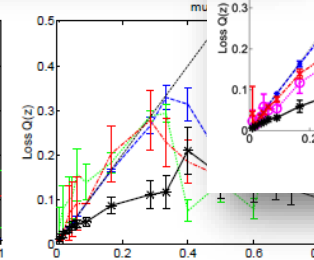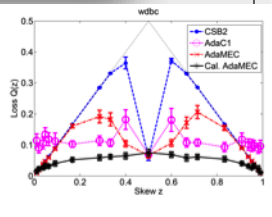- 21 degrees of cost imbalance.

# In summary…



**Average Brier Score Rank**

All except calibrated

All 4 Properties

AdaCost, AdaCostBeta2, CSB1, CSAda, CSB0, AdaC3, AdaC1, AdaC2, CSB2, AdaMEC, CGAda, AsymAda, CGAda-Calibrated, AdaMEC-Calibrated, AsymAda-Calibrated

AdaMEC, CGAda & AsymAda **outperform all others.**

Their **calibrated** versions **outperform** the **uncalibrated** ones

# In summary...

"Calibrated-AdaMEC" was one of the top methods.

1. Take <u>original</u> Adaboost.

2. Calibrate it (we use Platt scaling)

3. Shift the decision threshold.... $\dfrac{c_{FP}}{c_{FP} + c_{FN}}$

**Consistent** with all theory perspectives.

**No** extra **hyperparameters** added.

**No need to retrain** if cost ratio changes.

Consistently **top (or joint top)** in empirical comparisons.

# Methods & properties

| Method | FGD-consistent | Cost-consistent | Asymmetry-preserving | Calibrated estimates |
|---|---|---|---|---|
| AdaBoost (Freund & Schapire 1997) | ✓ | | ✓ | |
| AdaCost (Fan et al. 1999) | | | | |
| AdaCost($\beta_2$) (Ting 2000) | | | | |
| CSB0 (Ting 1998) | | | ✓ | |
| CSB1 (Ting 2000) | | | ✓ | **All algorithms** produce **uncalibrated** probability estimates! |
| CSB2 (Ting 2000) | | | ✓ | |
| AdaC1 (Sun et al. 2005, 2007) | | ✓ | | |
| AdaC2 (Sun et al. 2005, 2007) | ✓ | | ✓ | |
| AdaC3 (Sun et al. 2005, 2007) | | | | |
| CSAda (Mashnadi-Shirazi & Vasconselos 2007, 2011) | ✓ | ✓ | | |
| AdaDB (Landesa-Vázquez & Alba-Castro 2013) | ✓ | ✓ | | |
| AdaMEC (Ting 2000, Nikolaou & Brown 2015) | ✓ | ✓ | ✓ | |
| CGAda (Landesa-Vázquez & Alba-Castro 2012, 2015) | ✓ | ✓ | ✓ | |
| AsymAda (Viola & Jones 2002) | ✓ | ✓ | ✓ | |

So could we just calibrate these last three?  We use "Platt scaling".

# Q: What if we calibrate all methods?

A: In **theory**, …
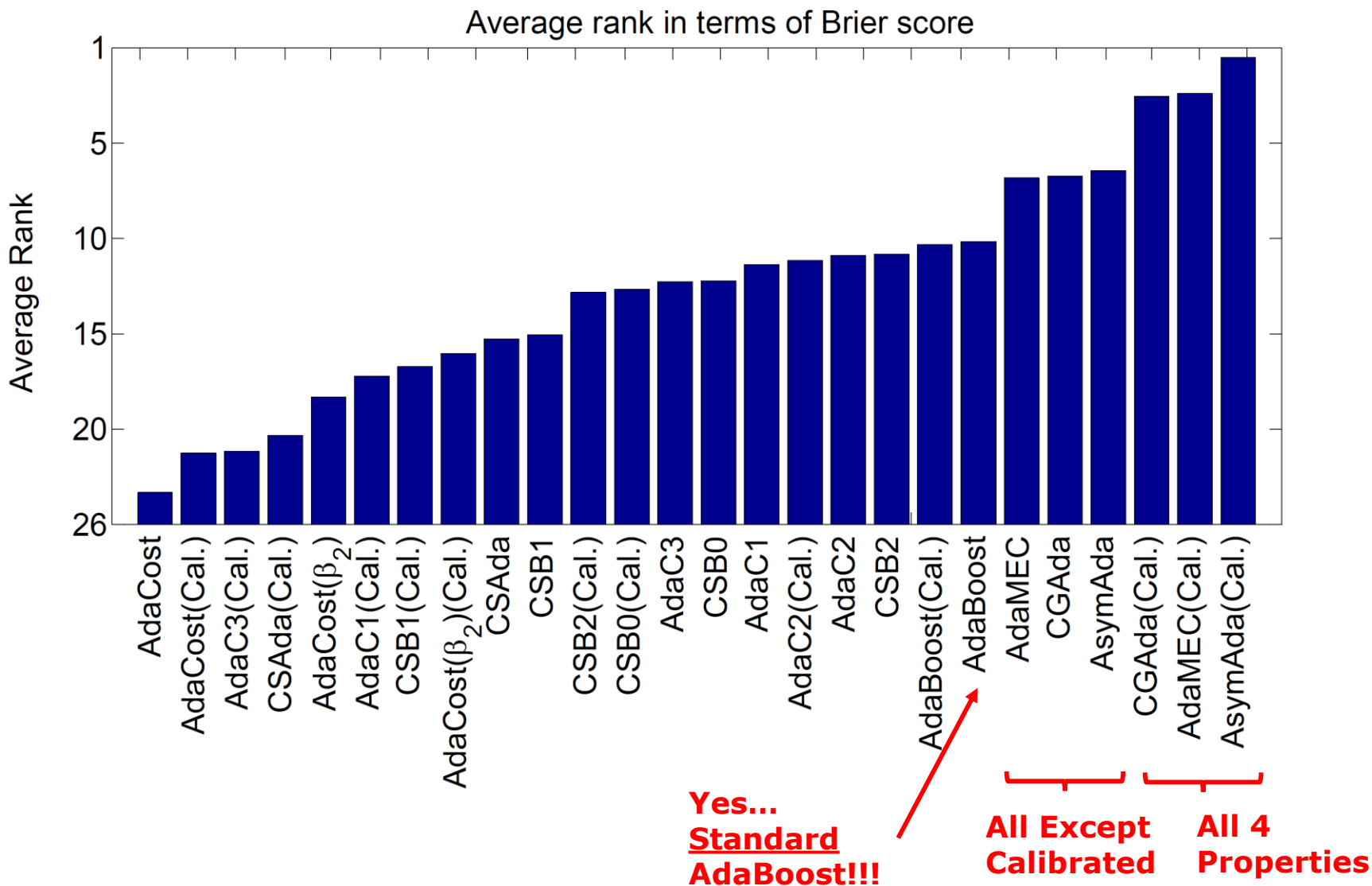
… calibration improves probability estimates.

… if a method is not cost-sensitive, will not make it.

… if the steps are not consistent, will not make them.

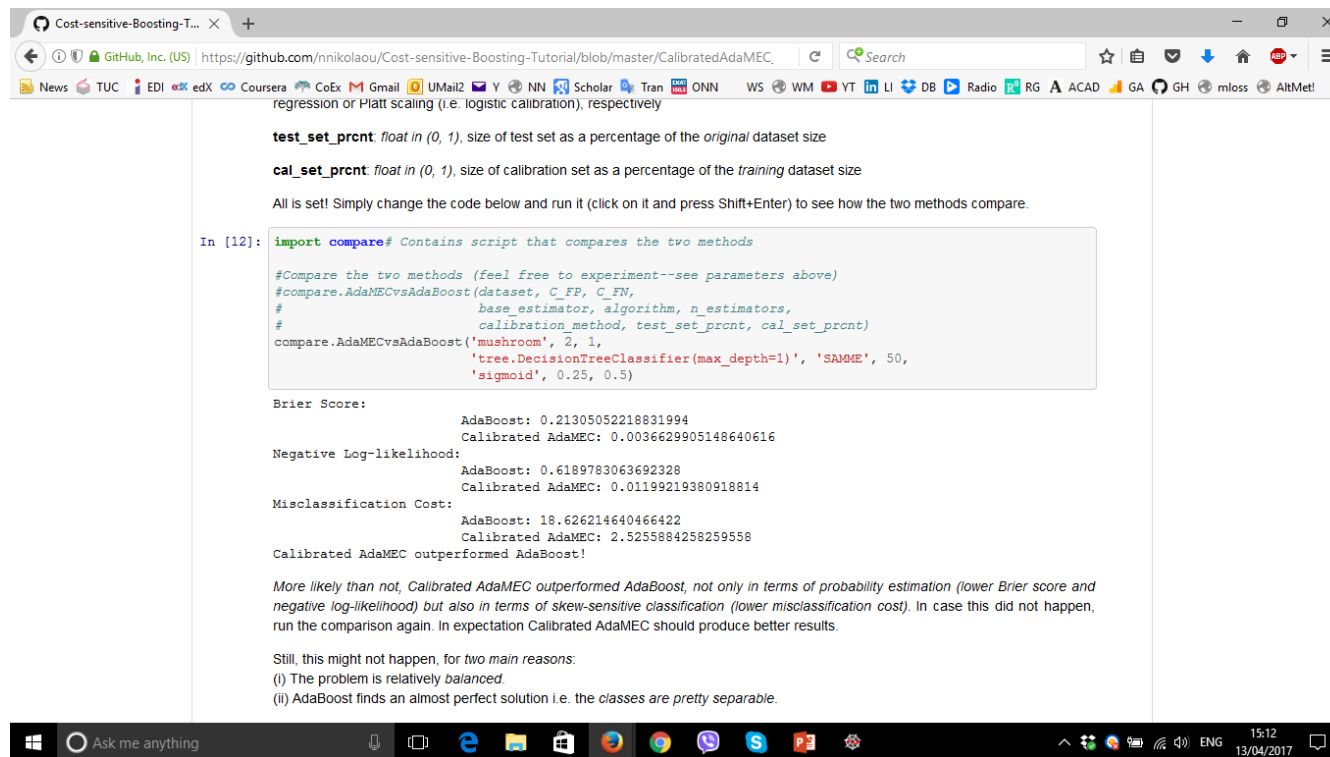… if class importance is swapped during training, will not correct.

# Results



Average rank in terms of Brier score

Yes… Standard AdaBoost!!!

All Except Calibrated

All 4 Properties

# Methods & properties

| Method | FGD-consistent | Cost-consistent | Asymmetry-preserving | Calibrated estimates |
|---|---|---|---|---|
| AdaBoost (Freund & Schapire 1997) | ✓ | | ✓ | |
| AdaCost (Fan et al. 1999) | | | | |
| AdaCost($\beta_2$) (Ting 2000) | | | | |
| CSB0 (Ting 1998) | | | ✓ | |
| CSB1 (Ting 2000) | | | ✓ | **All algorithms** produce **uncalibrated** probability estimates! |
| CSB2 (Ting 2000) | | | ✓ | |
| AdaC1 (Sun et al. 2005, 2007) | | ✓ | | |
| AdaC2 (Sun et al. 2005, 2007) | ✓ | | ✓ | |
| AdaC3 (Sun et al. 2005, 2007) | | | | |
| CSAda (Mashnadi-Shirazi & Vasconselos 2007, 2011) | ✓ | ✓ | | |
| AdaDB (Landesa-Vázquez & Alba-Castro 2013) | ✓ | ✓ | | |
| AdaMEC (Ting 2000, Nikolaou & Brown 2015) | ✓ | ✓ | ✓ | |
| CGAda (Landesa-Vázquez & Alba-Castro 2012, 2015) | ✓ | ✓ | ✓ | |
| AsymAda (Viola & Jones 2002) | ✓ | ✓ | ✓ | |

So could we just calibrate these last three?  We use "Platt scaling".

# Q: Sensitive to calibration choices?

A: Check it out on your own!

https://github.com/nnikolaou/Cost-sensitive-Boosting-Tutorial

# Results

Isotonic regression > Platt scaling, for larger datasets

Can do better than 50%-50% train-calibration split (problem dependent; see Part II)

(Calibrated) Real AdaBoost > (Calibrated) Discrete AdaBoost...

# In summary…

"Calibrated-AdaMEC" was one of the top methods.

    1. Take <u>original</u> Adaboost.

    2. Calibrate it (we use Platt scaling)

    3. Shift the decision threshold…. $\dfrac{c_{FP}}{c_{FP} + c_{FN}}$

**Consistent** with all theory perspectives.

**No** extra **hyperparameters** added.

**No need to retrain** if cost ratio changes.

Consistently **top (or joint top)** in empirical comparisons.

# Conclusions

We analyzed the cost-sensitive boosting literature

… **15+** variants over **20** years, from **4** different theoretical perspectives

"Cost sensitive" modifications to the **original** Adaboost are not needed...

**… if** the scores are properly calibrated,
**and** the decision threshold is shifted according to the cost matrix.

# Relevant publications

- N. Nikolaou and G. Brown, *Calibrating AdaBoost for Asymmetric Learning*, Multiple Classifier Systems, 2015

- N. Nikolaou, N. Edakunni, M. Kull, P. Flach and G. Brown, *Cost-sensitive Boosting algorithms: Do we really need them?*, Machine Learning Journal, Vol. 104, Issue 2, Sept 2016
  - Best Poster Award, INIT/AERFAI summer school in ML 2014
  - Plenary Talk ECML 2016 -- 12/129 eligible papers (9.3%)
  - Best Paper Award 2016, School of Computer Science, University of Manchester

- N. Nikolaou, *Cost-sensitive Boosting: A Unified Approach*, PhD Thesis, University of Manchester, 2016
  - Best Thesis Award 2017, School of Computer Science, University of Manchester

# Resources & code

- Easy-to-use but not so flexible 'Calibrated AdaMEC' python implementation (scikit-learn style):

    https://mloss.org/revision/view/2069/

- i-python tutorial for all this with interactive code for 'Calibrated AdaMEC', where every choice can be tweaked:

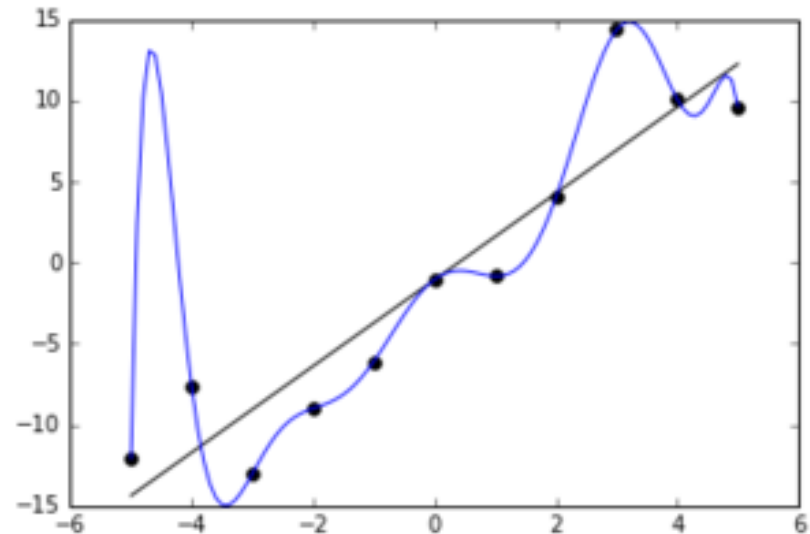    https://github.com/nnikolaou/Cost-sensitive-Boosting-Tutorial

# Connections to Deep Learning (1)

Both **Boosting** and **Deep Neural Networks** (DNNs) exhibit **very good generalization**…

..despite constructing **overparameterized** (drawn from a very rich family) **models**

**Too high richness** (capacity, complexity, degrees of freedom) of model → **overfitting**

# Connections to Deep Learning (2)

**Overfitting**: **fitting the training dataset 'too well'**, 'memorizing it' rather than 'learning from it', capturing noise as part of the concept to be learned thus **failing to generalize to new data (poor performance on test set)**



But both **Boosting** & **DNNs** can **improve fitting the test data even beyond the point of perfectly fitting the training data!**

''**Boosting the margin: a new explanation for the effectiveness of voting methods**'', Schapire et al. 1997
''**Understanding Deep Learning Requires Rethinking Generalization**'', Zhang et al, 2017
''**Opening the Black Box of Deep Neural Networks via Information**'', Shwartz-Ziv & Tishby, 2017

# Connections to Deep Learning (3)

The **good classification generalization of DNNs** has been justified through

- **margin maximization**:

''**Robust Large Margin Deep Neural Networks**'', Sokolic et al., 2017

[Note: As with Boosting]

- **properties of (Stochastic) GD**:

''**A Bayesian Perspective on Generalization and Stochastic Gradient Descent**'', Smith & Le, 2017

''**The Implicit Bias of Gradient Descent on Separable Data**'', Soudry et al., 2017

[Note: Boosting also a Gradient Descent process; stochasticity also applied/substituted by other mechanisms]

- **information theory**:

''**Opening the Black Box of Deep Neural Networks via Information**'', Shwartz-Ziv & Tishby, 2017

[Note: **We are currently applying similar ideas to justify generalization in Boosting-seems to work!**]

**Residual Networks (ResNets),** a state of the art DNN architecture has been directly **explained through boosting theory**

''**Learning Deep ResNet Blocks Sequentially using Boosting Theory**'', Huang et al., 2017

# Connections to Deep Learning (4)

**ResNets** also **very good classifiers** but **very poor probability estimators**

**"On Calibration of Modern Neural Networks"**, Guo et al., 2017

**CONJECTURE : Not a coincidence!** [direct analogy to boosting]

Similar behaviour in **other architectures**...

**"Understanding Deep Learning Requires Rethinking Generalization"**, Zhang et al, 2017

**"Regularizing Neural Networks by Penalizing Confident Output Distributions"**, Pereyra et al., 2017

**CONJECTURE: Also not a coincidence!** [implicit regularization afforded by GD optimization ≡ margin maximization: good for generalization but scores are distorted towards the extremes]

At any rate, when solving **probability estimation/cost-sensitive** problems using **DNNs** you **should calibrate their outputs**!

# End of Part I

# Questions?

# Part II:
# Calibrating Online Boosting

# Next Step: Online learning

**Examples** presented **one (or a few) @ a time**

Learner makes **predictions as examples are received**

Each '**minibatch**' used to **update model**, then discarded; **constant time & space complexity**

Why?
- Data arrive this way (**streaming**)
- Problem (e.g. data distribution) **changes over time**
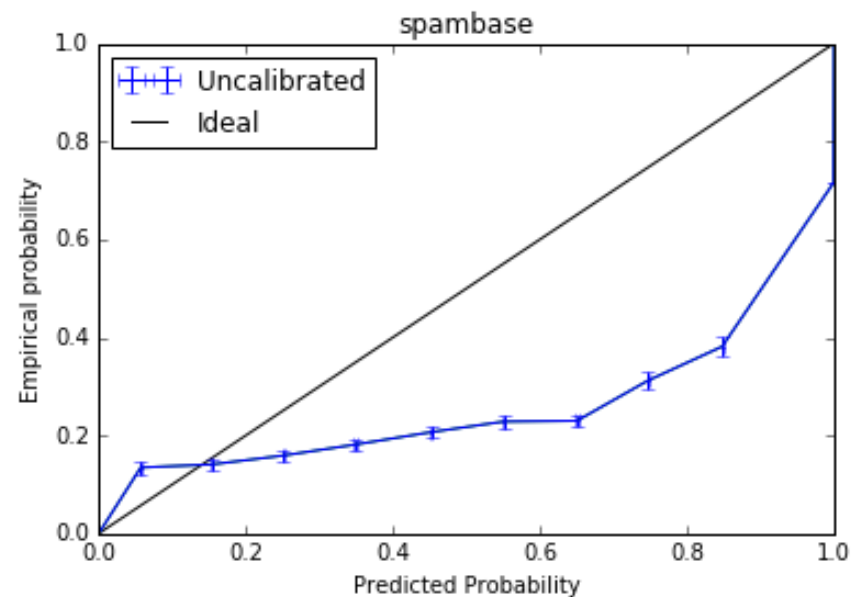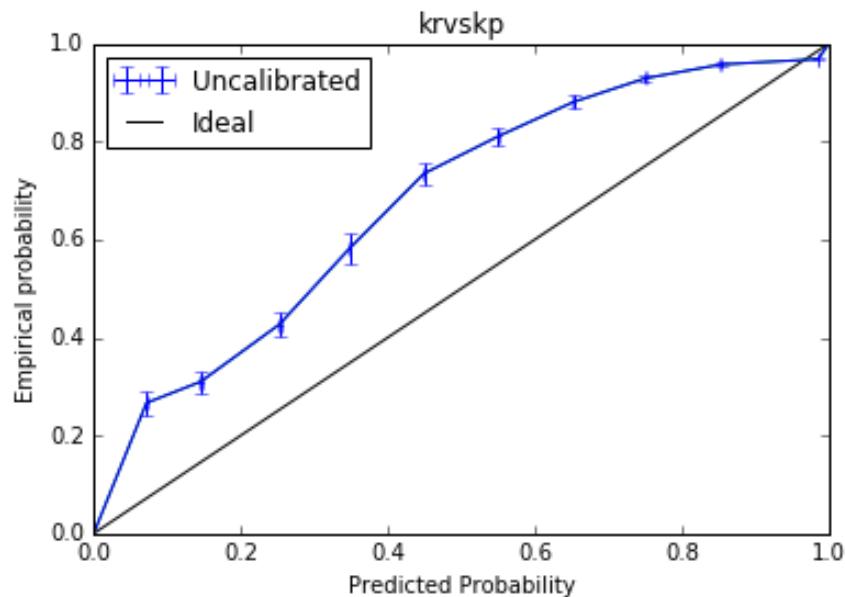- To **speed up learning** in big data applications

# Online learning

For each $minibatch$ $n$ do:
1. **Receive** $n$
2. **Predict label / class probability** of examples in $n$
3. **Get true label** of examples in $n$
4. **Evaluate** learner's performance on $n$
5. **Update** learner **parameters** accordingly

# Online Boosting (Oza, 2004)

Probability estimates -as in AdaBoost- are uncalibrated:

# How to calibrate online Boosting?

**Batch Learning**: **reserve part of the dataset** to train calibrator function (logistic sigmoid, if Platt scaling)

**Online learning**: **cannot do this**; on each minibatch we must **decide** whether to **train ensemble or calibrator**

How to make this decision?

# Naïve approach

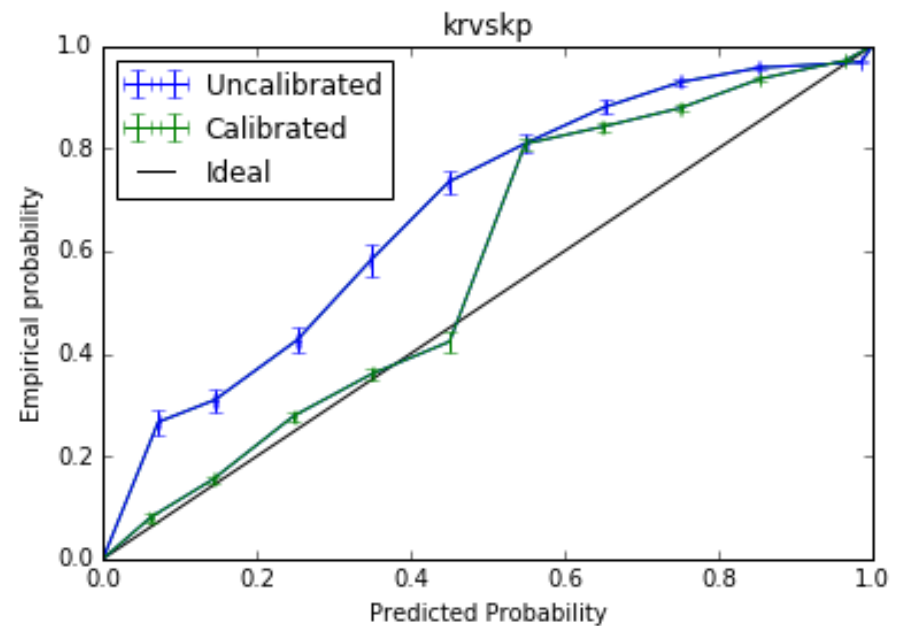Fixed Policy: calibrate **every $N$ rounds**

How to pick $N$?
- Will depend on **problem**
- Will depend on **ensemble hyperparameters**
- Will depend on **calibrator hyperparameters**
- **Might change** during training…

In batch learning can **choose via cross-validation**; **not here**

# Still, naïve better than nothing

Results with N = 2 (**not** necessarily best value):

# A more refined approach

- What if we could **learn** a good sequence of alternating between actions?



**Bandit Algorithms**

# Bandit optimization

A **set of actions (arms)** -on each round we choose one

Each action associated with a **reward distribution**

Each time an action taken we **sample** its reward distribution

**Sequence of actions** that **minimize cumulative regret?**

**Exploration vs. Exploitation**

In **online calibrated boosting**:

**Two actions**: **{ train , calibrate }**

**Reward**: **Increase in overall model likelihood** after action

# Thompson sampling

A **Bayesian** take on bandits for updating reward distribution

Assume **rewards are Gaussian**; start with **Gaussian prior**, then **update** using **self-conjugacy of Gaussian distribution**

Take action with **highest posterior reward**

# UCB policies

**'Optimism in the face of uncertainty'**

Choose not the action with best expected reward, but that with **highest upper bound on reward**

Bounds derived for arbitrary (UCB1, UCB1-Improved) or specific (KL-UCB) reward distributions

# Discounted rewards

**'Forgeting the past'**

Weigh past rewards less; protects from **<span style="color:red">non-stationarity</span>**
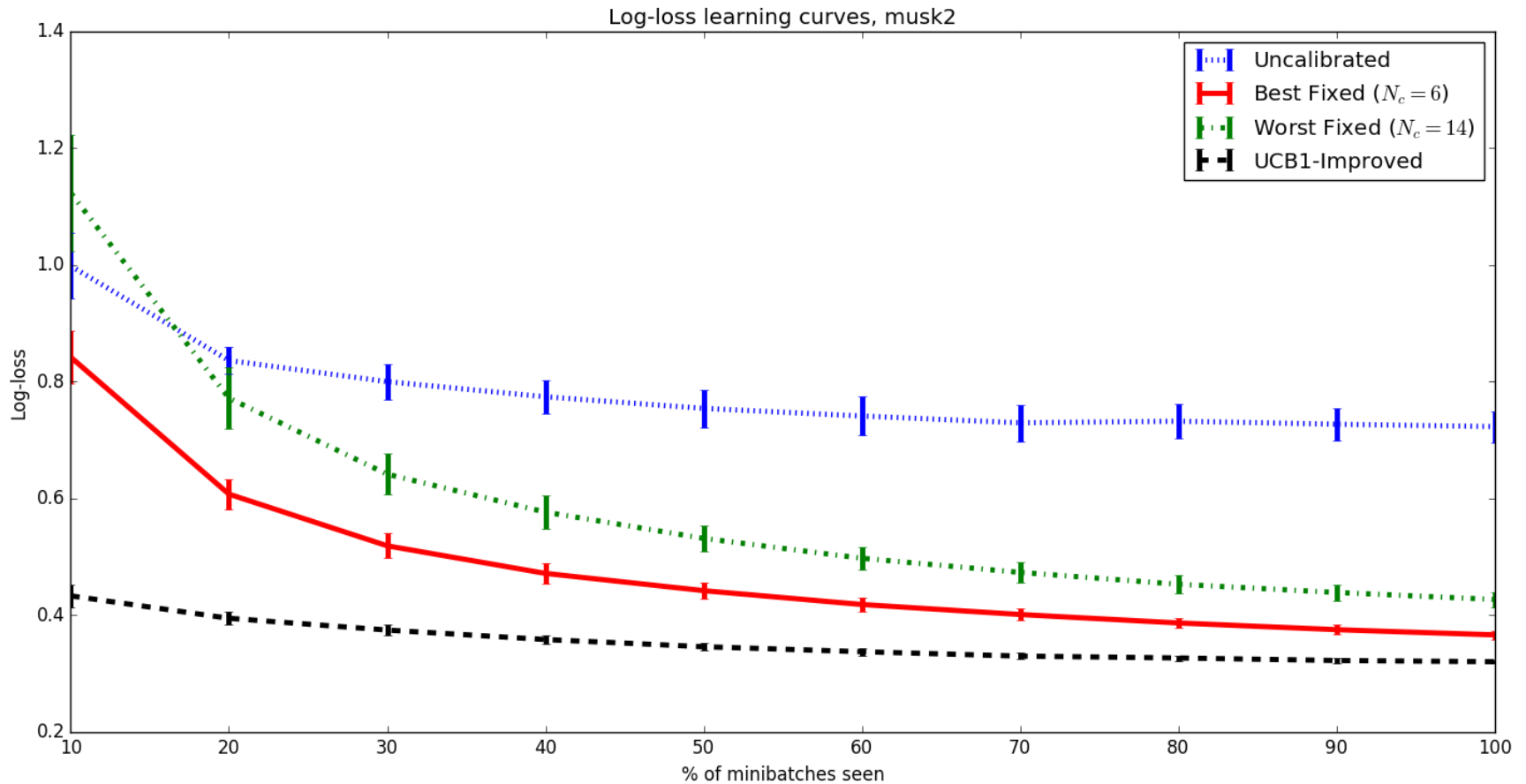
Why non-stationary?
- **Data distribution** might **change**…
- …most importantly: **reward distributions** will **change**:
  if we perform one action many times, the relative reward for
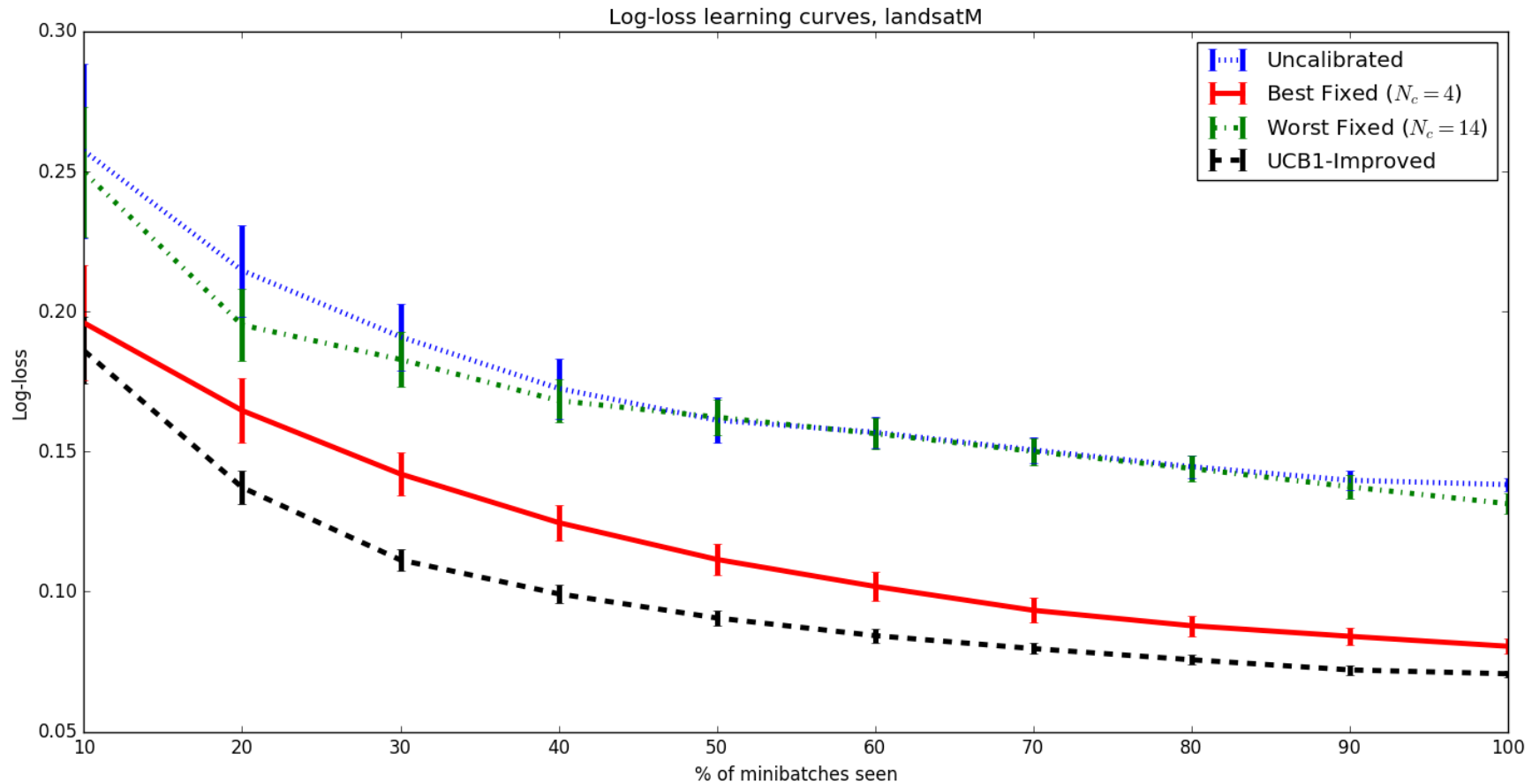  performing the other is expected to have increased

# Some initial results

- Uncalibrated

  vs. 'Every $N$ policies' $N \in \{2, 4, 6, 8, 10, 12, 14\}$

  vs. UCB1, UCB1-Improved, Gaussian Thompson Sampling

  vs. Discounted versions of above

- Initial results:
  - calibrating (even naive) > not calibrating
  - non-discounted UCB1 variants ≥ **best** 'Every N' policy
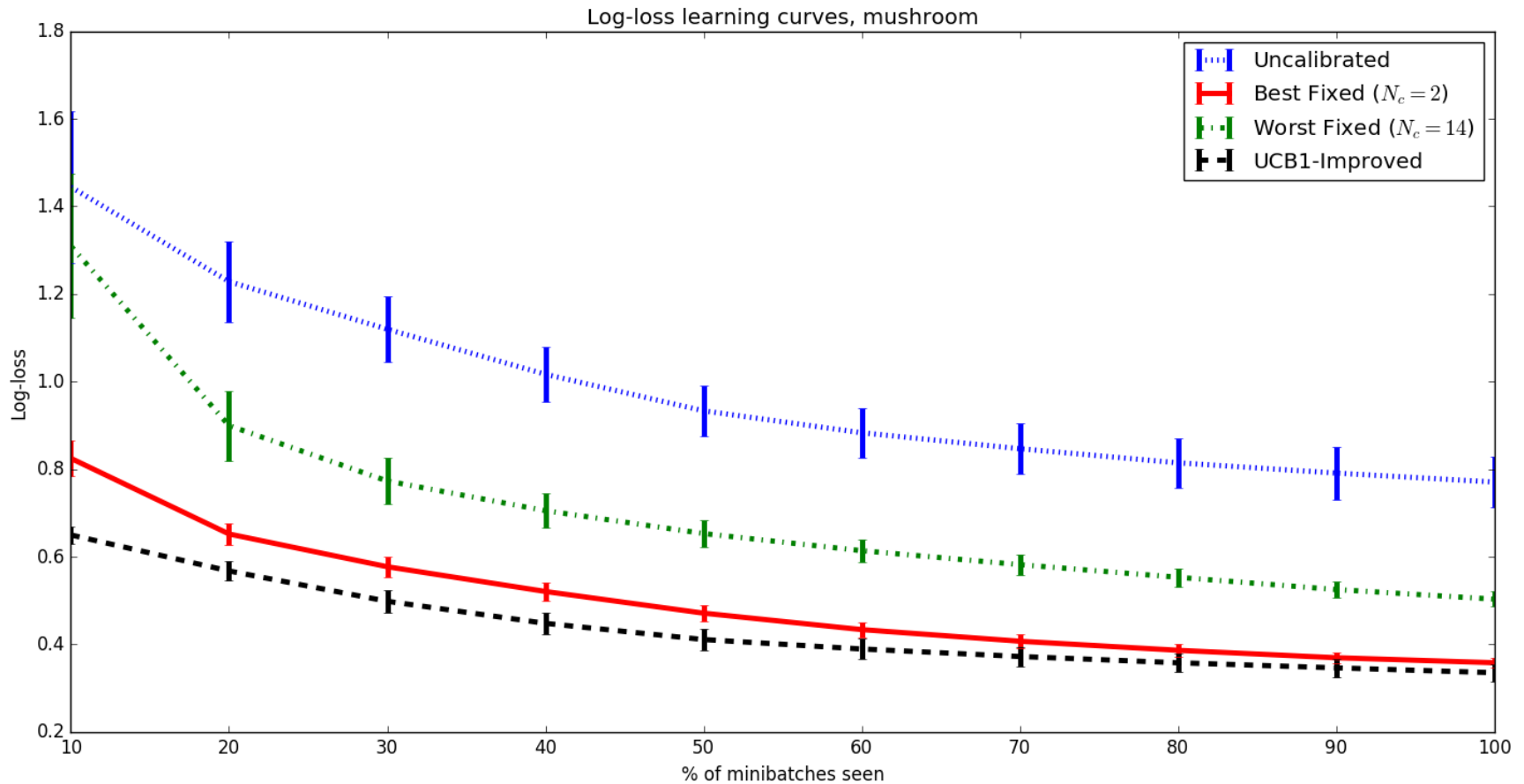  - discounted Thompson Sampling ≥ **best** 'Every N' policy
  - … plus no need to set $N$

# Log-loss learning curves (Impr. UCB1)



Log-loss learning curves, musk2

# Log-loss learning curves (Impr. UCB1)



Log-loss learning curves, landsatM
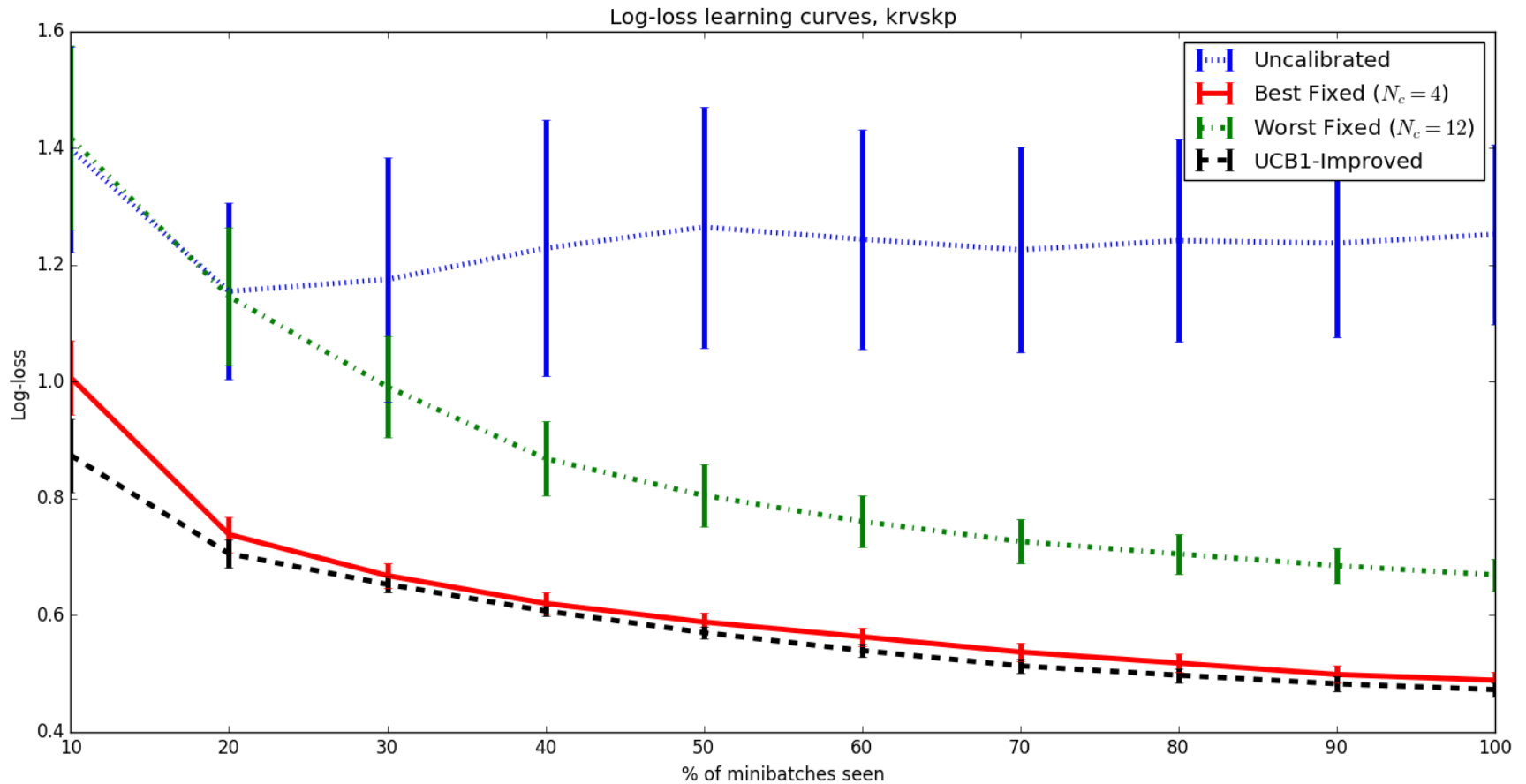
# Log-loss learning curves (Impr. UCB1)



Log-loss learning curves, mushroom

# Log-loss learning curves (Impr. UCB1)



Log-loss learning curves, krvskp

# Some Notes

Results shown for ensembles of M=10 Naïve Bayes weak learners

Similar results for
   **other bandit policies**
   **other weak learners**
   **regularized weak learners**
   **varying ensemble sizes**
   **presence of inherent non-stationarity**

Also **beats other Naïve policies** (mention)

# In summary…

Online Boosting **poor probability estimates**; some **calibration** can improve

**Learn** a good sequence of calibration / training actions using **bandits**

**Online**, **fast**, **at least as good as 'best naïve' + adaptive to non-stationarity**
Easy to **adapt to other problems** (e.g. cost-sensitive learning)
**Robust** to ensemble/calibrator **hyperparameters**

Extensions: e.g. **adversarial**, **contextual**, **more actions**, **refine calibration**, …

# Thank you! Ευχαριστώ!

## Questions?