Learning Sound Event Classifiers from Web Audio with Noisy Labels

Eduardo Fonseca¹, Manoj Plakal², Daniel P. W. Ellis², Frederic Font¹, Xavier Favory¹, and Xavier Serra¹



UniversitatMTGPompeu FabraMusic TechnologyBarcelonaGroup





Label noise in sound event classification

- Labels that fail to properly represent acoustic content in audio clip
- Why is label noise relevant?



Label noise effects: performance decrease / increased complexity

How to mitigate label noise?

How to mitigate label noise?



How to mitigate label noise?



Our contributions

1. FSDnoisy18k: a dataset to foster label noise research



FSDnoisy18k

🔞 Eduardo Fonseca; Mercedes Collado; Manoj Plakal; 🙆 Daniel P. W. Ellis; 🔞 Frederic Font; Xavier Favory; 🔞 Xavier Serra

FSDnoisy18k is an audio dataset collected with the aim of fostering the investigation of label noise in sound event classification. It contains 42.5 hours of audio across 20 sound classes, including a small amount of manually-labeled data and a larger quantity of real-world noisy data.

Our contributions

- 1. FSDnoisy18k: a dataset to foster label noise research
- 2. CNN baseline system
- 3. Evaluation of noise-robust loss functions





FSDnoisy18k: creation

- Freesound
 - → Audio content & metadata (tags)
- AudioSet Ontology
 - → 20 classes (labels)

- What free sound

{ III } AudioSet

FSDnoisy18k: creation

- Freesound
 - Audio content & metadata (tags)
- AudioSet Ontology
 - → 20 classes (labels)

- What free sound

{ III } AudioSet



FSDnoisy18k: creation

- Freesound
 - Audio content & metadata (tags)
- AudioSet Ontology
 - → 20 classes (labels)

- What free sound

{ III } AudioSet



Types of label noise

• singly-labeled data



Types of label noise

• singly-labeled data



Types of label noise

• singly-labeled data



- in-vocabulary (IV): events that are part of our target class set (closed-set)
- out-of-vocabulary (OOV): events not covered by the class set (open-set)



Observed label from the vocabulary:

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart /

Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / Walk, footsteps / Wind / Writing

Examples: clip #1 🔹

<u>True</u> label from the vocabulary:

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart / Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / **Walk, footsteps** / Wind / Writing





Observed label from the vocabulary:

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart /

Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / Walk, footsteps / Wind / Writing



True label from the vocabulary:

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart /

Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / Walk, footsteps / Wind / Writing

Missing labels: male speech / laughter / children shouting / chirp, tweet / chatter





<u>Observed</u> label from the vocabulary

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart /

Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / Walk, footsteps / Wind / Writing



<u>True</u> label from the vocabulary:

Acoustic guitar / Bass guitar / Clapping / Coin (dropping) / Crash cymbal / Dishes, pots, and pans / Engine / Fart / Fire / Fireworks / Glass / Hi-hat / Piano / Rain / Slam / Squeak / Tearing / Walk, footsteps / Wind / Writing

True label: electronic music



Label noise distribution in FSDnoisy18k



- most frequent types of label noise: OOV
- *some clips are incorrectly labeled, but still similar in terms of acoustics

FSDnoisy18k

- 20 classes / 18k clips / 42.5 h
- singly-labeled data
- variable clip duration: 300ms 30s
- proportion train_noisy / train_clean = 90% / 10%
- per-class varying degree of types and amount of label noise
- expandable
- <u>http://www.eduardofonseca.net/FSDnoisy18k/</u>



CNN baseline system



- Why?
 - → model-agnostic / minimal intervention / efficient

• Why?

- model-agnostic / minimal intervention / efficient
- Default loss function in multi-class setting: Categorical Cross-Entropy (CCE)

$$\mathcal{L}_{cce} = -\sum_{k=1}^{K} \underline{y(k)} \log(p(k))$$
predictions
target labels

- Why?
 - model-agnostic / minimal intervention / efficient
- Default loss function in multi-class setting: Categorical Cross-Entropy (CCE)

$$\mathcal{L}_{cce} = -\sum_{k=1}^{K} y(k) \log(p(k))$$

- CCE is sensitive to label noise: emphasis on *difficult* examples (weighting)
 - → beneficial for clean data
 - → detrimental for noisy data

- Soft bootstrapping
 - dynamically update target labels based on model's current state
 - → updated target label: convex combination

$$\mathcal{L}_{soft} = -\sum_{k=1}^{K} [\beta y(k) + (1 - \beta) p(k)] \log(p(k)), \quad \beta \in [0, 1]$$
predictions
target labels

Scott E. Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, Andrew Rabinovich, *Training Deep Neural Networks on Noisy Labels with Bootstrapping*. In ICLR 2015

- \mathcal{L} q loss intuition
 - → CCE: sensitive to noisy labels (weighting)
 - → Mean Absolute Error (MAE):

$$\mathcal{L}_{mae} = \sum_{k=1}^{K} |y(k) - p(k)|$$

- avoid weighting
- difficult convergence

Zhilu Zhang and Mert Sabuncu, *Generalized cross entropy loss for training deep neural networks with noisy labels*. In NeurIPS 2018

 \mathcal{L}_{q} loss intuition

 \rightarrow

- CCE: sensitive to noisy labels (weighting)
- Mean Absolute Error (MAE): $\mathcal{L}_{mae} = \sum_{k=1}^{n} |y(k) - p(k)|$ \rightarrow
 - avoid weighting
 - difficult convergence
- \mathcal{L}_q loss is a generalization of CCE and MAE:
 - negative Box-Cox transformation of softmax predictions \rightarrow

$$\mathcal{L}_{q} = \frac{1 - \left(\sum_{k=1}^{K} y(k)p(k)\right)^{q}}{q}, \quad q \in (0, 1]$$

$$q = 1 \rightarrow \mathcal{L}_{q} = \mathsf{MAE} \quad ; \quad q \rightarrow 0 \rightarrow \mathcal{L}_{q} = \mathsf{CCE}$$

Zhilu Zhang and Mert Sabuncu, Generalized cross entropy loss for training deep neural networks with noisy labels. In NeurIPS 2018

Experiments



- supervision by user-provided tags can be useful for sound event classification
- \mathcal{L}_q works well for sound classification tasks with OOV (and some IV) noises

Experiments



- boost by using \mathcal{L} q on noisy set: 1.9% (little engineering effort)
- boost by adding curated data to noisy set: 5.1% (significant manual effort)

Summary & takeaways

- FSDnoisy18k
 - open dataset for investigation of label noise
 - → 20 classes / 18k clips / 42.5 h / singly-labeled data
 - → small amount of manually-labelled data and a large amount of noisy data
 - → label noise characterization
- CNN baseline system
 - → large amount of Freesound audio & tags feasible for training sound recognizers
- Noise-robust loss functions
 - → efficient way to improve performance in presence of noisy labels
 - \mathcal{L}_q is top-performing loss

If you are interested in label noise...





Learning Sound Event Classifiers from Web Audio with Noisy Labels

Thank you!

http://www.eduardofonseca.net/FSDnoisy18k/

https://zenodo.org/record/2529934

https://github.com/edufonseca/icassp19

Eduardo Fonseca¹, Manoj Plakal², Daniel P. W. Ellis², Frederic Font¹, Xavier Favory¹, and Xavier Serra¹



UniversitatMTGPompeu FabraMusic TechnologyBarcelonaGroup





Why this vocabulary?

- data availability
- classes "suitable" for the study of label noise
 - classes described with tags also used for other audio materials
 - Bass guitar, Crash cymbal, Engine, ...
 - → field-recordings: several sound sources expected
 - only the most predominant(s) tagged: Rain, Fireworks, Slam, Fire, ...
 - \rightarrow pairs of related classes:
 - Squeak & Slam / Wind & Rain