# Knowing what you don't know:
## Learning to abstain and beyond

Harikrishna Narasimhan

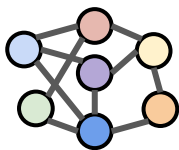Google Research

# Contributors
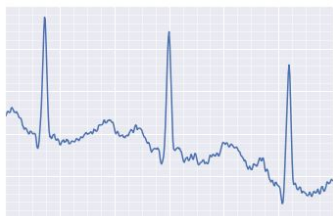

Wittawat
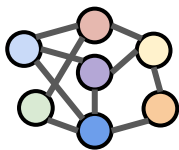Jitkrittum


Ankit S.
Rawat


Aditya K.
Menon


Sanjiv
Kumar

Google

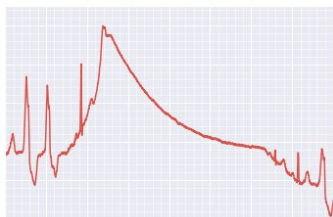# Standard classification paradigm

- Standard classification → **single model** for all samples

- However, it may be challenging to model the entire input space



Input sample                    Decision maker                    Prediction

# Learning to reject

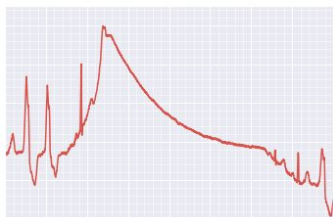- Model can **give up** on a sample, incurring some **cost**



Input sample                Decision maker                Prediction

# Learning to defer to an expert

- Model can **defer** to an **expert**, incurring some **cost**
  - e.g., human expert



Input sample          Decision maker          **Expert / Specialist**          Prediction

# Learning to defer to an expert

- Model can **defer** to an **expert**, incurring some **cost**
  - e.g., human expert, powerful learning model



Input sample          Decision maker          **Expert / Specialist**          Prediction

# Learning to abstain on outliers

- Model can abstain on samples it deems to be **out-of-distribution (OOD)**



Input sample          Decision maker          Prediction

**Goal**: learn the base classifier, **and** the abstention rule

Google

# Cost of abstention: *classical version*

We will denote a joint classifier $h$: $X \rightarrow$ [n] $\cup$ { 🙇 }. In the simplest case, one may associate a constant cost $c$ to abstaining on a sample

$$\mathbb{1}\left(\hat{y} \neq y, \hat{y} \neq 🙇\right) \; + \; c \cdot \mathbb{1}\left(\hat{y} = 🙇\right)$$

Usual error when not abstaining

Constant cost when abstaining

# Chow's rule: a surprisingly competitive baseline

C. Chow. On optimum recognition error and reject tradeoff.
*IEEE Transactions on Information Theory*, 16(1):41−46, 1970.

Bayes-optimal rejection rule: abstain on a sample when

$$\max_{y} \mathbb{P}(y \mid x) < 1 - c$$

# Chow's rule: a surprisingly competitive baseline

C. Chow. On optimum recognition error and reject tradeoff.
*IEEE Transactions on Information Theory*, 16(1):41−46, 1970.

Bayes-optimal rejection rule: abstain on a sample when

$$\max_{y} \widehat{\mathbb{P}}(y \mid x) \; < \; 1 - c$$

In practice: max softmax probability
from a standard classifier

# When Chow's rule fails and ways to remedy it!

- Learning to reject
  - classical Chow's rule is very competitive

- Learning to defer to an expert
  - remedy: expert-aware Chow's rule

- Learning to abstain on outliers
  - remedy: outlier-aware Chow's rule

Google

# Cost of abstention: *when deferring to an expert*

In the learning to defer paradigm, the cost of invoking the expert:

$$\mathbb{1}\big(\hat{y} \neq y,\ \hat{y} \neq 🤷\big)\ +\ c_{\mathrm{exp}}(x, y) \cdot \mathbb{1}\big(\hat{y} = 🤷\big)$$

Usual error when not abstaining

Cost of invoking the expert

# Expert cost: fixed cost + expert's error rate

A natural candidate for the expert cost would include both a fixed cost and the penalty when the expert makes a mistake

$$c_{\exp}(x, y) = c_0 + \mathbf{1}(y \neq h_{\exp}(x))$$

Fixed cost    Expert prediction

(e.g. monetary cost)

Google

# Chow's rule can be sub-optimal for this setting



Degrades with more abstentions!

Synthetic dataset
Base model: linear features
Expert model: quadratic features

# Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

$$\max_{y} \mathbb{P}(y \mid x) \; < \; \mathbb{E}_{y\mid x}\big[\mathbf{1}(y = h_{\exp}(x))\big] \; - \; c_0$$

~ Base classifier's confidence

~ Expert's confidence

# Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

When the expert's confidence is highly non-uniform, this is substantially different from Chow's rule

$$\max_y \mathbb{P}(y \mid x) \; < \; \mathbb{E}_{y|x}[\mathbf{1}(y = h_{\exp}(x))] \; - \; c_0$$

~ Base classifier's confidence

~ Expert's confidence

# Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

$$\max_y \widehat{\mathbb{P}}(y \mid x) \; < \; \widehat{\mathbb{E}}_{y|x}[\mathbf{1}(y = h_{\exp}(x))] \; - \; c_0$$

~ Base classifier's confidence

~ Expert's confidence

Unlike classical Chow, we need to estimate the expert's confidence

# Separate model for expert's confidence

Raghu et al. '19 suggest training separate model to estimate expert's confidence (using a sample annotated with the expert's predictions)

$$\max_y \widehat{\mathbb{P}}(y \mid x) \; < \; \widehat{\mathbb{E}}_{y|x}[\mathbf{1}(y = h_{\exp}(x))] \; - \; c_0$$

~Softmax probabilities from base classifier

~ Separate model to estimate expert's confidence

# Separate model for expert's confidence

- This approach has appealing properties:

    ✓  Simple to compute

    ✓  Approximates the Bayes deferral rule

    !  Separate models to estimate base and expert confidence

# Cost-sensitive softmax cross-entropy (CSS)

- Mozannar & Sontag '20 suggest training a joint model with an additional label $\perp$

# Cost-sensitive softmax cross-entropy (CSS)

- Mozannar & Sontag '20 suggest training a joint model with an additional label ⊥

- Minimize a **cost-sensitive** softmax cross-entropy (**CSS**) loss

$$\ell_{\mathrm{css}}(x, y, \bar{f}(x)) = -\log\left(\bar{p}_y(x)\right) - \mathbf{1}(y = h_{\exp}(x)) \cdot \log\left(\bar{p}_\perp(x)\right) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

Classification loss
to train base model

Loss to estimate
expert's confidence

Takes into account
fixed cost $c_0$

*The original loss of Mozannar and Sontag uses a slightly tighter formulation; see our paper for details

# The case for CSS

- The CSS loss has a number of appealing characteristics:

  ✓ **Joint model** for both base classifier and expert's confidence

  ✓ Optimal solution matches the **Bayes-optimal classifier**

  ✓ **Empirically effective** on several benchmarks

  ❗ when fixed cost $c_0 = 0$...

# The case against CSS?

- CSS strongly **underfits** when there is non-zero fixed deferral cost $c_0$!



CIFAR 100
ResNet8 base
ResNet32 expert

# A label smoothing perspective

- CSS equivalently applies high level of **label smoothing**:

$$\ell_{\mathrm{css}}(x, y, \bar{f}(x)) = -\log\left(\bar{p}_y(x)\right) - \mathbf{1}(y = h_{\exp}(x)) \cdot \log\left(\bar{p}_\perp(x)\right) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

- Encourages predictions to become highly uniform
- Low separation between true label and competing labels

Treat **all** labels as candidate positive

Google

# A label smoothing perspective

- CSS equivalently applies high level of **label smoothing**:

$$\ell_{\text{css}}(x, y, \bar{f}(x)) = -\log\left(\bar{p}_y(x)\right) - \mathbf{1}(y = h_{\exp}(x)) \cdot \log\left(\bar{p}_\perp(x)\right) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

  - Encourages predictions to become highly uniform

  - Low separation between true label and competing labels

Treat **all** labels as candidate positive

- Not apparent when $c_0$ = 0 (as in prior work)!

  - $c_0 > 0$ is crucial in practical settings (e.g. when the expert is a larger model)

Google

# Solution: Set $c_0 = 0$ during training; include it in a post-hoc step

- Train base model with $c_0 = 0$, i.e., by minimizing:

$$\ell_{\text{css}}(x, y, \bar{f}(x)) = -\log\left(\bar{p}_y(x)\right) - \mathbf{1}(y = h_{\exp}(x)) \cdot \log\left(\bar{p}_\perp(x)\right) - \cancel{c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))}$$

| $\bar{p}_1$ | $\bar{p}_2$ | ... | $\bar{p}_L$ |
|---|---|---|---|

$\boxed{\bar{p}_\perp}$

Class probabilities

Probability that the expert is correct

Google

# Solution: Set $c_0 = 0$ during training; include it in a post-hoc step

Construct a **post-hoc rejector** to include $c_0$ (that mimics the Bayes-optimal rule):

$$\max_y \overline{p}_y(x) < \overline{p}_\perp(x) - c_0$$

Probability that the
expert is correct

Deferral cost

Google

# **Proposal: two-step plug-in approach** [Narasimhan et al '22]



Training sample → **Base model training**: Minimise joint `expert-aware` loss with $c_0 = 0$ → Logits $\bar{f}_1, \ldots, \bar{f}_L$ / Defer logit $\bar{f}_\perp$ → **Compute statistics**: Class probabilities $\bar{p}_y(x)$, Expert confidence $\bar{p}_\perp(x)$ → $\bar{p}_y(x), \bar{p}_\perp(x)$ → **Final classifier**: Use rejection rule: $\max_y \bar{p}_y(x) < \bar{p}_\perp(x) - c_0$ for any $c_0$

# Experimental setup

- **Specialist** expert
  - Model allowed to defer to a "specialist" expert trained on a subset of labels
- Baselines

  - **Chow**: confidence thresholding based only on the deferral cost $c_0$
  - **CSS**: in-training loss of Mozannar & Sontag (2020) with $c_0$ included
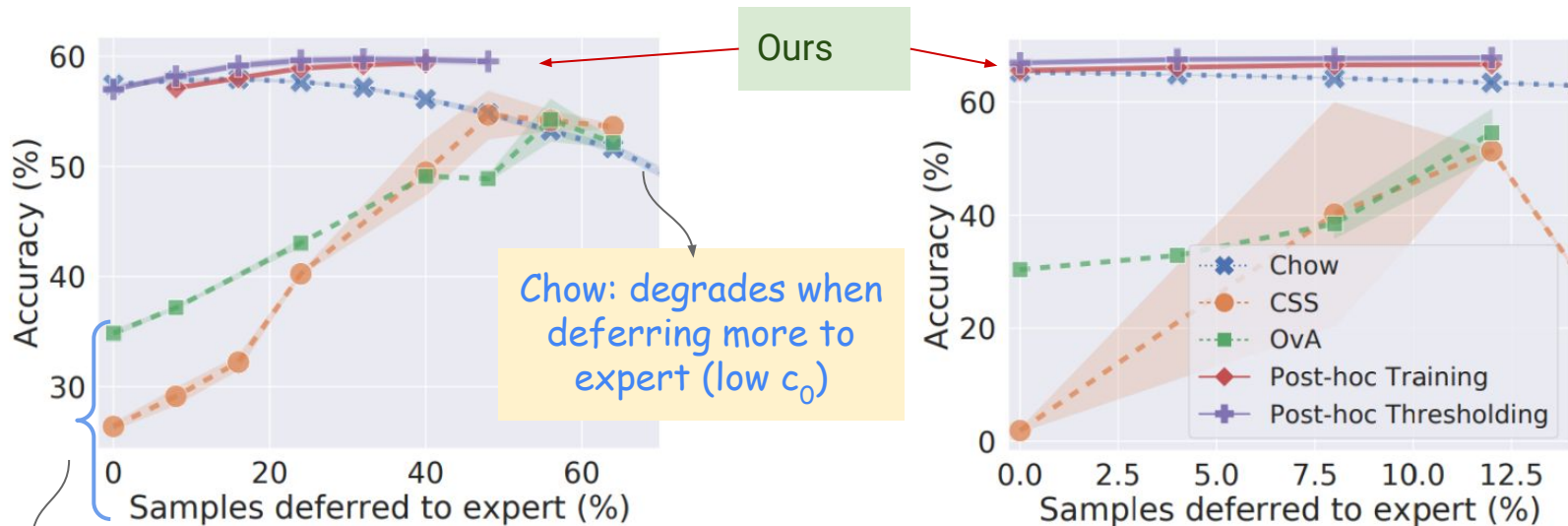  - **OvA**: in-training loss of Verma and Nalisnick (2022) with $c_0$ included

Ignores expert error

Underfits when $c_0$ is large

Google

# Experimental results: expert-aware abstention



Ours

Chow: degrades when deferring more to expert (low $c_0$)

Joint losses: degrade when deferring less to expert (high $c_0$)

CIFAR 100
ResNet8 base
ResNet56 expert
(expert trained on first 10 coarse labels)

ImageNet
MobileNet-v2 base
EfficientNet-B0 expert
(expert trained on "dog" synset)

Legend:
- Chow
- CSS
- OvA
- Post-hoc Training
- Post-hoc Thresholding

# When Chow's rule fails and ways to remedy it!

- Learning to reject
  - classical Chow's rule is very competitive

- Learning to defer to an expert
  - remedy: expert-aware Chow's rule

- **Learning to abstain on outliers**
  - remedy: outlier-aware Chow's rule

# Learning to abstain on outliers

Abstain on "out-of-distribution" samples that come from distribution different from the one used for training



**Inlier samples**

**Outlier samples**

Google

# Chow's rule (or the MSP scorer) is a popular baseline!

Thresholding the maximum softmax probability (MSP) from a standard classifier is a common baseline in this literature [Hendrycks et al. '17; Vaze et al. '22].

$$\max_{y \in [L]} \widehat{\mathbb{P}}_{\text{in}}(y \mid x) \; < \; t$$

# Chow's rule can fail for outlier detection



$P_{in}(y=1|x) \approx 0.5$: Chow's rule will **abstain** on these, despite them being inlier

$P_{in}(y=1|x) \approx 1$: Chow's rule will **not abstain** on these, despite them being outliers.

Inlier samples

Outlier samples

# Cost of abstention: *when abstaining on outliers*

We need to account for both inlier **and** outlier abstentions.

$$\mathbb{P}_{\text{in}}\left(\hat{y} \neq y,\ \hat{y} \neq \text{🤐}\right) \ +\ \alpha \cdot \mathbb{P}_{\text{in}}\left(\hat{y} = \text{🤐}\right) \ +\ \beta \cdot \mathbb{P}_{\text{out}}\left(\hat{y} \neq \text{🤐}\right)$$

Error on inlier samples
(when not abstaining)

Cost of abstaining
on inlier samples

Cost of *not* abstaining
on outlier samples

# Outlier–aware Chow's rule

Bayes-optimal rule: abstain on a sample when [Narasimhan et al. '23]

$$\max_{y} \mathbb{P}_{\text{in}}(y \mid x) \; < \; 1 - \alpha + \beta \cdot \frac{\mathbb{P}_{\text{out}}(x)}{\mathbb{P}_{\text{in}}(x)}$$

Inlier class
probabilities

Outlier-to-inlier
density ratio

# Outlier–aware Chow's rule

Bayes-optimal rule: abstain on a sample when [Narasimhan et al. '23]

$$\max_{y} \widehat{\mathbb{P}}_{\text{in}}(y \mid x) \; < \; 1 - \alpha + \beta \cdot \frac{\widehat{\mathbb{P}_{\text{out}}(x)}}{\mathbb{P}_{\text{in}}(x)}$$
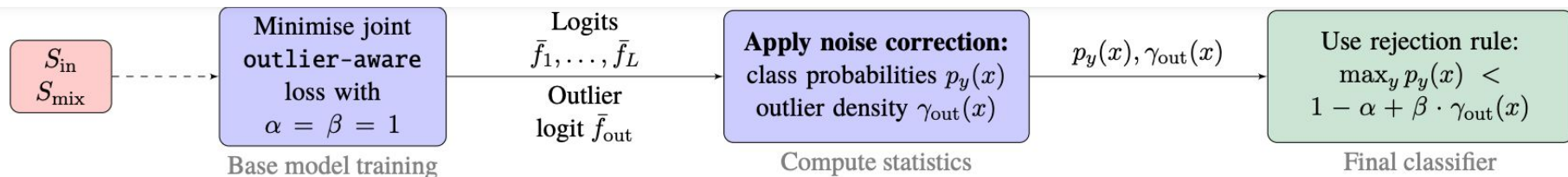
Inlier class
probabilities

Outlier-to-inlier
density ratio

We need to estimate both the inlier probabilities and the density ratio

# Our proposal: Two-step plug-in approach [Narasimhan et al '23]

Given: labeled inlier sample $S_{\text{in}}$, and an unlabeled mix of inlier and outlier samples $S_{\text{mix}}$



$S_{\text{in}}$
$S_{\text{mix}}$

Minimise joint **outlier-aware** loss with $\alpha = \beta = 1$
Base model training

Logits $\bar{f}_1, \ldots, \bar{f}_L$
Outlier logit $\bar{f}_{\text{out}}$

**Apply noise correction:** class probabilities $p_y(x)$ outlier density $\gamma_{\text{out}}(x)$
Compute statistics

$p_y(x), \gamma_{\text{out}}(x)$

Use rejection rule: $\max_y p_y(x) < 1 - \alpha + \beta \cdot \gamma_{\text{out}}(x)$
Final classifier

# Experimental results: outlier-aware abstention

# When Chow's rule fails and ways to remedy it!

- ## Learning to reject (L2R)
  - Classical Chow's rule is very competitive

- ## Learning to defer to an expert (L2D)
  - Chow may fail; use *expert-aware* Chow

- ## Learning to abstain on outliers (OOD)
  - Chow may fail; use *outlier-aware* Chow

| L2R | $\max_{y} \mathbb{P}(y|x) \quad < \quad 1 - c$ |
|-----|-----------------------------------------------|
| L2D | $\max_{y} \mathbb{P}(y|x) \quad < \quad \mathbb{E}_{y|x}[\mathbf{1}(y = h_{\exp}(x))] - c_0$ |
| OOD | $\max_{y} \mathbb{P}_{\text{in}}(y|x) < \quad 1 - \alpha + \beta \cdot \dfrac{\mathbb{P}_{\text{out}}(x)}{\mathbb{P}_{\text{in}}(x)}$ |

Narasimhan et al. "Post-hoc Estimators for Learning to Defer to an Expert". NeurIPS 2022.

Narasimhan et al. "Learning to Reject Meets OOD Detection: Are All Abstentions Created Equal?". Manuscript, 2023. [arXiv:2301.12386]

**Thank you!**

Google