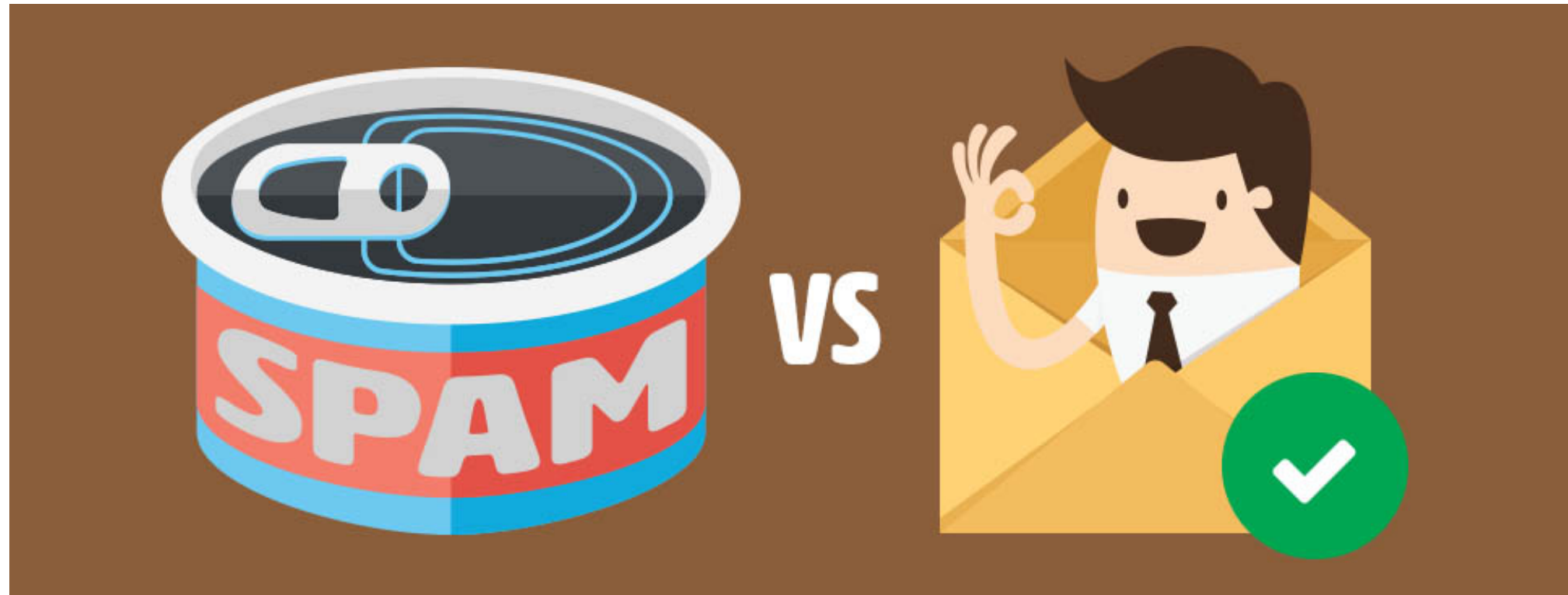


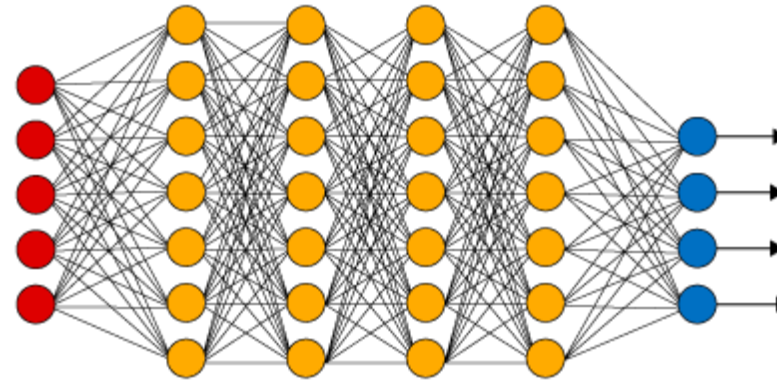
How effective is your classifier?

Revisiting the role of metrics in machine learning

SANMI KOYEJO

Joint work with Ran Li, Xiaoyan Wang, Gaurush Hiranandani, Shant Boodaghians, and Ruta Mehta





- Accuracy = 95%
- \$\$\$

- Users complain that most real emails are labelled spam
- ~90% of all email is spam*
- Suggests that accuracy is the wrong metric as it gives equal weight to all errors

HOW SHOULD YOU
MEASURE
PERFORMANCE?

It depends... on the relative cost/benefit of different kinds of errors.

The **metric** is a quantitative description of tradeoffs.

Error analysis for binary classification

		Ground truth	
		Spam	Not Spam
Predicted	Spam	TP	FP
	Not Spam	FN	TN

- $\text{Accuracy} = \text{TP} + \text{TN} = 1 - \text{FP} - \text{FN}$
- To improve user calibration, try evaluating and/or optimizing weighted accuracy e.g.

$$\phi(h) = 1 - 0.1 \text{FP} - \text{FN}$$

The confusion matrix

		Ground truth		$= \mathbf{C}(h)$
		Y = 1	Y = 0	
Predicted	$h(x) = 1$	TP	FP	
	$h(x) = 0$	FN	TN	

Beyond Accuracy, more general metrics are nested functions

$$\phi(h) = \psi(\mathbf{C}(h))$$

- Metrics are used to compare classifiers, or can be optimized directly
- The classifier performance metric can be approximated from data.

Lots of real-world examples

$$\phi(h) = a_1 \text{TP} + a_2 \text{FP} + a_3 \text{FN} + a_4 \text{TN}$$

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \text{ TNR} = \frac{\text{TN}}{\text{FP} + \text{TN}}, \text{ Prec} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \text{ FNR} = \frac{\text{FN}}{\text{FN} + \text{TP}}, \text{ NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}}.$$

$$\text{AM} = \frac{1}{2} \left(\frac{\text{TP}}{\pi} + \frac{\text{TN}}{1 - \pi} \right) = \frac{(1 - \pi)\text{TP} + \pi\text{TN}}{2\pi(1 - \pi)}, \quad F_\beta = \frac{(1 + \beta^2)\text{TP}}{(1 + \beta^2)\text{TP} + \beta^2\text{FN} + \text{FP}} = \frac{(1 + \beta^2)\text{TP}}{\beta^2\pi + \gamma},$$

$$\text{JAC} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}} = \frac{\text{TP}}{\pi + \text{FP}} = \frac{\text{TP}}{\gamma + \text{FN}}, \quad \text{WA} = \frac{w_1\text{TP} + w_2\text{TN}}{w_1\text{TP} + w_2\text{TN} + w_3\text{FP} + w_4\text{FN}}.$$

$$\pi = \text{TP} + \text{FN}, \quad \gamma = \text{TP} + \text{FP}$$

The Netflix logo, consisting of the word "NETFLIX" in a bold, red, sans-serif font.The Amazon logo, featuring the word "amazon" in a black, lowercase, sans-serif font with a curved orange arrow underneath it.The LinkedIn logo, featuring the word "Linked" in a black, sans-serif font followed by a blue square containing the white letters "in".The Google logo, featuring the word "Google" in its multi-colored, sans-serif font.

Metrics in ranking and recommendation

“Results show that improvements in RMSE often do not translate into [top-N ranking] accuracy improvements. In particular, a naive non-personalized algorithm can outperform some common recommendation approaches and almost match the accuracy of sophisticated algorithms”

P. Cremonesi, Y. Koren, and R. Turrin. "Performance of recommender algorithms on top-n recommendation tasks." Recsys, 2010.

sklearn.metrics : Metrics

Regression metrics

See the [Regression metrics](#) section of the user guide for further details.

<code>metrics.explained_variance_score</code> (y_true, y_pred)	Explained variance regres
<code>metrics.mean_absolute_error</code> (y_true, y_pred)	Mean absolute error regre
<code>metrics.mean_squared_error</code> (y_true, y_pred[, ...])	Mean squared error regre
<code>metrics.mean_squared_log_error</code> (y_true, y_pred)	Mean squared logarithmic
<code>metrics.median_absolute_error</code> (y_true, y_pred)	Median absolute error reg
<code>metrics.r2_score</code> (y_true, y_pred[, ...])	R^2 (coefficient of determ

Multilabel ranking metrics

See the [Multilabel ranking metrics](#) section of the user guide for further details

<code>metrics.coverage_error</code> (y_true, y_score[, ...])	Coverage error mea
<code>metrics.label_ranking_average_precision_score</code> (...)	Compute ranking-ba
<code>metrics.label_ranking_loss</code> (y_true, y_score)	Compute Ranking Ic

Clustering metrics

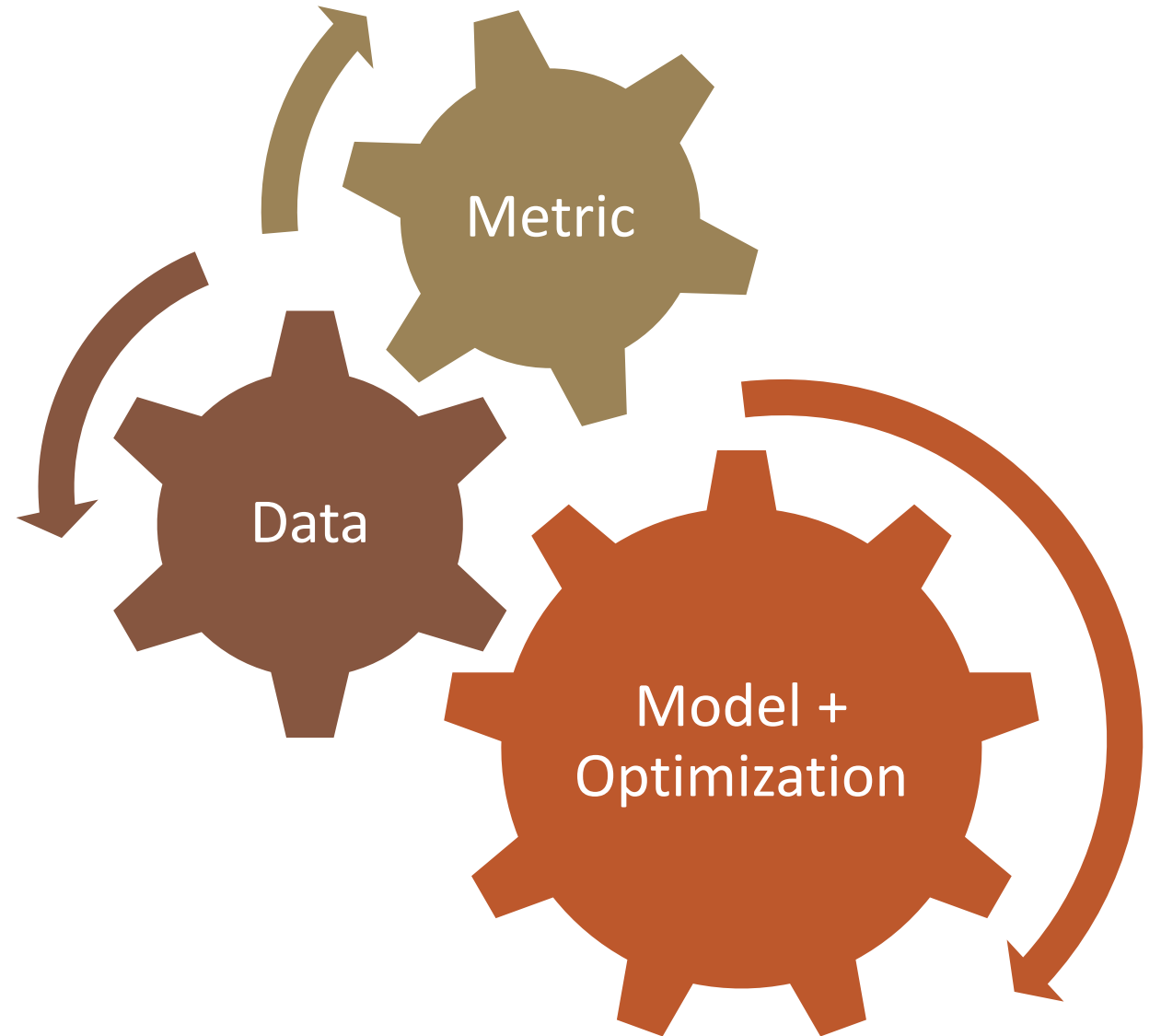
See the [Clustering performance evaluation](#) section of the user guide for further details.

The `sklearn.metrics.cluster` submodule contains evaluation metrics for cluster analysis results. There are two forms of evaluation:

- supervised, which uses a ground truth class values for each sample.
- unsupervised, which does not and measures the 'quality' of the model itself.

<code>metrics.adjusted_mutual_info_score</code> (...[, ...])	Adjusted Mutual Information between two clusterings.
<code>metrics.adjusted_rand_score</code> (labels_true, ...)	Rand index adjusted for chance.
<code>metrics.calinski_harabaz_score</code> (X, labels)	Compute the Calinski and Harabaz score.
<code>metrics.davies_bouldin_score</code> (X, labels)	Computes the Davies-Bouldin score.
<code>metrics.completeness_score</code> (labels_true, ...)	Completeness metric of a cluster labeling given a ground truth.
<code>metrics.cluster.contingency_matrix</code> (...[, ...])	Build a contingency matrix describing the relationship between labels.
<code>metrics.fowlkes_mallows_score</code> (labels_true, ...)	Measure the similarity of two clusterings of a set of points.
<code>metrics.homogeneity_completeness_v_measure</code> (...)	Compute the homogeneity and completeness and V-Measure scores at once.
<code>metrics.homogeneity_score</code> (labels_true, ...)	Homogeneity metric of a cluster labeling given a ground truth.
<code>metrics.mutual_info_score</code> (labels_true, ...)	Mutual Information between two clusterings.
<code>metrics.normalized_mutual_info_score</code> (...[, ...])	Normalized Mutual Information between two clusterings.
<code>metrics.silhouette_score</code> (X, labels[, ...])	Compute the mean Silhouette Coefficient of all samples.
<code>metrics.silhouette_samples</code> (X, labels[, metric])	Compute the Silhouette Coefficient for each sample.
<code>metrics.v_measure_score</code> (labels_true, labels_pred)	V-measure cluster labeling given a ground truth.

Metric selection
is challenging,
but it's a
required step in
modern ML



Metric choice has a large impact on real-world machine learning performance.

1

Given a complex metric, how can we efficiently construct classifiers that (approximately) optimize it?

2

Given a new classification problem, which metric should you use to measure performance?

One simple trick...

A RE-WEIGHTING
STRATEGY

Multiclass classification

		Ground truth			
		Y = 1	Y = 2	...	Y=K
Predicted	h(x) = 1	C ₁₁	C ₁₂		C _{1K}
	h(x) = 2	C ₂₁	C ₂₂		C _{2K}
	⋮				...
	h(x) = K	C _{K1}	C _{K2}		C _{KK}

$$C(h)$$

Standard metric is Accuracy

$$\begin{aligned}\phi(h) &= c_{11} + c_{22} + \dots + c_{KK} \\ &= \langle I, C(h) \rangle\end{aligned}$$

$$s_i(x) \approx p(y = i|x)$$

e.g. logistic regression, RF, DNN, ...

$$h(x) = \operatorname{argmax}_{i \in [K]} s_i$$

Standard Prediction Strategy

$$\max_h \langle A, C(h) \rangle \quad \left| \quad s_i(x) \approx p(y = i|x) \quad \left| \quad h(x) = \operatorname{argmax}_{i \in [K]} \mathbf{a}_i^\top \mathbf{s}(x) \right.$$

e.g. logistic regression, RF, DNN, ...

Proposed Postprocessing Strategy

A small experiment

$$\eta_k(x) \propto e^{\mathbf{w}_k^\top \mathbf{x}}$$

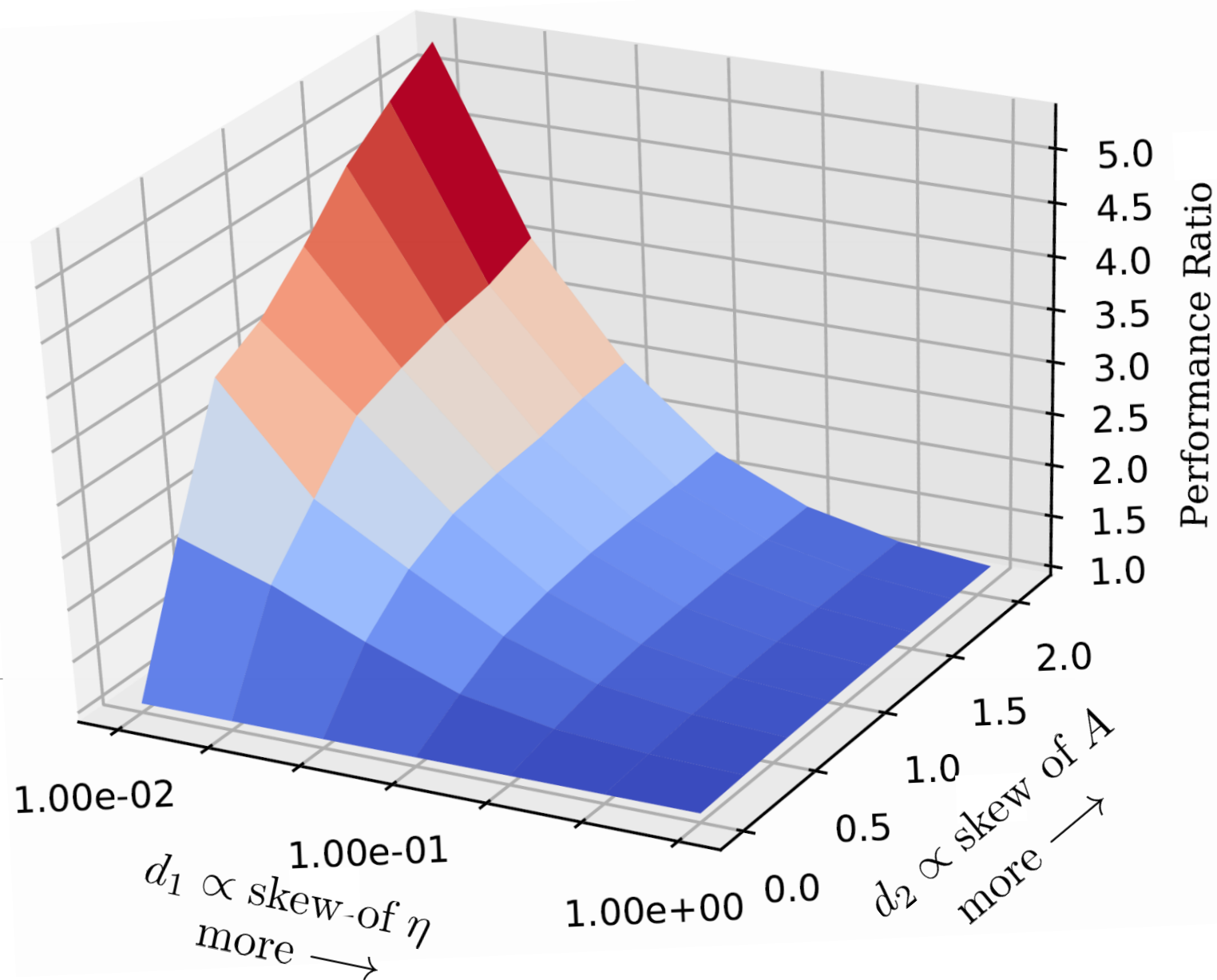
$$\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}); \quad \mathbf{w}_k = d_1 |k - K| \mathbf{1}$$

$$A_{j,j} = e^{-d_2 j}$$

1. Generate random data from model
2. Fit a logistic regression model
3. Post-process predictions

$$\text{Performance Ratio} = \frac{\text{Perf. of weighted postprocess}}{\text{Perf. of std. prediction}}$$

Simple
re-weighting
can have a
huge effect!



$$\max_h \psi(C(h))$$

$$s_i \approx p(y = i|x)$$

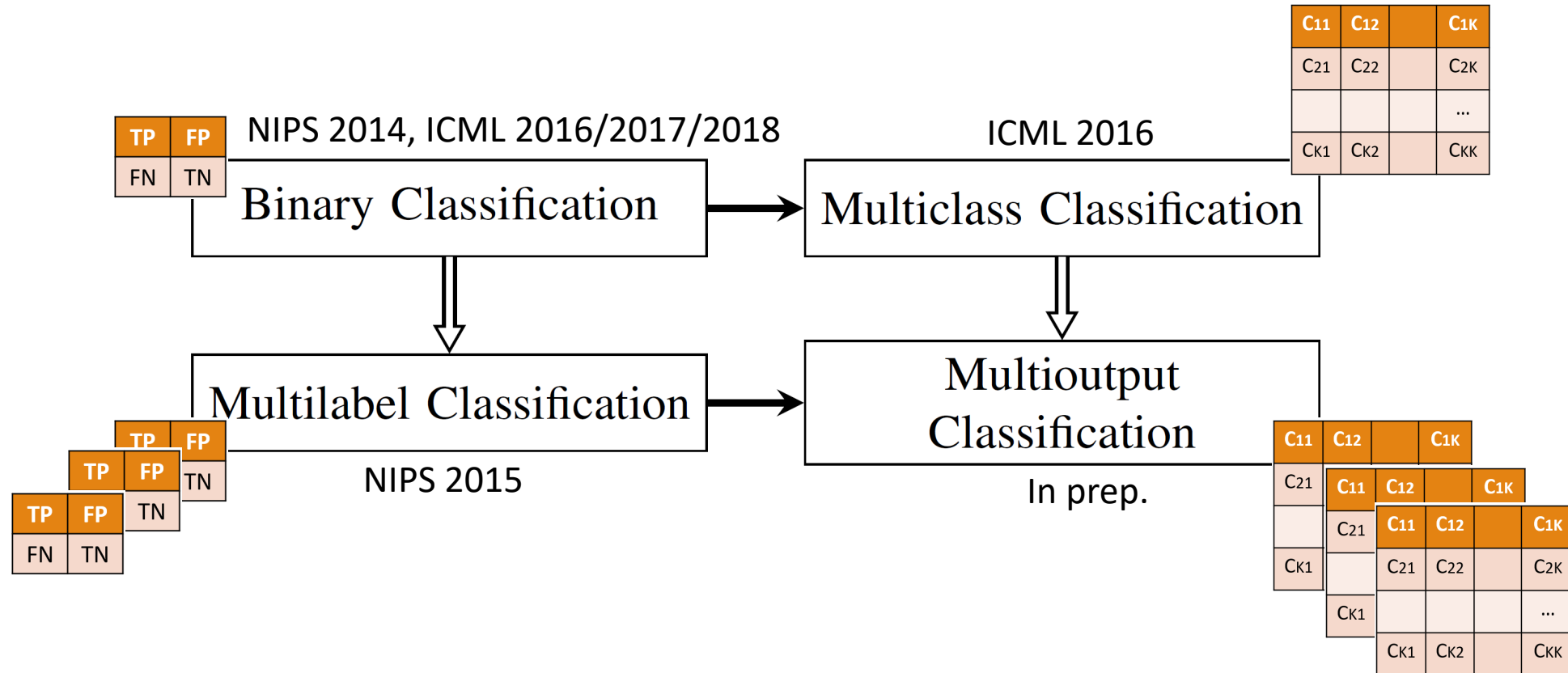
any calibrated classifier

$$h(x) = \operatorname{argmax}_{i \in [K]} \mathbf{b}_i^\top \mathbf{s}$$

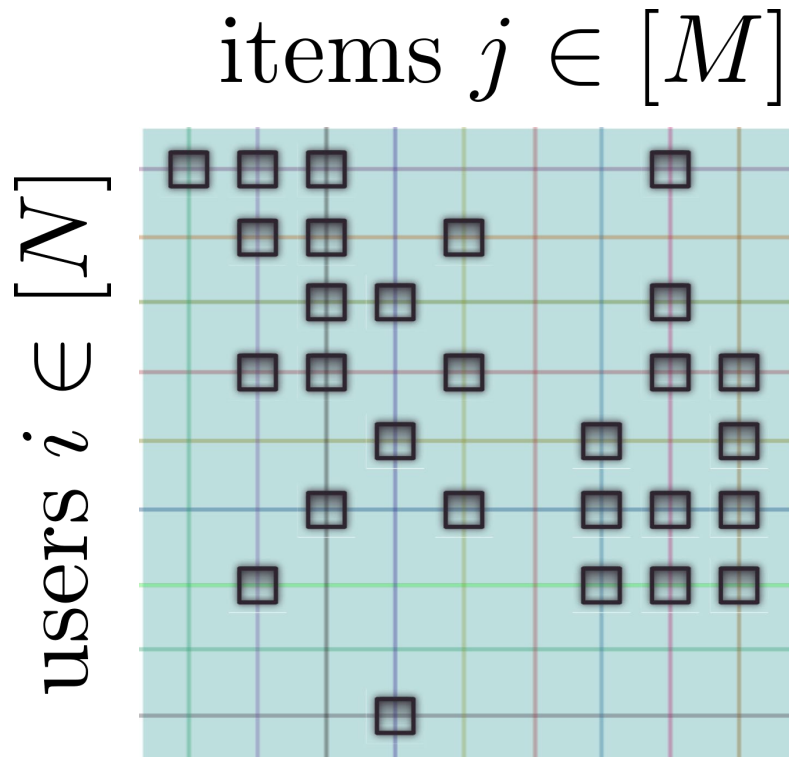
$$B = \nabla \psi|_{C=C^*}$$

Same strategy works for more complex metrics

Applies to more general settings



An application to recommender systems



User assigns rating to each item.

$$r_{i,j} \in [K]$$

Solve this as simultaneous (over items)
multiclass classification problem i.e.
multioutput classification

$$\phi(h) = \sum_{i=1}^K \sum_{j=1}^K |i - j| C_{i,j}$$

Postprocessed OrdRec

$$s_i(x) \approx p(y = i|x)$$

$$\operatorname{argmax}_{i \in [K]} s_i(x)$$

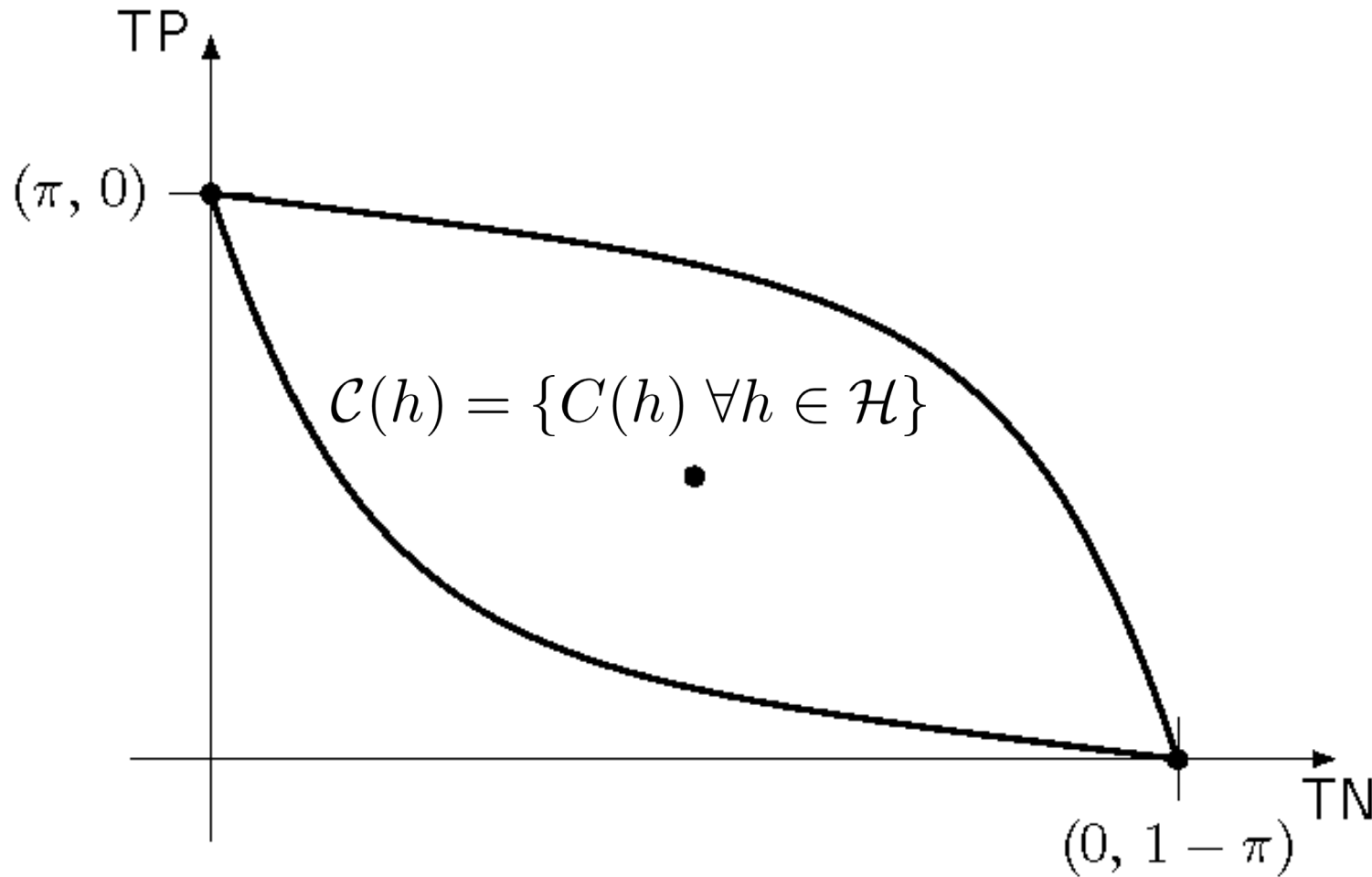
$$\operatorname{argmax}_{i \in [K]} \mathbf{a}_i^\top \mathbf{s}(x)$$

Koren, Yehuda, and Joe Sill.
"OrdRec: an ordinal model for
predicting personalized item
rating distributions." *Recsys*
2011.

AVERAGE	ORDREC	C-ORDREC
MICRO	0.8603 ± 0.0010	0.8640 ± 0.0009
MACRO	0.8577 ± 0.0032	0.8643 ± 0.0022
INSTANCE	0.8565 ± 0.0014	0.8619 ± 0.0011

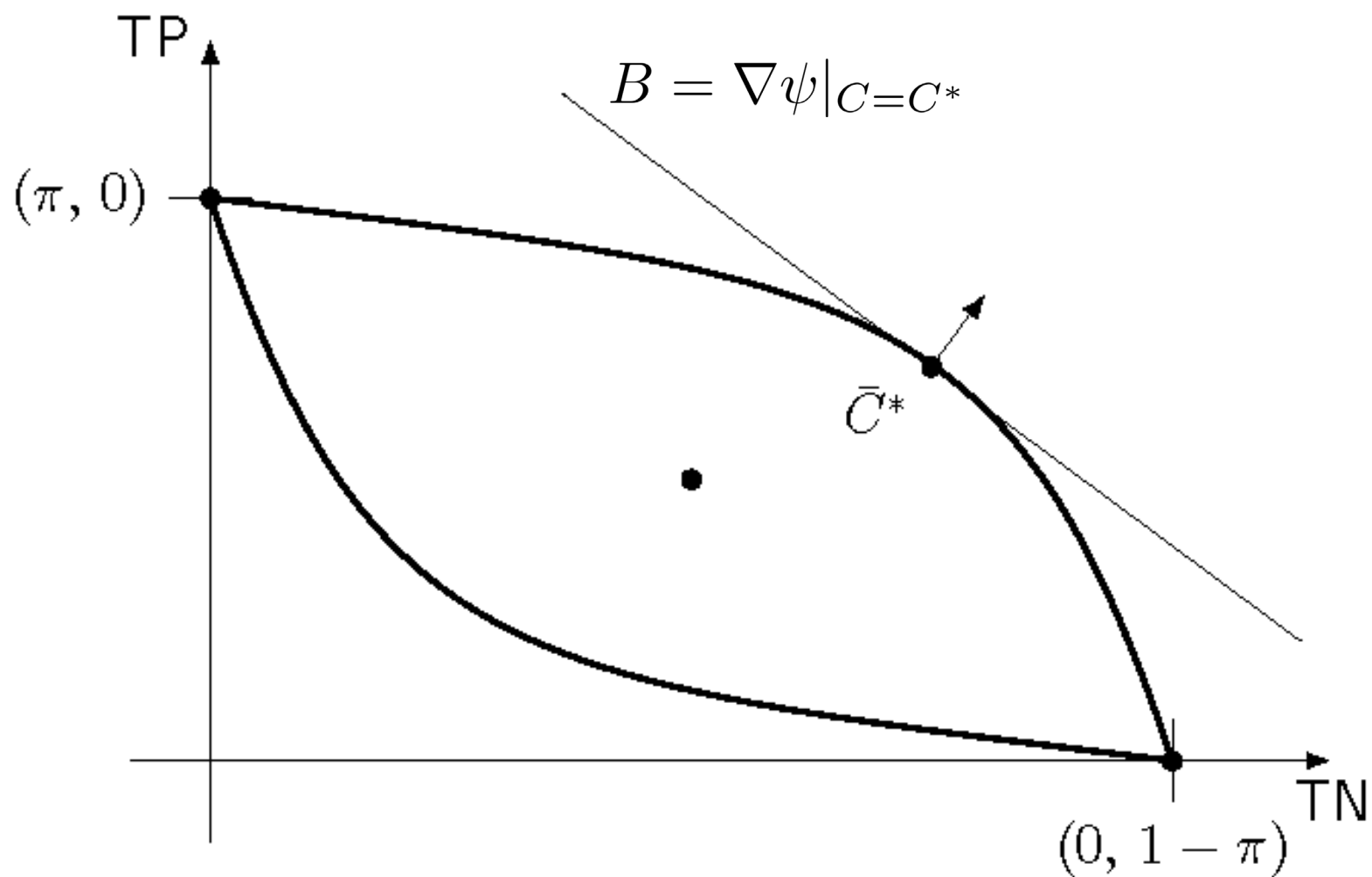
When & Why
does re-
weighting
work?

THE GEOMETRY OF
CONFUSION



$$TP + FN = \pi, \quad TN + FP = 1 - \pi$$

- Set of feasible confusion matrices is a bounded convex set
- Optimization properties will depend on how gradient field of the metric interacts with the feasible set
- Any monotonic metric will be optimized at the boundary



- All points on the boundary are determined by the support function
- This characterization is exhaustive i.e. characterizes ALL metrics that are consistently optimizable via linear post-processing

$$s_i(x) \rightarrow p(y = i|x)$$

$$B \rightarrow \nabla \psi|_{C=C^*}$$

$$\operatorname{argmax}_{i \in [K]} \mathbf{b}_i^\top \mathbf{s} \rightarrow h^*(x)$$

This classification strategy is consistent

Binary classification with general metrics

$$s(x) \approx p(y = i|x)$$

Logistic regression w/ MLE
Holder densities w/ kernel approx.

$$g_\delta(x) = \text{sign} \left(s(x) - \hat{\delta} \right)$$

Plug-in classifier

$$\hat{h}_n(x) = \operatorname{argmax}_{\delta \in [0,1]} \phi_n(g_\delta)$$

Threshold search

$$|\phi(h^*) - \phi(\hat{h}_n)| \leq O\left(\frac{\log n}{n}\right)$$

Which metric should you use?

THE BINARY CLASSIFICATION CASE

A solid orange horizontal bar at the bottom of the slide.

Recall: Lots of real world examples

$$\phi(h) = a_1 \text{TP} + a_2 \text{FP} + a_3 \text{FN} + a_4 \text{TN}$$

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \text{ TNR} = \frac{\text{TN}}{\text{FP} + \text{TN}}, \text{ Prec} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \text{ FNR} = \frac{\text{FN}}{\text{FN} + \text{TP}}, \text{ NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}}.$$

$$\text{AM} = \frac{1}{2} \left(\frac{\text{TP}}{\pi} + \frac{\text{TN}}{1 - \pi} \right) = \frac{(1 - \pi)\text{TP} + \pi\text{TN}}{2\pi(1 - \pi)}, \quad F_\beta = \frac{(1 + \beta^2)\text{TP}}{(1 + \beta^2)\text{TP} + \beta^2\text{FN} + \text{FP}} = \frac{(1 + \beta^2)\text{TP}}{\beta^2\pi + \gamma},$$

$$\text{JAC} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}} = \frac{\text{TP}}{\pi + \text{FP}} = \frac{\text{TP}}{\gamma + \text{FN}}, \quad \text{WA} = \frac{w_1\text{TP} + w_2\text{TN}}{w_1\text{TP} + w_2\text{TN} + w_3\text{FP} + w_4\text{FN}}.$$



Source: Fox news
Thanks to Willem Marai



► JAMES PANERO | THE NEW CRITERION EXECUTIVE EDITOR

IS THE METRIC SYSTEM COMPLETELY MADE UP?

5 OF 25 HOSPITAL PATIENTS WHO AUTHORITIES SAY WERE DELIBERATELY GIVEN OVERDOSES OF PAINKILLER



Limited formal guidance

Academia:

Use the standard metric in your application area

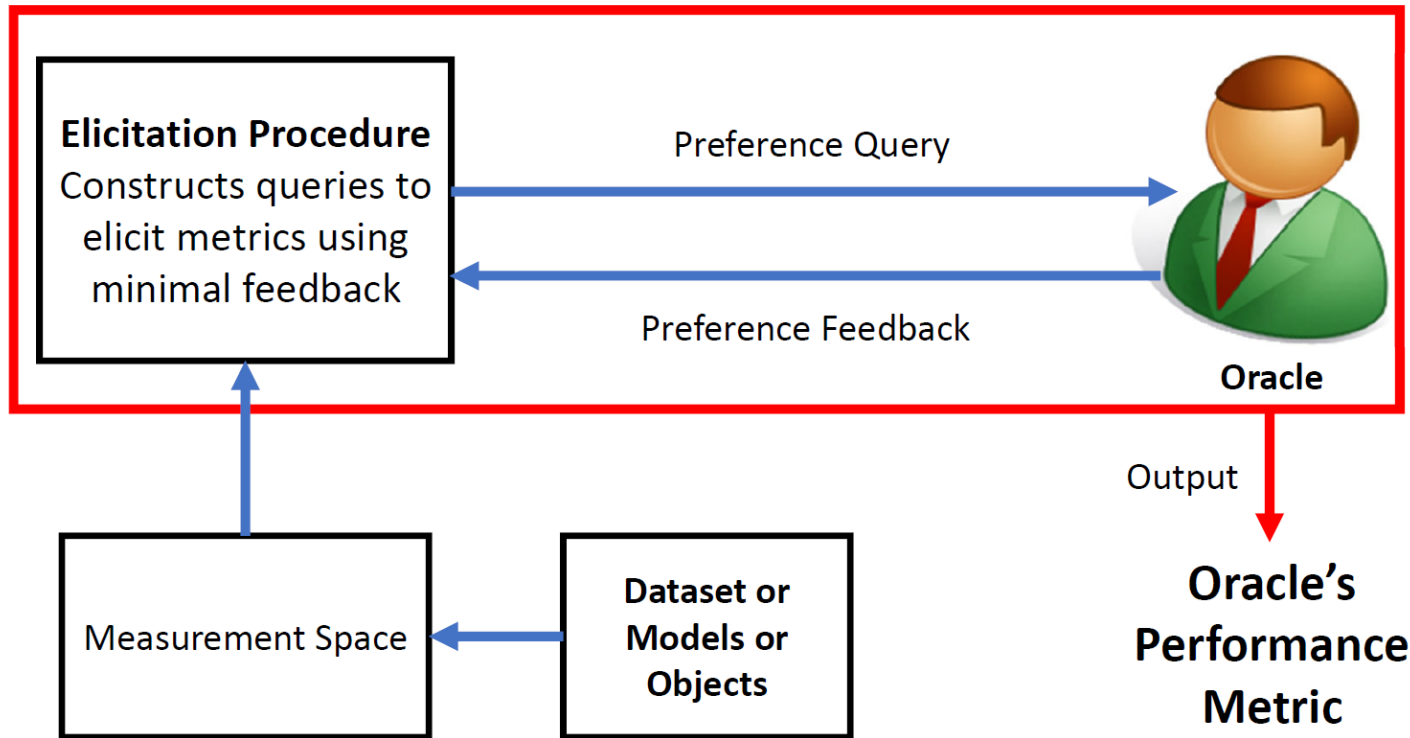
- Accuracy
- Top-K accuracy
- F1 measure

Industry:

Hire a consultant or economist

- User survey
- A/B tests





Metric Elicitation

QUERY AN EXPERT
(OR USER POPULATION)
TO QUANTIFY UTILITY

Querying the oracle

Humans (even domain experts) are often inaccurate when asked to quantify value

We're much more accurate when asked for pairwise preferences

$$\phi(\text{neural network diagram}) = ?$$

$$\phi(\text{neural network diagram})$$

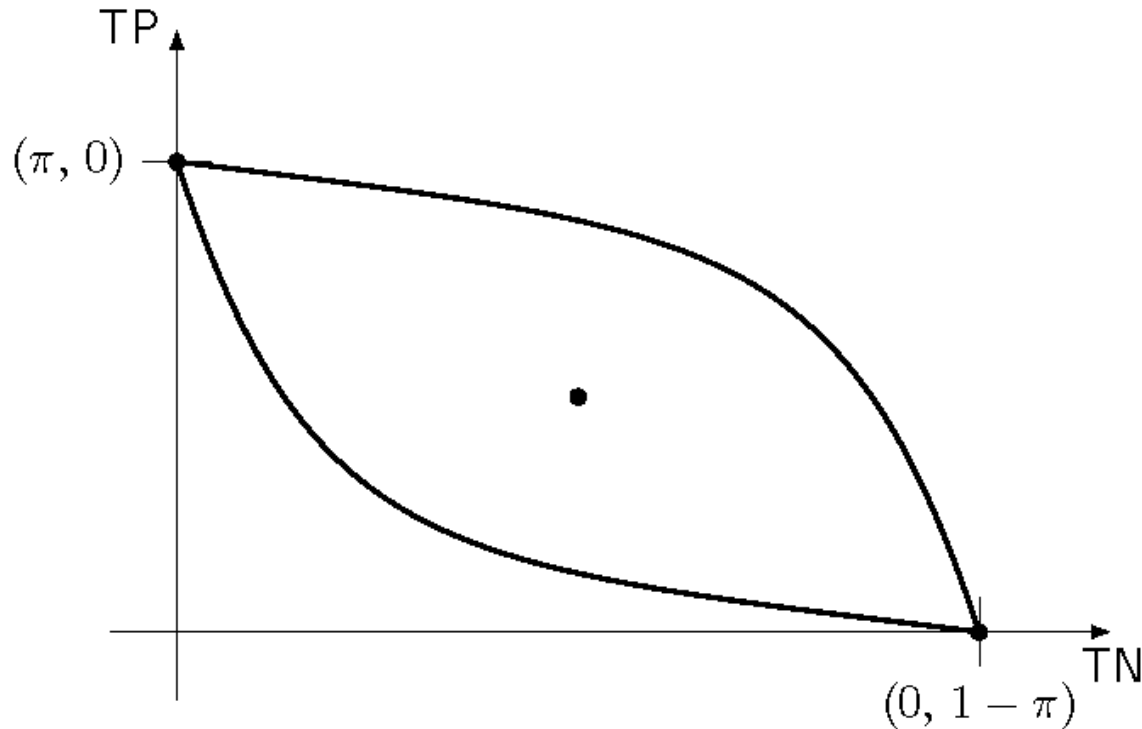
vs.

$$\phi(\text{2D plot with points and lines})?$$

TOO MANY QUERIES
MAY RESULT IN FATIGUE

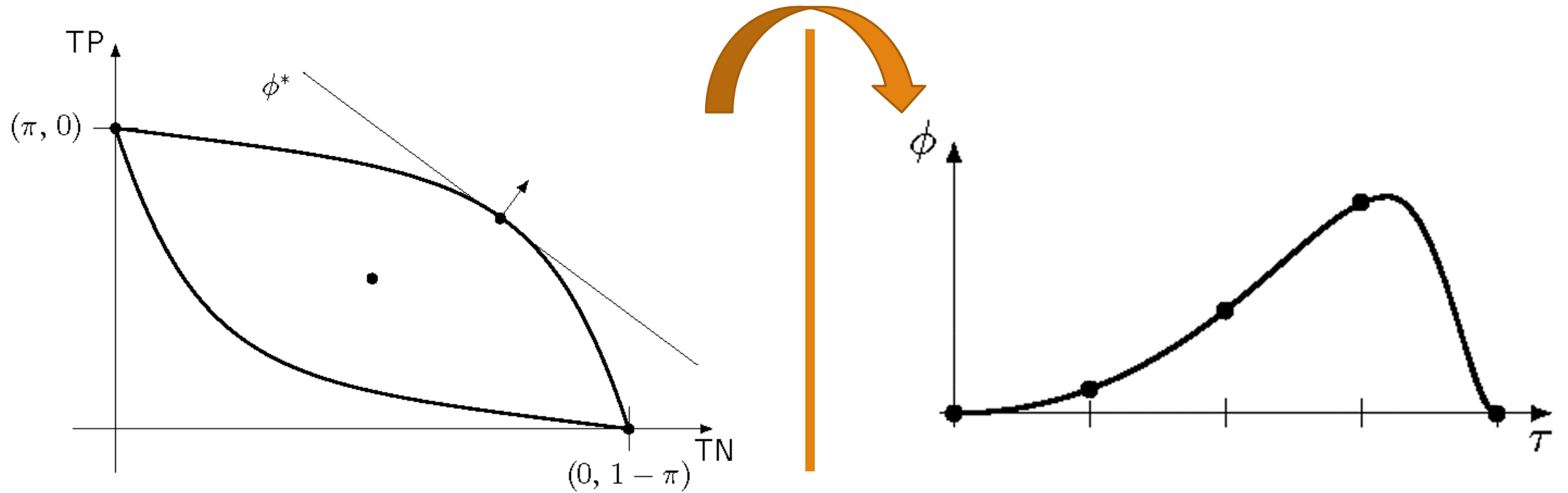
Goal:
accurately elicit
metric using
fewest pairwise
queries

Back to binary classification with weighted metrics



$$TP + FN = \pi, \quad TN + FP = 1 - \pi$$

- Set of feasible confusion matrices is a bounded convex set (related to the ROC curve)
- Weighted metrics correspond to linear functions
- Thus, eliciting weights requires finding a tangent to the boundary curve (equiv. point along the boundary) that recovers oracle tradeoff



unroll the boundary curve... find the maximum point

Metric recovery with finite queries

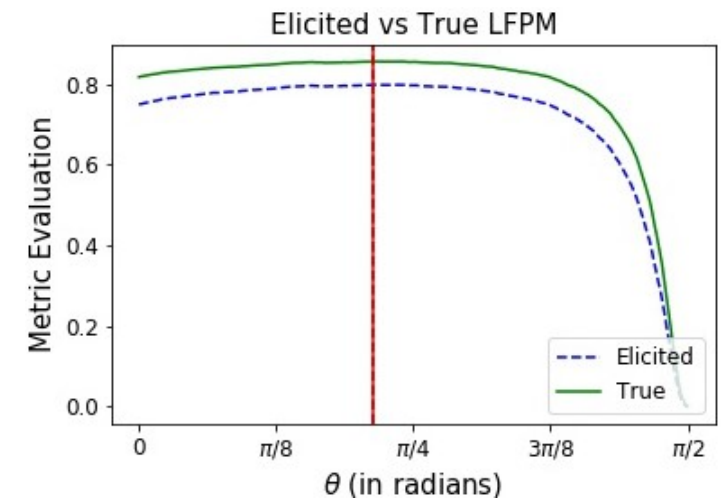
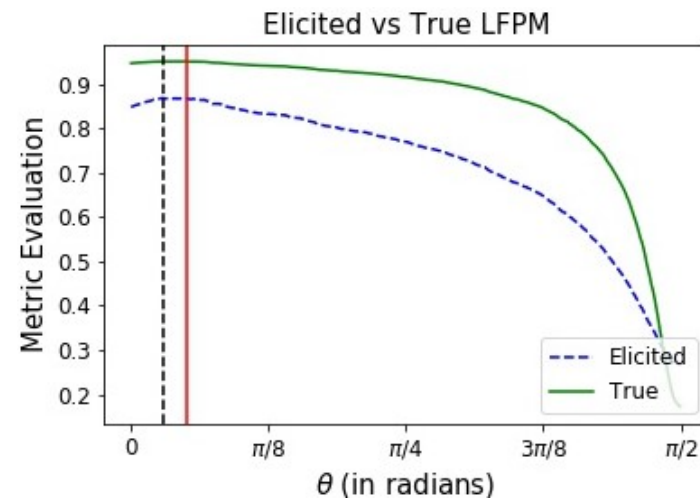
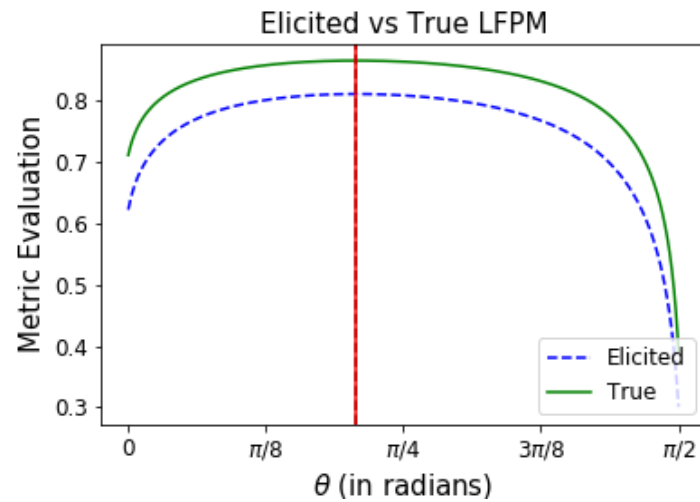
This algorithm provably recovers the oracle's weighted metric:

$$\phi^* \left(\text{Diagram} \right) = 1 - \left(a_1^* \text{FP} \left(\text{Diagram} \right) + a_2^* \text{FN} \left(\text{Diagram} \right) \right)$$

Guaranteed to be ϵ accurate after $\mathcal{O} \left(\log \left(\frac{1}{\epsilon} \right) \right)$ queries

Achieves the theoretical optimal elicitation rate

Stable to system noise e.g. noisy responses from the oracle



Empirical evaluation

Metric selection is a
crucial step in
constructing effective
machine learning
models

Simple algorithms can be used to
elicit an oracle metric (i.e. from an
expert or user population) using a
few queries

Conclusion

Metric choice has a large impact on real-world machine learning performance.

1

Re-weighted post-processing is efficient for optimizing complex metrics.

2

Can reduce metric elicitation for binary classifiers to binary search with bounded query complexity.

01

Extensions to other machine learning problems e.g. ranking, regression, probabilistic density modeling ...

02

Query mechanisms, HCI

03

Noise tolerance, robust elicitation

04

Applications to addressing bias and fairness in machine learning

Lots of work in progress

Thank you

QUESTIONS?

SANMI@ILLINOIS.EDU
@SANMIKOYEJO