

CSC2412: Exponential Mechanism & Private PAC Learning

Sasho Nikolov

Classification Basics

The learning problem

Problem: develop an algorithm that classifies avocados into ⁺¹ripe and ⁻¹unripe.

We have a big data set of avocado data. For each avocado, we have:

- colour, firmness, size, shape, skin texture, ... *features*
- ripe or not *label*

From this data, we want to classify unseen avocados.

*Learn a rule to predict
label from features*



The learning problem, formally

Model: all possible settings to the features

- Known data universe \mathcal{X} and an unknown probability distribution D on \mathcal{X}

- Known concept class C and an unknown concept $c \in C$

All allowed rules to map features to label. | E.g. all functions that depend on ≤ 3 features

Assumption (realizability):
some $c \in C$ can produce all correct labels

- We get a dataset $X = \{(x_1, c(x_1)), \dots, (x_n, c(x_n))\}$, where each x_i is an independent sample from D .

features → label

Goal: Learn c from X .

an approximation of c

E.g. an avocado is ripe
iff colour =  and firmness = medium } c

The goal, formally

The error of a concept $c' \in C$ is

$$L_{D,c}(c') = \mathbb{P}_{x \sim D}(c'(x) \neq c(x)).$$

↳ loss of c' (w.r.t. D, c)

on input $X = ((x_1, c(x_1)), \dots, (x_n, c(x_n)))$ sampled iid from D

We want an algorithm \mathcal{M} that outputs some $c' \in C$ and satisfies

↳ output of $\mathcal{M}(X)$ misclassifies $\leq \alpha$ fraction of the population

$$\mathbb{P}(L_{D,c}(\mathcal{M}(X)) \leq \alpha) \geq 1 - \beta.$$

↳ taken over randomness in choosing X and any randomness of \mathcal{M}

↳ Fraction of the population labeled incorrectly by c'

Probably Approximately Correct learning (PAC) [Valiant]

Empirical risk minimization

Issue: We want to find $\arg \min_{c' \in C} L_{D,c}(c')$, but we do not know D, c .

↳ approximate minimizers are also ok

Solution: Instead we solve $\arg \min_{c' \in C} L_X(c')$, where

$$L_X(c) = 0$$

is the empirical error.

$$L_X(c') = \frac{|\{i : c'(x_i) \neq c(x_i)\}|}{n}$$

fraction of pts in X
misclassified by c'

$L_{D,c}$ population loss
(unknown)

L_X empirical loss
(known)

Theorem (Uniform convergence) → Pop and emp. loss are close for $\forall c' \in C$
Suppose that $n \geq \frac{\ln(|C|/\beta)}{2\alpha^2}$. Then, with probability $\geq 1 - \beta$,

Hoeffding's inequality
(exercise)

$$\max_{c' \in C} L_{D,c}(c') - L_X(c') \leq \alpha.$$

$$\forall c' \in C \quad L_{D,c}(c') \leq L_X(c') + \alpha$$

other versions for
infinite C , e.g. VC-dimension

Private learning

In private PAC learning, we require that

- when X is a sample of iid labeled data points, we learn ^{an approximately} ~~the~~ correct concept, as in standard PAC learning;

\mathcal{M} that on input X outputs $c' \in C$

- the learning algorithm is ϵ -differentially private for any labeled data set $X \in (\mathcal{X} \times C)^n$.

$\forall X, X'$ neighbouring $\forall S \subseteq C$

$$\mathbb{P}(\mathcal{M}(X) \in S) \leq e^\epsilon \cdot \mathbb{P}(\mathcal{M}(X') \in S)$$

Privacy must hold even if data X not iid

Want to do ERM w/ ϵ -DP

i.e. (approximately) minimize

$$L_X(c') = \frac{|\{i : c(x_i) \neq c'(x_i)\}|}{n} \quad \text{over } c' \in \mathcal{C}$$

How can we use Laplace noise mechanism for this?

$L_X(c')$ is a counting query Exercise: analyze this ↓

We could release answers to all counting queries

$$\mathcal{C} = \{c_1, \dots, c_k\} \quad \{L_X(c_1), L_X(c_2), \dots, L_X(c_k)\}$$

Exponential mechanism

We want to solve $\underline{\arg \min}_{c' \in C} L_X(c')$.

How do we minimize with differential privacy?

Sample concepts with less error with higher probability

$$\mathbb{P}(\mathcal{M}(X) = c') \propto \exp\left(-\frac{\varepsilon n}{2} L_X(c')\right)$$

↓
proportional to

Exponential Mechanism

General set-up: score function $u: \mathcal{X}^n \times \mathcal{Y} \rightarrow \mathbb{R}$

Goal: given X , find $\arg \max_{y \in \mathcal{Y}} u(X, y)$

Sensitivity

$$\Delta u = \max_{y \in \mathcal{Y}} \max_{X \sim X'} |u(X, y) - u(X', y)|.$$

How different can the score be between neighbouring X, X'

The mechanism $\mathcal{M}_{\text{exp}}(X)$ which outputs a random Y so that

$$\mathbb{P}(Y = y) = \frac{e^{\epsilon u(X, y) / 2\Delta u}}{\sum_{z \in \mathcal{Y}} e^{\epsilon u(X, z) / 2\Delta u}} \rightarrow \text{normalizing factor}$$

is ϵ -differentially private

$u(X, y)$ = "how good of an output is y for the dataset X ?"

$$\mathbb{P}(Y = y) = \frac{e^{\varepsilon u(X, y) / 2\Delta u}}{\sum_{z \in \mathcal{Y}} e^{\varepsilon u(X, z) / 2\Delta u}}$$

Enough to show : $\forall x \sim x'$
 $\forall y \in \mathcal{Y}$ $\frac{\mathbb{P}(M(x) = y)}{\mathbb{P}(M(x') = y)} \leq e^\varepsilon$

$$\frac{\mathbb{P}(M(x) = y)}{\mathbb{P}(M(x') = y)} = \exp\left(\frac{\varepsilon(u(x, y) - u(x', y))}{2\Delta u}\right) \cdot \frac{\sum_{z \in \mathcal{Y}} e^{\varepsilon u(x', z) / 2\Delta u}}{\sum_{z \in \mathcal{Y}} e^{\varepsilon u(x, z) / 2\Delta u}}$$

$$\leq e^{\varepsilon/2} \cdot e^{\varepsilon/2} = e^\varepsilon$$



Accuracy of the exponential mechanism

≥ 1 output y^*
achieves $\text{OPT}(X)$

$$\text{OPT}(X) = \max_{y \in \mathcal{Y}} u(X, y)$$

Then, for the output $Y = \mathcal{M}_{\text{exp}}(X)$,

$$\mathbb{P}(u(X, Y) \leq \text{OPT}(X) - t) \leq$$

$$\frac{e^{\varepsilon(\text{OPT}(X) - t)/2\Delta u}}{\sum_{z \in \mathcal{Y}} e^{\varepsilon u(X, z)/2\Delta u}} \cdot |\{y : u(X, y) \leq \text{OPT}(X) - t\}|$$

$\geq e^{\varepsilon u(X, y^*)/2\Delta u}$

$$\leq \frac{e^{\varepsilon(\text{OPT} - t)/2\Delta u} \cdot (|\mathcal{Y}| - 1)}{e^{\varepsilon \text{OPT}/2\Delta u}} \leq e^{-\varepsilon t/2\Delta u} \cdot |\mathcal{Y}|$$

$$\leq e^{-\varepsilon t/2\Delta u} \cdot |\mathcal{Y}|$$

Private Learning

Unknown distribution D on known \mathcal{X}

Unknown c in a known concept class \mathcal{C}

Data set $X = \{(x_1, c(x_1)), \dots, (x_n, c(x_n))\}$ where $x_1, \dots, x_n \sim_{\text{iid}} D$

$$L_{D,c}(c') = \mathbb{P}_{x \sim D} (c(x) \neq c'(x)) \quad \left| \quad L_X(c') = \frac{|\{i : c(x_i) \neq c'(x_i)\}|}{n}$$

If $n \geq \frac{\ln(|\mathcal{C}|/\beta)}{2d^2}$ then w/ prob $\geq 1-\beta$,
 $\forall c' \in \mathcal{C} : L_{D,c}(c') \leq L_X(c') + d$

Exponential mechanism: sample $y \in \mathcal{Y}$ w/ prob.
proportional to $\exp(\varepsilon u(x, y) / 2\Delta u)$

$$\Delta u = \max_{y \in \mathcal{Y}} \max_{x \sim x'} |u(x, y) - u(x', y)|$$

Putting things together

A concept class C can be learned by an ϵ -differentially private mechanism when the sample size is

$$n \geq \max \left\{ \frac{4 \ln(2|C|/\beta)}{\epsilon \alpha}, \frac{2 \ln(2|C|/\beta)}{\alpha^2} \right\}$$

Use exp. mechanism with $\mathcal{Y} = C$

with $u(X, c') = -L_X(c')$

With prob $\geq 1 - \beta/2$, c' output by $M_{\text{exp}}(X)$ has $L_X(c') \leq \frac{\alpha}{2}$. By unif. convergence, w/ prob $\geq 1 - \beta/2$, $L_{D, c}(c') - L_X(c') \leq \frac{\alpha}{2}$

Putting things together

$$u(X, c') = -L_X(c'); \quad \Delta u = \frac{1}{n}; \quad \text{sample } c' \in C \text{ w/ prob}$$

prop. to $\exp(-\epsilon n L_X(c')/2)$

$$\text{OPT}(X) = \max_{c' \in C} -L_X(c')$$

↓
 ϵ -DP

$$= -\min_{c' \in C} L_X(c') = 0$$

$$\mathbb{P}(L_X(\mathcal{M}(X)) \geq \frac{\alpha}{2}) = \mathbb{P}(u(X, \mathcal{M}(X)) \leq \text{OPT}(X) - \frac{\alpha}{2})$$

$$\leq e^{-\epsilon \alpha^n / 4} \cdot |C| \leq \beta/2 \quad \text{if } n \geq \frac{4 \ln(2|C|/\beta)}{\epsilon \alpha}$$

+ by unif. conv. (with $\alpha/2, \beta/2$) w/ prob $1 - \beta/2$, $L_{D,c}(\mathcal{M}(X)) \leq L_X(\mathcal{M}(X)) + \frac{\alpha}{2}$

w/ prob $\geq 1 - \beta$
 $\frac{\alpha}{2} + \frac{\alpha}{2} = \alpha$.