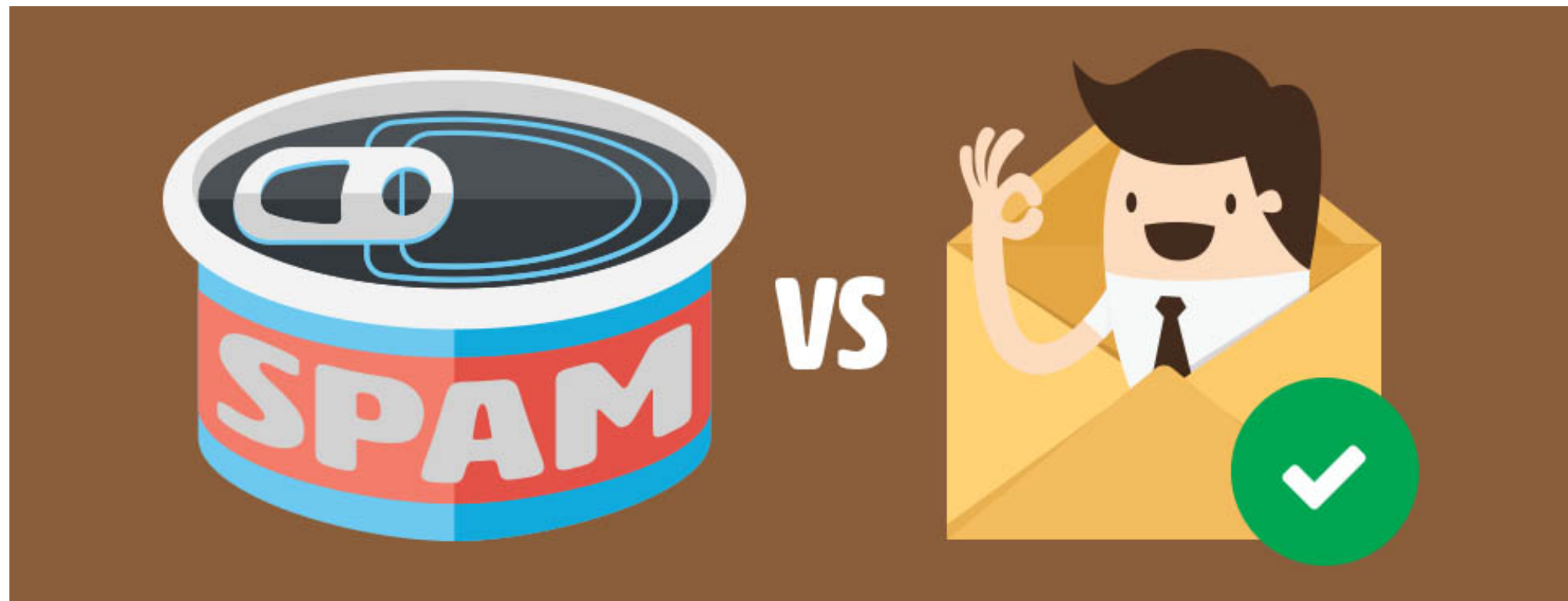
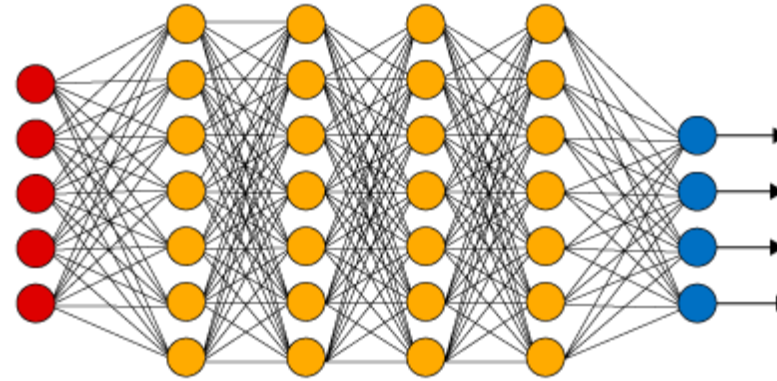


How effective is your classifier? Revisiting the role of metrics in machine learning

SANMI KOYEJO
CS @ ILLINOIS

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

Error analysis

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

- Accuracy = TP + TN = 1 - FP - 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	
		Y = 1	Y = 0
Predicted	h(x) = 1	TP	FP
	h(x) = 0	FN	TN

 = $\mathbf{C}(h)$

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 word "NETFLIX" is written in a bold, red, sans-serif font.The word "amazon" is written in a black, lowercase, sans-serif font with a yellow curved arrow underneath it.The word "Google" is written in its multi-colored, lowercase, sans-serif font.The word "Linked" is in black, lowercase, sans-serif font, followed by a blue square containing the white lowercase letters "in".

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.

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

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

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

Proposed Postprocessing Strategy

A small experiment

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

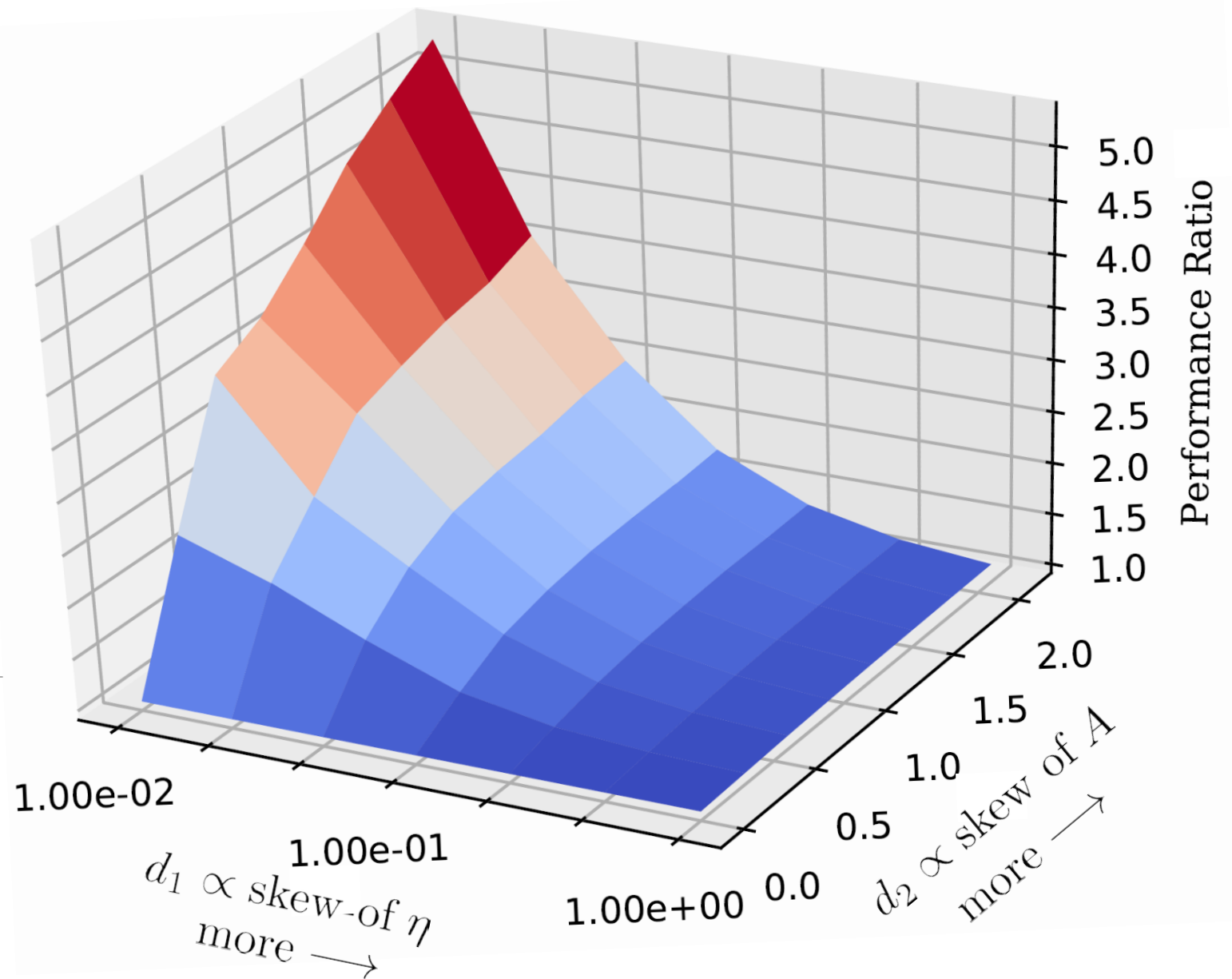
$$\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}); \quad \mathbf{w}_k = d_1 |k - K| \mathbf{1}$$

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

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!



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

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

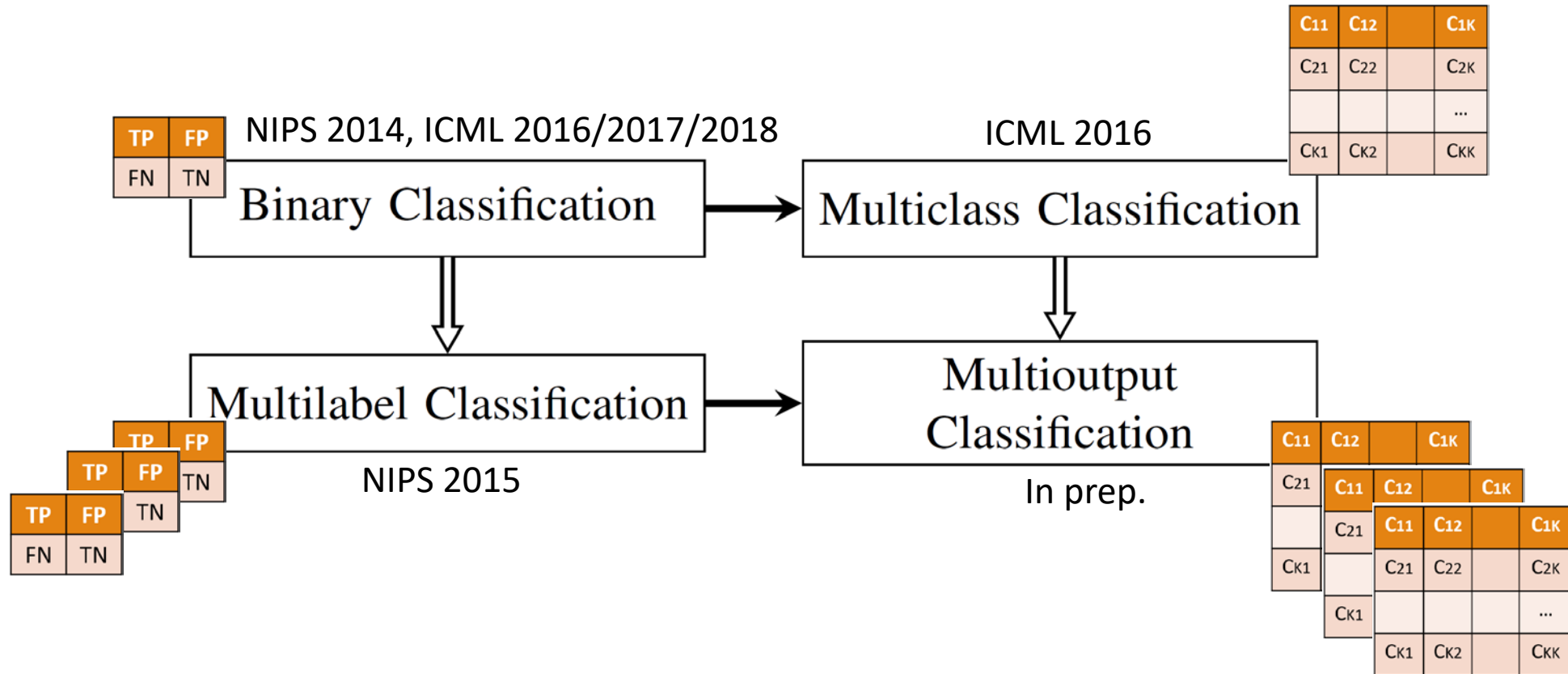
any calibrated classifier

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

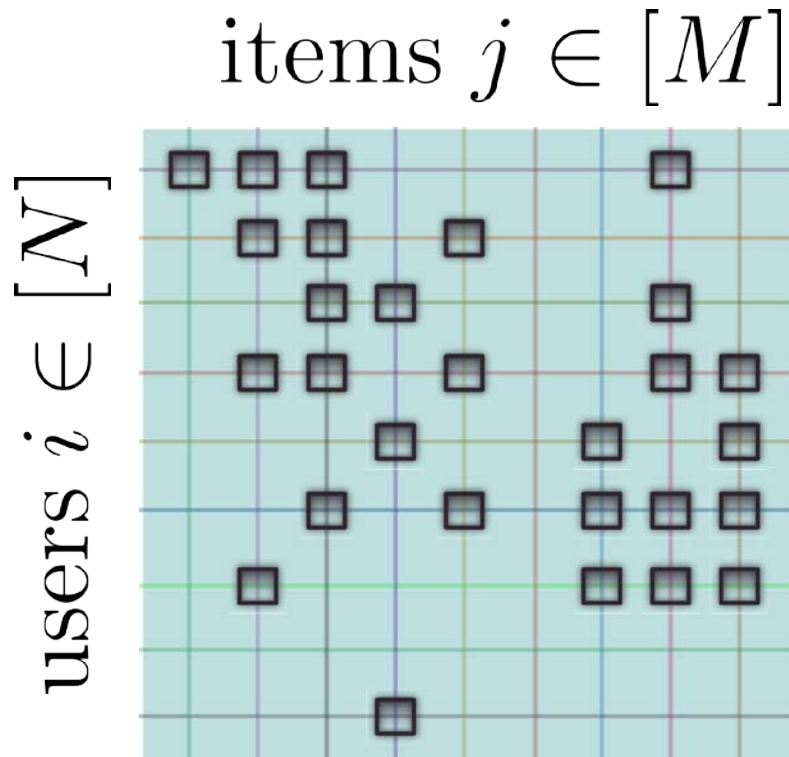
$$B = \nabla \psi|_{C=C^*}$$

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

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

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

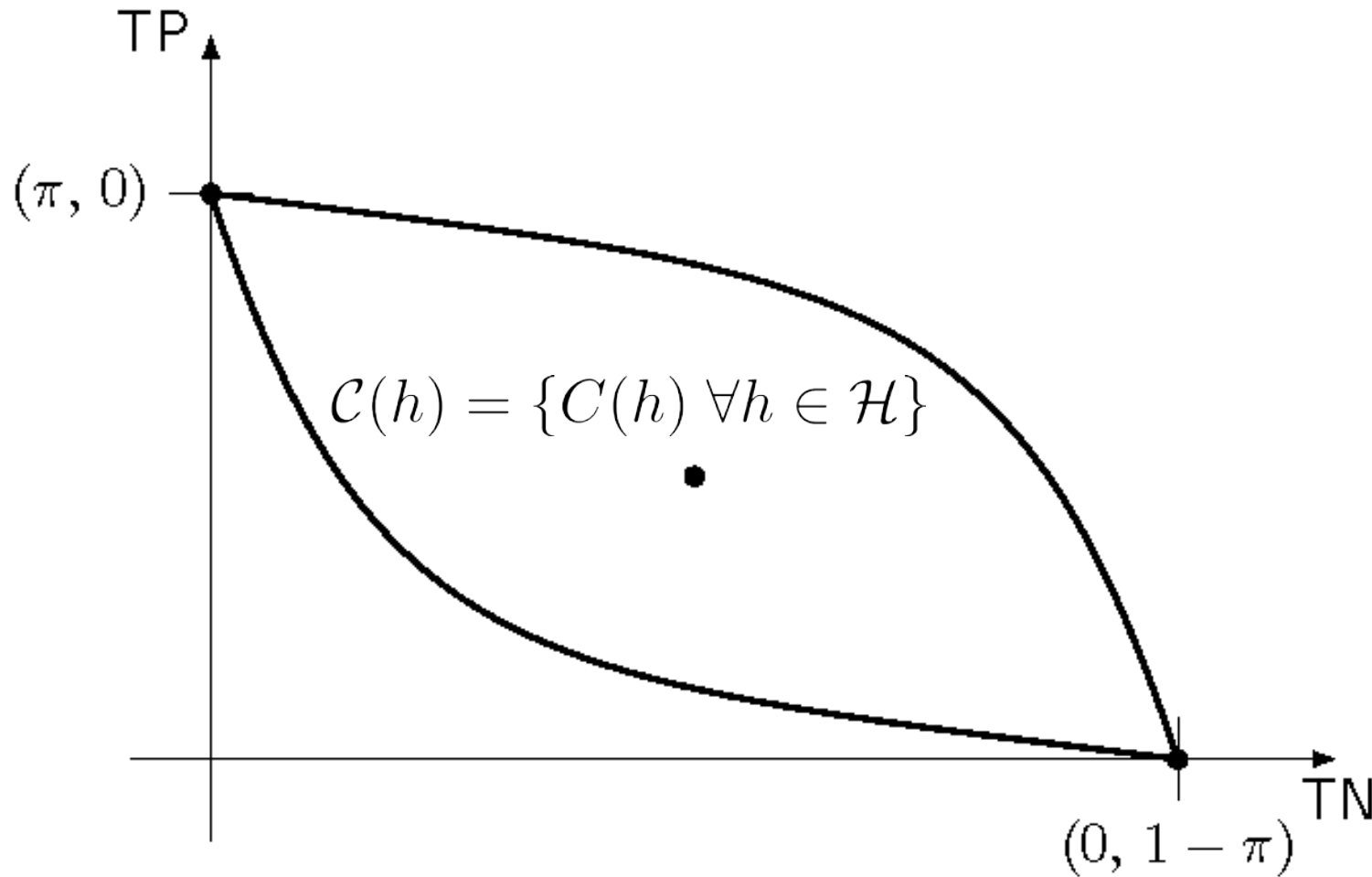
$$\operatorname{argmax}_{i \in [K]} \mathbf{a}_i^\top \mathbf{s}(x)$$

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

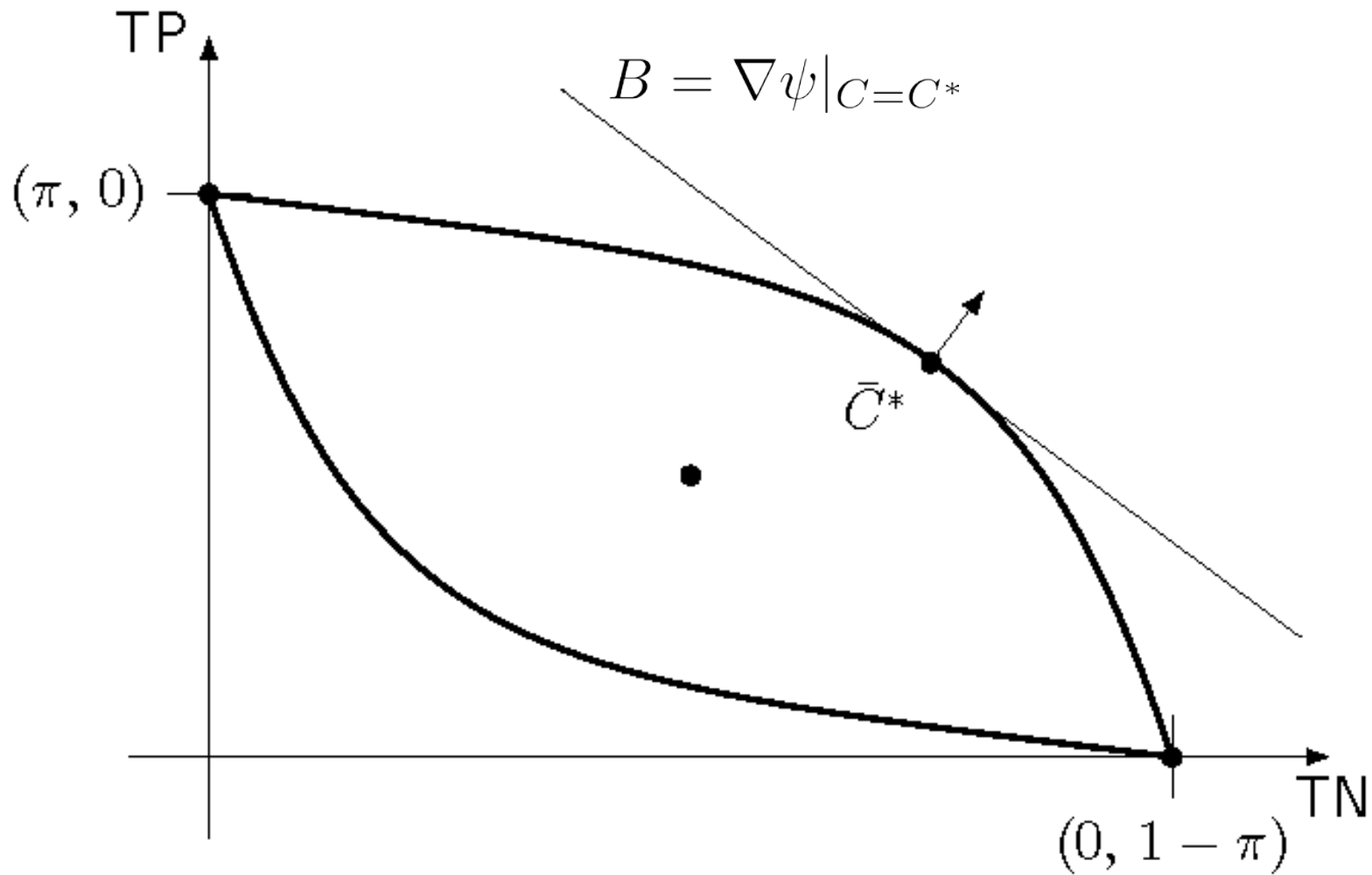
When & Why
does re-
weighting
work?

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

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

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



Recall: Lots of real world examples

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

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

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

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



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



Our Approach

Query an “expert” to determine the real-world value of a classifier i.e. the ideal evaluation metric

Pairwise queries

Experts give inaccurate results for value queries

More accurate results for comparison queries

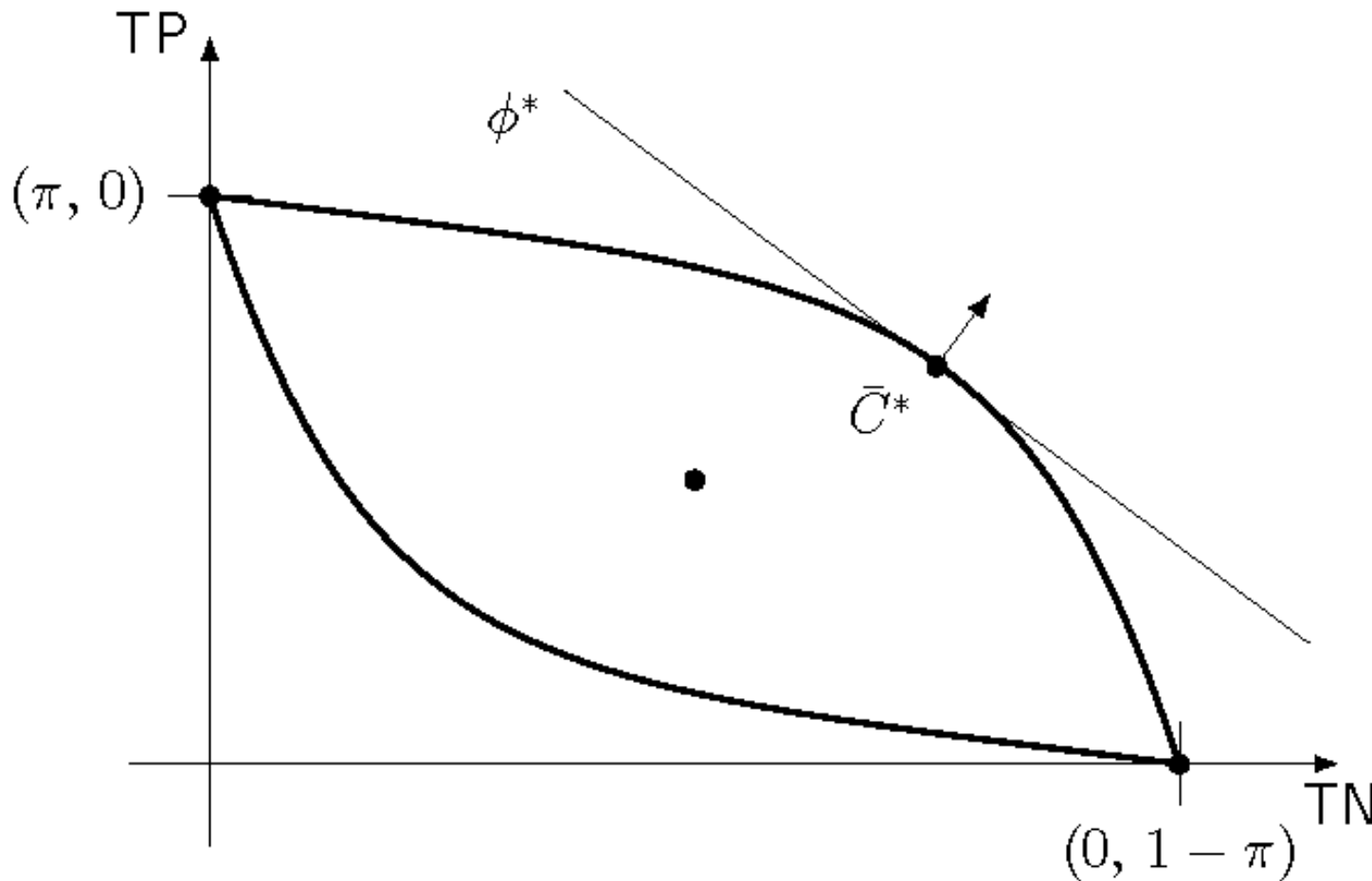
$$\phi(h) = ?$$

$$\phi(h_1) \text{ vs. } \phi(h_2) ?$$

THE "ORACLE" CARES
ABOUT WORST CASE
QUERY COMPLEXITY

Speed Matters!

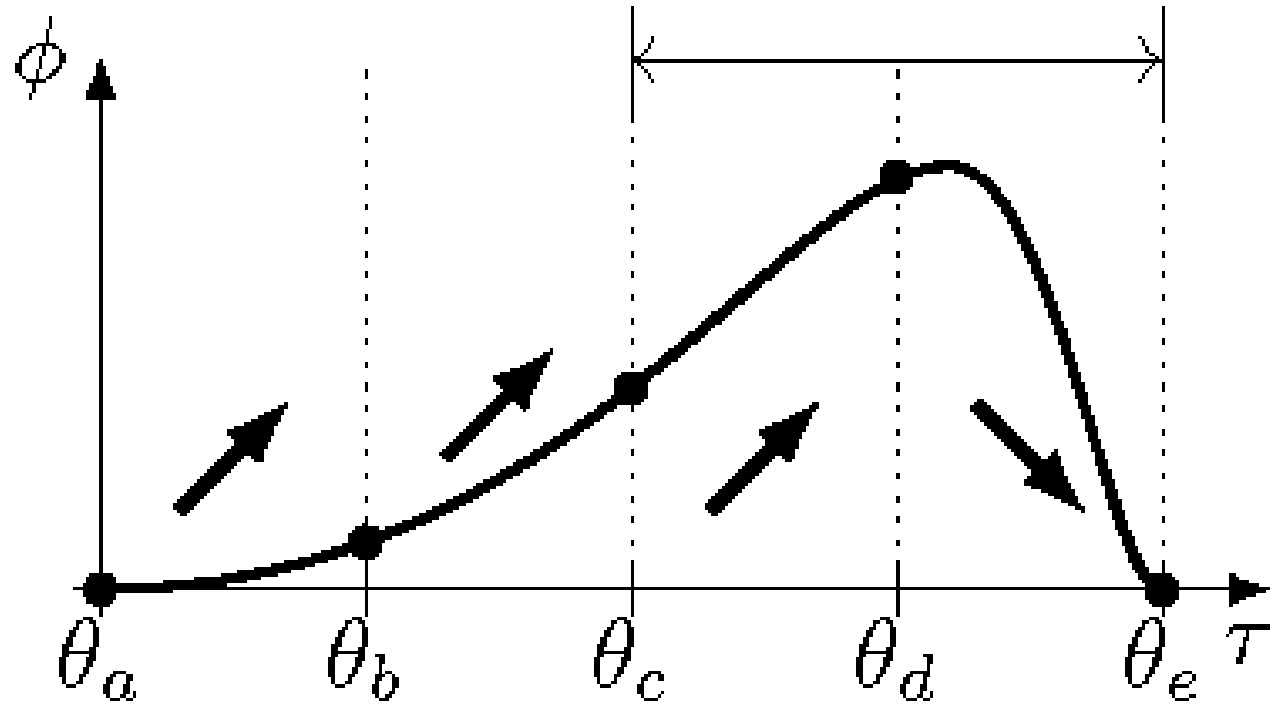
Exploiting the geometry...



- Only need to query classifiers on the boundary – since we already know optimal is within this subset
- Boundary is one-dimensional, parameterized by “angle”

Using binary search

- Under weak conditions, metric is unimodal with respect to boundary
- Thus, can simple binary search to find the optimal confusion matrix
- Simultaneously recovers gradient of the optimal metric



Guaranteed recovery with finite queries

For the linear case, when algorithm terminates, we recover

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

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

If no additional assumptions, this matches lower bound

Stable to system noise e.g. noisy responses from the “expert”

Conclusion

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

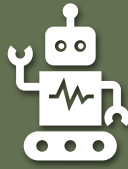
1

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.

Measurement is
at the core of
empirical
research



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



Faster elicitation using alternative query mechanisms



Noise tolerance, robust elicitation

Thank you

QUESTIONS?

SANMI@ILLINOIS.EDU