Google DeepMind

Conformal prediction under ambiguous ground truth

David Stutz November 6th 2023

inter-observer variability

benign

Motivation: Ambiguity in Classification

- High-stakes and security-critical applications
- Rich structure of (hierarchical) classes
- Rare classes or long-tailed class distribution
- True ground truth unknown or uncertain



Wang et al. Learning to Model the Tail, 2017; Karimi et al., Deep learning with noisy labels: exploring techniques and remedies in medical image analysis, 2020; Bates et al., Distribution-Free, Risk-Controlling Prediction Sets, 2021; Northcutt et al., Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks, 2021.

Talk Outline

Conformal prediction:

- Notation and background
- Last talk: conformal training

Monte Carlo conformal prediction:

- Where does out ground truth come from?
- What if annotators disagree?
- How to calibrate against ambiguous ground truth?

Paper: arxiv.org/abs/2307.09302

Conformal Prediction

For model $\pi_{\theta,y} \approx p(y|x)$ construct confidence sets $C(x) \subseteq [K] = \{1, \dots, K\}$ such that

 $p(y \in C(x)) \geq 1 - lpha$ (coverage guarantee)

• confidence level α user-specified

Conformal Prediction

For model $\pi_{\theta,y} \approx p(y|x)$ construct confidence sets $C(x) \subseteq [K] = \{1, \dots, K\}$ such that

 $p(y \in C(x)) \geq 1 - lpha$ (coverage guarantee)

- confidence level α user-specified
- inefficiency = average confidence set size |C(x)|
- requires exchangeability, independent of model and distribution



Split Conformal Prediction

Split conformal prediction with two steps: prediction and calibration:

1. Prediction (test time): define how confidence sets are constructed

$$C(x):=\{k\in [K]: E(x,k):=\pi_{ heta,k}(x)\geq au\}$$

with $E(x,k) := \pi_{\theta,k}(x)$ called conformity scores.



Mauricio Sadinle, Jing Lei, and Larry Wasserman. Least ambiguous set-valued classifiers with bounded error levels. Journal of the American Statistical Association (JASA), 114(525):223–234, 2019.

Split Conformal Prediction

Split conformal prediction with two steps: prediction and calibration:

1. Prediction (test time): define how confidence sets are constructed

$$C(x):=\{k\in [K]: E(x,k):=\pi_{ heta,k}(x)\geq heta\}$$

with $E(x,k) := \pi_{\theta,k}(x)$ called conformity scores.

2. Calibration: define threshold au on N held-out calibration examples as



Conformal p-values

Alternative view (will be important later):

1. We test the null hypothesis that k is the true label of test example x:

 $H_k: y=k$

2. Compute a p-value for this hypothesis using:

$$ho_k = rac{\sum_{i=1}^N \, \delta[E(x_i,y_i) \leq E(x,k)] + 1}{N+1}$$

3. Construct confidence set

$$C(x) = \{k \in [K]:
ho_k \geq lpha\}$$

Public

Example Results

Inefficiency \downarrow for different methods (82% base accuracy):

Dataset, $lpha$	Thr	APS	RAPS
CIFAR10, 0.05	1.64	2.06	1.74
CIFAR10, 0.01	2.93	3.30	3.06

Different conformity scores

γ

Yaniv Romano, Matteo Sesia, and Emmanuel J. Candes. Classification with valid and adaptive coverage. In Advances in Neural Information Processing Systems (NIPS), 2020. Anastasios Nikolas Angelopoulos, Stephen Bates, Michael I. Jordan, Jitendra Malik: Uncertainty Sets for Image Classifiers using Conformal Prediction. ICLR 2021

Last Talk: Conformal Training

Conformal prediction is typically applied after training:

• Training loss and calibration objectives are not aligned!



cross-entropy loss

Last Talk: Conformal Training

Conformal training simulates split conformal prediction on each mini-batch:

- Allows to optimize arbitrary losses $\mathcal L$ on the confidence sets
- Independent of conformal predictor at test time + preserves coverage



→ More details: <u>arxiv.org/abs/2110.09192</u>

Talk Outline

Conformal prediction:

- Notation and background
- Last talk: conformal training

Monte Carlo conformal prediction:

- Where does our ground truth for calibration come from?
- What if this ground truth is uncertain because annotators disagree?
- How can we handle this during calibration?

Paper: arxiv.org/abs/2307.09302

Obtaining Calibration Labels

Need conformity scores of the true labels $E(x_i, y_i)$ for $x_i, y_i \sim p(x_i, y_i)$:



Obtaining Calibration Labels

Need conformity scores of the true labels $E(x_i, y_i)$ for $x_i, y_i \sim p(x_i, y_i)$:



Calibration Against Majority Voted Labels

Need conformity scores of the true labels $E(x_i, y_i)$ for $x_i, y_i \sim p(x_i, y_i)$:



- We have access to majority voted labels $y_{
 m vote} \sim p_{
 m vote}(y|x)$
- For this example, clearly $p_{ ext{vote}}
 eq p$
- But we need " $p_{\mathrm{vote}} = p$ " to guarantee coverage w.r.t. p

Proprietary + Confidential

A Simple Example







 In practice, we never observe these true labels (we cannot calibrate against them or obtain coverage against them)



- Ambiguity is captured in the true posteriors p(yert x)
- In practice, we usually do not observe the true posteriors either



- The "majority voted" label $y_{
 m vote} \sim p_{
 m vote}(y|x)$ ignores uncertainty
- We can calibrate and obtain coverage against $p_{
 m vote}
 eq p$

A Serious Example

Observation



b¹: {Pyogenic granuloma (Low)} {Hemangioma (Med)}
{Melanoma (High)}
b² {Angiokeratoma of skin (Low)} {Atypical Nevus (Med)}
b³: {Hemangioma (Med)} {Melanocytic Nevus (Low), Melanoma (High), O/E - ecchymoses present (Low)}
b⁴: {Hemangioma (Med), Melanoma (High), Skin Tag (Low)}
b⁵: {Melanoma (High)}
b⁶: {Hemangioma (Med)} {Melanoma (High)} {Melanocytic Nevus (Low)}

Conditions, Low/Med/High risk conditions

Annotations

A Serious Example

Observation



Annotations

b¹: {Pyogenic granuloma (Low)} {Hemangioma (Med)}
{Melanoma (High)}
b² {Angiokeratoma of skin (Low)} {Atypical Nevus (Med)}
b³: {Hemangioma (Med)} {Melanocytic Nevus (Low), Melanoma (High), O/E - ecchymoses present (Low)}
b⁴: {Hemangioma (Med), Melanoma (High), Skin Tag (Low)}
b⁵: {Melanoma (High)}

b⁶: {<u>Hemangioma</u> (Med)} {Melanoma (High)} {Melanocytic Nevus (Low)}

Conditions, Low/Med/High risk conditions

Majority vote

<u>Hemangioma</u> = benign

A Serious Example

Observation



Annotations

b¹: {Pyogenic granuloma (Low)} {Hemangioma (Med)}
{<u>Melanoma</u> (High)}
b² {Angiokeratoma of skin (Low)} {Atypical Nevus (Med)}
b³: {Hemangioma (Med)} {Melanocytic Nevus (Low), Melanoma (High), O/E - ecchymoses present (Low)}
b⁴: {Hemangioma (Med), <u>Melanoma</u> (High), Skin Tag (Low)}
b⁵: {<u>Melanoma</u> (High)}
b⁶: {Hemangioma (Med)} {<u>Melanoma</u> (High)} {Melanocytic Nevus (Low)}

Conditions, Low/Med/High risk conditions

Ignores a cancerous condition

Majority vote

Hemangioma = benign

Embracing Ambiguity in Conformal Prediction

Use annotations directly – for example, in terms of aggregated frequencies:



- Aggregating annotations is our best option to approximate the true p (we can only be as good in this tasks as our expert annotators are)
- How can we calibrate for and evaluate coverage w.r.t. $p_{
 m agg}$?

Aggregated Coverage for a Single Example

Call estimates of $\lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:





 $\lambda = (0, 0, 0.32, 0.46, 0.02, 016, 0.04, 0, 0, 0)$

C(x) = {cat, dog} – do we have coverage?

Majority-voted coverage	1	
Aggregated coverage	0.62 = 0.46 + 0.16	

"Covered plausibility mass"

Call estimates of $\lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:

```
p_{\mathrm{agg}}(y \in C(x))
Guarantee coverage "against annotations"
```

Call estimates of $\lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:

$$p_{\mathrm{agg}}(y\in C(x)) = \mathbb{E}_{p_{\mathrm{agg}}}[\delta[y\in C(x)]]$$

Binary event, express as expectation

Call estimates of $\lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:

$$egin{aligned} p_{ ext{agg}}(y \in C(x)) &= \mathbb{E}_{p_{ ext{agg}}}[\delta[y \in C(x)]] \ &= \mathbb{E}_{x,y \sim p(x)p_{ ext{agg}}(y|x)}[\delta[y \in C(x)]] \ & \uparrow \end{aligned}$$

Decompose joint probability

Call estimates of $\lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:

$$p_{\mathrm{agg}}(y \in C(x)) = \mathbb{E}_{p_{\mathrm{agg}}}[\delta[y \in C(x)]]$$

 $= \mathbb{E}_{x,y \sim p(x)p_{\mathrm{agg}}(y|x)}[\delta[y \in C(x)]]$
 $= \mathbb{E}_{x \sim p(x)}[\mathbb{E}_{y \sim p_{\mathrm{agg}}(y|x)}[\delta[y \in C(x)]]]$
Distribute coverage across plausibilities $\longrightarrow \sum_{k} \lambda_k \delta[k \in C(x)]$

Call estimates of $\ \lambda_{ik} = p_{
m agg}(y=k|x_i) pprox p(y|x_i)$ plausibilities:

$$p_{
m agg}(y \in C(x)) = \mathbb{E}_{p_{
m agg}}[\delta[y \in C(x)]]$$

 $= \mathbb{E}_{x,y \sim p(x)p_{
m agg}(y|x)}[\delta[y \in C(x)]]$
 $= \mathbb{E}_{x \sim p(x)}[\mathbb{E}_{y \sim p_{
m agg}(y|x)}[\delta[y \in C(x)]]]$
Distribute coverage across plausibilities $\longrightarrow \sum_{k} \lambda_k \delta[k \in C(x)]$

- → Coverage is marginal over examples and labels!
- → If $p_{agg} = p$, this is coverage wrt. the true labels!

Coverage distributed across examples and labels



Calibrating with Voted Labels



Calibrate Against Sampled Labels

Basic idea:

- Use plausibilities for calibration: $\lambda_{ik} = p_{
 m agg}(y=k|x_i) pprox p(y|x_i)$
- Repeat each calibration example M times
- Standard calibration using the *augmented* calibration set

$$\{E(x_i,y_{ij})\}_{i\in[N],j\in[M]}$$
 with $y_{ij}\sim p_{\mathrm{agg}}(y_{ij}=k|x_i)=\lambda_{ik}$

Problem:

• Invalidates coverage by breaking exchangeability:

 $p(z_{11},z_{12},z_{13},\ldots,z_{21},\ldots,z_{NM},z)$ for $z_{ij}=(x_i,y_{ij})$ and test example z I know the first M examples are repeated

Monte Carlo Conformal Prediction

Solution:

- Use plausibilities for calibration: $\lambda_{ik} = p_{
 m agg}(y=k|x_i) pprox p(y|x_i)$
- Repeat each calibration example M times
- Calibrate using the *augmented* calibration set

 $\{E(x_i,y_{ij})\}_{i\in[N],j\in[M]}$ with $y_{ij}\sim p_{\mathrm{agg}}(y_{ij}=k|x_i)=\lambda_{ik}$

• Adjust quantile computation to

$$rac{\lfloor lpha(N+1)
floor}{N} \longrightarrow rac{\lfloor lpha M(N+1)
floor - M+1}{MN}$$

$$ho_k = rac{\sum_{i=1}^N \sum_{j=1}^M \delta[E(x_i,y_{ij}) \leq E(x,k)] + 1}{M \cdot N + 1}$$

$$ho_k = rac{\sum_{i=1}^N \sum_{j=1}^M \delta[E(x_i, y_{ij}) \leq E(x, k)] + 1}{M \cdot N + 1}$$



$$ho_k = rac{\sum_{i=1}^N \sum_{j=1}^M \delta[E(x_i,y_{ij}) \leq E(x,k)] + 1}{M \cdot N + 1}$$
 $ar{
ho}_k = rac{1}{M} \sum_{j=1}^M
ho_k^j \longrightarrow
ho_k^j \stackrel{!}{=} rac{\sum_{i=1}^N \delta[E(x_i,y_{ij}) \leq E(x,k)] + 1}{N+1}$

Consider the p-values computed for standard conformal prediction:



→ <u>Vovk and Wang</u> establish coverage $1 - 2\alpha$ when averaging p-values

Coverage Beyond 1-2lpha

Monte Carlo conformal prediction:

- Can be re-formulated as averaging M p-values
- This establishes a $1-2lpha\,$ coverage guarantee
- Can improve to $(1-lpha)(1-\delta)$ for $\,\delta>0\,$ with additional calibration split

Remarks:

- Empirically, we always observe coverage 1-lpha
- Without ambiguity, we recover standard conformal prediction (any M)
- Ambiguous examples: we improve coverage by sacrificing efficiency
- Unambiguous examples: it behaves like standard conformal prediction

Results in Dermatology



Results in Dermatology



Qualitative Results in Dermatology





Monte Carlo CP w.r.t. $p_{ m agg}$









Conclusion for Monte Carlo CP

= conformal prediction based on sampled labels from annotators/plausibilities.

- The labels we have access to are usually voted labels, from $p_{
 m vote}$
- In ambiguous settings, voted labels can deviate from true labels:

$$p_{ ext{vote}}
eq p$$

- Monte Carlo conformal prediction samples labels from $p_{
 m agg} pprox p$
- The best we can do: "calibrate wrt. to annotators"
- Establishes coverage guarantees for multi-label classification and calibration with data augmentation

Paper: arxiv.org/abs/2307.09302 | Contact: davidstutz.de / dstutz@google.com