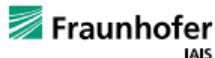


# Supervised methods for quantification

**Mirko Bunse**, Alejandro Moreo, and Fabrizio Sebastiani

LQ @ ECML-PKDD 2024 – September 13<sup>th</sup>

Partner institutions:



Institutionally funded by:



Ministerium für  
Kultur und Wissenschaft  
des Landes Nordrhein-Westfalen



## Problem statement



**Given:** a labeled training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n \sim \mathbb{P}^n$  where

- $\mathcal{X}$  is the feature space (e.g.,  $\mathcal{X} = \mathbb{R}^d$ )
- $\mathcal{Y} = \{1, 2, \dots, C\}$  is the set of class labels

# Problem statement

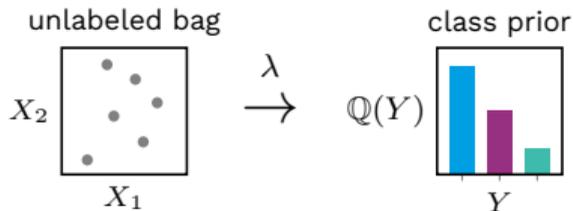


**Given:** a labeled training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n \sim \mathbb{P}^n$  where

- $\mathcal{X}$  is the feature space (e.g.,  $\mathcal{X} = \mathbb{R}^d$ )
- $\mathcal{Y} = \{1, 2, \dots, C\}$  is the set of class labels

**Find:** a quantifier  $\lambda : \bigcup_{m=1}^{\infty} \mathcal{X}^m \rightarrow \Delta^{C-1}$  where

- $\bigcup_{m=1}^{\infty} \mathcal{X}^m$  is the space of unlabeled data bags of any size  $m$
- $\Delta^{C-1} = \{\mathbf{p} \in \mathbb{R}^C : \mathbf{p}_i \geq 0 \forall i, \sum_{i=1}^C \mathbf{p}_i = 1\}$  is the space of class prevalences
- for any bag  $B \sim \mathbb{Q}^m$ , we want to achieve that  $\lambda(B) = \mathbb{Q}(Y)$



# Prior probability shift



We typically want to achieve  $\lambda(B) = \mathbb{Q}(Y)$  when otherwise unknown

## Definitions:

- $\forall \mathbf{x} \in B : \mathbf{x} \sim \mathbb{Q}(\mathbf{x})$  where  $\mathbb{Q}(\mathbf{x}) = \sum_{y=1}^C \mathbb{Q}(\mathbf{x}, y)$  (law of total probability)
- $\forall (\mathbf{x}, y) \in D : (\mathbf{x}, y) \sim \mathbb{P}(\mathbf{x}, y)$

<sup>1</sup> Kull and Flach, “Patterns of dataset shift”, 2014.



## Prior probability shift

We typically want to achieve  $\lambda(B) = \mathbb{Q}(Y)$  when otherwise unknown

### Definitions:

- $\forall \mathbf{x} \in B : \mathbf{x} \sim \mathbb{Q}(\mathbf{x})$  where  $\mathbb{Q}(\mathbf{x}) = \sum_{y=1}^C \mathbb{Q}(\mathbf{x}, y)$  (law of total probability)
- $\forall (\mathbf{x}, y) \in D : (\mathbf{x}, y) \sim \mathbb{P}(\mathbf{x}, y)$

### Identically & independently distributed (IID) data:

- $\mathbb{Q}(X, Y) = \mathbb{P}(X, Y)$
- we could estimate  $\mathbb{Q}(Y) = \mathbb{P}(Y)$

<sup>1</sup> Kull and Flach, “Patterns of dataset shift”, 2014.



## Prior probability shift

We typically want to achieve  $\lambda(B) = \mathbb{Q}(Y)$  when otherwise unknown

### Definitions:

- $\forall \mathbf{x} \in B : \mathbf{x} \sim \mathbb{Q}(\mathbf{x})$  where  $\mathbb{Q}(\mathbf{x}) = \sum_{y=1}^C \mathbb{Q}(\mathbf{x}, y)$  (law of total probability)
- $\forall (\mathbf{x}, y) \in D : (\mathbf{x}, y) \sim \mathbb{P}(\mathbf{x}, y)$

### Identically & independently distributed (IID) data:

- $\mathbb{Q}(X, Y) = \mathbb{P}(X, Y)$
- we could estimate  $\mathbb{Q}(Y) = \mathbb{P}(Y)$

### Prior probability shift (PPS):

- $\mathbb{Q}(X | Y) = \mathbb{P}(X | Y)$
- $\mathbb{Q}(Y) \neq \mathbb{P}(Y)$

typical assumption in quantification

More types of data set shift exist.<sup>1</sup>

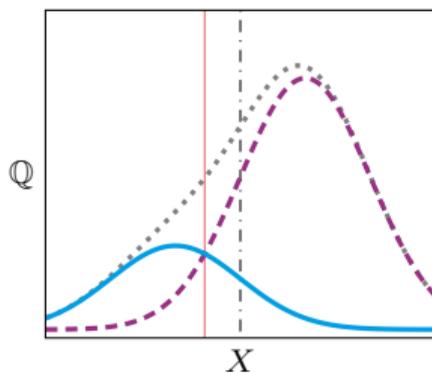
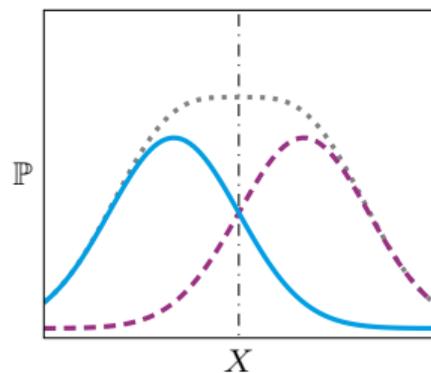
<sup>1</sup> Kull and Flach, "Patterns of dataset shift", 2014.

# Prior probability shift



training data:  $\mathbb{P}(Y) = (\frac{1}{2}, \frac{1}{2})$

testing bag:  $\mathbb{Q}(Y) = (\frac{1}{4}, \frac{3}{4})$



- $\mathbb{Q}(X | Y = 1) = \mathbb{P}(X | Y = 1)$
- - -  $\mathbb{Q}(X | Y = 2) = \mathbb{P}(X | Y = 2)$
- ⋯  $\mathbb{Q}(X, Y)$  or  $\mathbb{P}(X, Y)$
- - - Bayes-optimal classifier for  $\mathbb{P}$
- Bayes-optimal classifier for  $\mathbb{Q}$

We cannot learn a classifier (solely) from  $\mathbb{P}$  that is also optimal for  $\mathbb{Q}$ .

# Classification versus quantification



For both tasks, we are given  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$

## Classification:

- find  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- reason about individual data items
- (typically) assume IID data

## Quantification:

- find  $\lambda : \bigcup_{m=1}^{\infty} \mathcal{X}^m \rightarrow \Delta^{C-1}$
- reason about bags of data
- (typically) assume PPS

PPS requires quantifiers that are more sophisticated than Classify & Count.<sup>2</sup>

<sup>2</sup> Forman, “Quantifying counts and costs via classification”, 2008.



1. Problem statement

## 2. Desirable properties of quantifiers

3. Binary quantifiers

4. Multi-class quantifiers

5. Numerical optimization

6. Loss functions & data representations

7. Beyond linear systems of equations

# Fisher consistency



## Definition (Fisher consistency for PPS):

If a quantifier had access to the entire population, it would return the correct class prevalences:

$$\underbrace{\lambda'(\mathbb{Q}(X))}_{\text{population analogue of } \lambda(B)} = \mathbb{Q}(Y) \quad \forall \mathbb{Q} : \underbrace{\mathbb{Q}(X | Y) = \mathbb{P}(X | Y)}_{\text{for any } \mathbb{Q} \text{ with PPS}}$$

# Fisher consistency



## Definition (Fisher consistency for PPS):

If a quantifier had access to the entire population, it would return the correct class prevalences:

$$\underbrace{\lambda'(Q(X))}_{\text{population analogue of } \lambda(B)} = Q(Y) \quad \forall Q : \underbrace{Q(X | Y) = \mathbb{P}(X | Y)}_{\text{for any } Q \text{ with PPS}}$$

## Notes:

- can also be defined for other types of data set shift
- is different from unbiasedness and different from asymptotical consistency
- does not indicate good performance on finite samples
- hence, not a sufficient but certainly **a necessary criterion** for quantifier selection

**Tip:** write down this definition; there might be a small assignment!

## Estimation error



Since data is limited, we cannot solely rely on Fisher consistency.

**Empirical evaluation:** test quantifiers on data

- employ suitable protocols (as discussed the previous part of this tutorial)
- employ a representative collection of data sets

## Estimation error



Since data is limited, we cannot solely rely on Fisher consistency.

**Empirical evaluation:** test quantifiers on data

- employ suitable protocols (as discussed the previous part of this tutorial)
- employ a representative collection of data sets

**Asymptotical consistency:** look for desirable asymptotical behaviour; with any bound of the type

$$\|\lambda(B) - \mathbf{p}^*\| \leq f(\lambda, |D|, |B|, \delta)$$

prefer those quantifiers  $\lambda$  that achieve a small upper bound with a high probability  $1 - \delta$

# Resource efficiency



## User perspective:

- little waiting times for predictions
- without requiring excessive hardware

<sup>3</sup> Google, *Environmental Report*, 2024.

# Resource efficiency



## User perspective:

- little waiting times for predictions
- without requiring excessive hardware

## Environmental perspective:

- greenhouse gas emissions: use little computation and green energy
- Google:<sup>3</sup> “reducing emissions may be challenging due to increasing energy demands from the greater intensity of AI compute.”  
(their emissions increased by 48%, as compared to 2019, despite their goal of reducing emissions by 50% in 2030)

<sup>3</sup> Google, *Environmental Report*, 2024.

# Resource efficiency



## User perspective:

- little waiting times for predictions
- without requiring excessive hardware

## Environmental perspective:

- greenhouse gas emissions: use little computation and green energy
- Google:<sup>3</sup> “reducing emissions may be challenging due to increasing energy demands from the greater intensity of AI compute.”  
(their emissions increased by 48%, as compared to 2019, despite their goal of reducing emissions by 50% in 2030)

## Implications on quantification research:

- reduce resource consumption
- report on resource consumption (prediction times, memory consumption, GHG emissions, ...)

<sup>3</sup> Google, *Environmental Report*, 2024.



1. Problem statement
2. Desirable properties of quantifiers

## 3. Binary quantifiers

4. Multi-class quantifiers
5. Numerical optimization
6. Loss functions & data representations
7. Beyond linear systems of equations

# Binary Adjusted Classify & Count



## Preliminaries:

- let  $\mathcal{Y} = \{1, 2\}$  (binary quantification)
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

# Binary Adjusted Classify & Count



## Preliminaries:

- let  $\mathcal{Y} = \{1, 2\}$  (binary quantification)
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\Rightarrow \mathbb{Q}(h(X) = 1) = \sum_{i \in \mathcal{Y}} \mathbb{Q}(h(X) = 1 \mid Y = i) \cdot \mathbb{Q}(Y = i) \quad (\text{law of total probability})$$

# Binary Adjusted Classify & Count



## Preliminaries:

- let  $\mathcal{Y} = \{1, 2\}$  (binary quantification)
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\begin{aligned}\Rightarrow \mathbb{Q}(h(X) = 1) &= \sum_{i \in \mathcal{Y}} \mathbb{Q}(h(X) = 1 \mid Y = i) \cdot \mathbb{Q}(Y = i) \quad (\text{law of total probability}) \\ &= \underbrace{\text{TPR}}_{\mathbb{Q}(h(X) = 1 \mid Y = 1)} \cdot \mathbb{Q}(Y = 1) + \underbrace{\text{FPR}}_{\mathbb{Q}(h(X) = 1 \mid Y = 2)} \cdot (1 - \mathbb{Q}(Y = 1))\end{aligned}$$

# Binary Adjusted Classify & Count



## Preliminaries:

- let  $\mathcal{Y} = \{1, 2\}$  (binary quantification)
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\begin{aligned}\Rightarrow \mathbb{Q}(h(X) = 1) &= \sum_{i \in \mathcal{Y}} \mathbb{Q}(h(X) = 1 \mid Y = i) \cdot \mathbb{Q}(Y = i) \quad (\text{law of total probability}) \\ &= \underbrace{\text{TPR}}_{\mathbb{Q}(h(X) = 1 \mid Y = 1)} \cdot \mathbb{Q}(Y = 1) + \underbrace{\text{FPR}}_{\mathbb{Q}(h(X) = 1 \mid Y = 2)} \cdot (1 - \mathbb{Q}(Y = 1)) \\ \Rightarrow \mathbb{Q}(Y = 1) &= \frac{\mathbb{Q}(h(X) = 1) - \text{FPR}}{\text{TPR} - \text{FPR}}\end{aligned}$$

where:

- $\mathbb{Q}(h(X) = 1)$  can be estimated by counting the predictions  $h(\mathbf{x}) \forall \mathbf{x} \in B$
- TPR and FPR can be estimated with the training data  $D$  (due to PPS)

# Binary Adjusted Classify & Count



**Definition (binary ACC):**

$$\hat{Q}(Y = 1) = \frac{\hat{Q}(h(X) = 1) - \hat{FPR}}{\hat{TPR} - \hat{FPR}}$$

is Fisher-consistent,<sup>4</sup> where

- $\hat{Q}(h(X) = 1) = \frac{1}{|B|} \sum_{\mathbf{x} \in B} \mathbb{1}_{h(\mathbf{x})=1}$
- $\hat{TPR} = \frac{1}{|D_1|} \sum_{\mathbf{x} \in D_1} \mathbb{1}_{h(\mathbf{x})=1}$
- $\hat{FPR} = \frac{1}{|D_2|} \sum_{\mathbf{x} \in D_2} \mathbb{1}_{h(\mathbf{x})=1}$
- $D_i = \{(\mathbf{x}, y) \in D : y = i\} \forall i \in \mathcal{Y}$

<sup>4</sup> Tasche, “Fisher consistency for prior probability shift”, 2017.

# Binary Adjusted Classify & Count



**Definition (binary ACC):**

$$\hat{Q}(Y = 1) = \frac{\hat{Q}(h(X) = 1) - \text{F}\hat{\text{P}}\text{R}}{\text{T}\hat{\text{P}}\text{R} - \text{F}\hat{\text{P}}\text{R}}$$

is Fisher-consistent,<sup>4</sup> where

- $\hat{Q}(h(X) = 1) = \frac{1}{|B|} \sum_{\mathbf{x} \in B} \mathbb{1}_{h(\mathbf{x})=1}$
- $\text{T}\hat{\text{P}}\text{R} = \frac{1}{|D_1|} \sum_{\mathbf{x} \in D_1} \mathbb{1}_{h(\mathbf{x})=1}$
- $\text{F}\hat{\text{P}}\text{R} = \frac{1}{|D_2|} \sum_{\mathbf{x} \in D_2} \mathbb{1}_{h(\mathbf{x})=1}$
- $D_i = \{(\mathbf{x}, y) \in D : y = i\} \quad \forall i \in \mathcal{Y}$

**Definition (binary probabilistic ACC / PACC):**

Replace each occurrence of  $\mathbb{1}_{h(\mathbf{x})=1}$  with the soft classification  $s(\mathbf{x}) \in [0, 1]$

**Problem:**  $\hat{Q}(Y = 1)$  might be undefined or outside of  $[0, 1]$

<sup>4</sup> Tasche, “Fisher consistency for prior probability shift”, 2017.

## Counter example: Classify & Count



### Assignment [2 min]:

What would happen if we simply returned  $\hat{\mathbb{Q}}(h(X) = 1)$  as our estimate of  $\hat{\mathbb{Q}}(Y = 1)$ ?

## Counter example: Classify & Count



### Assignment [2 min]:

What would happen if we simply returned  $\hat{\mathbb{Q}}(h(X) = 1)$  as our estimate of  $\hat{\mathbb{Q}}(Y = 1)$ ?

**Answer:** on the population level, we would obtain

$$h(\mathbb{Q}(X)) = \mathbb{Q}(h(X) = 1)$$

## Counter example: Classify & Count



### Assignment [2 min]:

What would happen if we simply returned  $\hat{\mathbb{Q}}(h(X) = 1)$  as our estimate of  $\hat{\mathbb{Q}}(Y = 1)$ ?

**Answer:** on the population level, we would obtain

$$\begin{aligned}h(\mathbb{Q}(X)) &= \mathbb{Q}(h(X) = 1) \\ &= \text{TPR} \cdot \mathbb{Q}(Y = 1) + \text{FPR} \cdot (1 - \mathbb{Q}(Y = 1))\end{aligned}$$

## Counter example: Classify & Count



### Assignment [2 min]:

What would happen if we simply returned  $\hat{\mathbb{Q}}(h(X) = 1)$  as our estimate of  $\hat{\mathbb{Q}}(Y = 1)$ ?

**Answer:** on the population level, we would obtain

$$\begin{aligned}h(\mathbb{Q}(X)) &= \mathbb{Q}(h(X) = 1) \\ &= \text{TPR} \cdot \mathbb{Q}(Y = 1) + \text{FPR} \cdot (1 - \mathbb{Q}(Y = 1)) \\ &\neq \mathbb{Q}(Y = 1)\end{aligned}$$

if  $\text{TPR} \neq 1$  or if  $\text{FPR} \neq 0$ .

Hence, CC is **not** Fisher-consistent under PPS.



1. Problem statement
2. Desirable properties of quantifiers
3. Binary quantifiers

## 4. Multi-class quantifiers

5. Numerical optimization
6. Loss functions & data representations
7. Beyond linear systems of equations

## Example: from binary to multi-class (P)ACC



### Preliminaries:

- let  $\mathcal{Y} = \{1, 2, \dots, C\}$
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\Rightarrow \mathbb{Q}(h(X) = i) = \sum_{j \in \mathcal{Y}} \mathbb{Q}(h(X) = i \mid Y = j) \cdot \mathbb{Q}(Y = j) \quad \forall i \in \mathcal{Y} \quad (\text{just like before})$$

## Example: from binary to multi-class (P)ACC



### Preliminaries:

- let  $\mathcal{Y} = \{1, 2, \dots, C\}$
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\begin{aligned}\Rightarrow \mathbb{Q}(h(X) = i) &= \sum_{j \in \mathcal{Y}} \mathbb{Q}(h(X) = i \mid Y = j) \cdot \mathbb{Q}(Y = j) \quad \forall i \in \mathcal{Y} \quad (\text{just like before}) \\ &= \mathbf{M}_i^\top \mathbf{p}\end{aligned}$$

where

$$\mathbf{M}_i = \begin{pmatrix} \mathbb{Q}(h(X) = i \mid Y = 1) \\ \mathbb{Q}(h(X) = i \mid Y = 2) \\ \vdots \\ \mathbb{Q}(h(X) = i \mid Y = C) \end{pmatrix} \quad \mathbf{p} = \begin{pmatrix} \mathbb{Q}(Y = 1) \\ \mathbb{Q}(Y = 2) \\ \vdots \\ \mathbb{Q}(Y = C) \end{pmatrix}$$

## Example: from binary to multi-class (P)ACC



### Preliminaries:

- let  $\mathcal{Y} = \{1, 2, \dots, C\}$
- assume a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$

$$\begin{aligned}\Rightarrow \mathbb{Q}(h(X) = i) &= \sum_{j \in \mathcal{Y}} \mathbb{Q}(h(X) = i \mid Y = j) \cdot \mathbb{Q}(Y = j) \quad \forall i \in \mathcal{Y} \quad (\text{just like before}) \\ &= \mathbf{M}_i^\top \mathbf{p}\end{aligned}$$

where

$$\mathbf{M}_i = \begin{pmatrix} \mathbb{Q}(h(X) = i \mid Y = 1) \\ \mathbb{Q}(h(X) = i \mid Y = 2) \\ \vdots \\ \mathbb{Q}(h(X) = i \mid Y = C) \end{pmatrix} \quad \mathbf{p} = \begin{pmatrix} \mathbb{Q}(Y = 1) \\ \mathbb{Q}(Y = 2) \\ \vdots \\ \mathbb{Q}(Y = C) \end{pmatrix} \quad \mathbf{M} = \begin{pmatrix} \mathbf{M}_1^\top \\ \mathbf{M}_2^\top \\ \vdots \\ \mathbf{M}_C^\top \end{pmatrix} \quad \mathbf{q} = \begin{pmatrix} \mathbb{Q}(h(X) = 1) \\ \mathbb{Q}(h(X) = 2) \\ \vdots \\ \mathbb{Q}(h(X) = C) \end{pmatrix}^\top$$

such that  $\mathbf{q} = \mathbf{M}\mathbf{p}$

## Example: from binary to multi-class (P)ACC



### Synopsis:

- we have just seen how multi-class ACC and PACC yield systems of equations
- we have also seen how binary ACC and PACC differ in their computation of  $T\hat{P}R$ ,  $F\hat{P}R$ , and  $\hat{Q}(h(X) = 1)$

$$\text{e.g., } \hat{Q}(h(X) = 1) = \begin{cases} \frac{1}{|B|} \sum_{\mathbf{x} \in B} \mathbb{1}_{h(\mathbf{x})=1} & (\text{ACC}) \\ \frac{1}{|B|} \sum_{\mathbf{x} \in B} s(\mathbf{x}) & (\text{PACC}) \end{cases}$$

(they represent the data differently, either through  $h(\mathbf{x})$  or  $s(\mathbf{x})$ )

- we have not yet discussed how  $\mathbf{q} = \mathbf{M}\mathbf{p}$  can be solved

## Example: from binary to multi-class (P)ACC



### Synopsis:

- we have just seen how multi-class ACC and PACC yield systems of equations
- we have also seen how binary ACC and PACC differ in their computation of  $\hat{TPR}$ ,  $\hat{FPR}$ , and  $\hat{Q}(h(X) = 1)$

$$\text{e.g., } \hat{Q}(h(X) = 1) = \begin{cases} \frac{1}{|B|} \sum_{\mathbf{x} \in B} \mathbb{1}_{h(\mathbf{x})=1} & (\text{ACC}) \\ \frac{1}{|B|} \sum_{\mathbf{x} \in B} s(\mathbf{x}) & (\text{PACC}) \end{cases}$$

(they represent the data differently, either through  $h(\mathbf{x})$  or  $s(\mathbf{x})$ )

- we have not yet discussed how  $\mathbf{q} = \mathbf{M}\mathbf{p}$  can be solved

### Next steps [30 min]:

- generalize these concepts towards arbitrary data representations [10 min]
- discuss ways of solving  $\mathbf{q} = \mathbf{M}\mathbf{p}$
- define concrete representations (in addition to those of ACC and PACC)

# General systems of linear equations



**More generally:** any data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$  yields another system of equations  $\mathbf{q} = \mathbf{M}\mathbf{p}$  via

$$\mathbb{Q}(\Phi(X) = z) = \sum_{i \in \mathcal{Y}} \mathbb{Q}(\Phi(X) = z \mid Y = i) \cdot \mathbb{Q}(Y = i) \quad \forall z \in \mathcal{Z}$$

and any of these systems can suit the purpose of finding  $\mathbf{p}$ .

<sup>5</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# General systems of linear equations



**More generally:** any data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$  yields another system of equations  $\mathbf{q} = \mathbf{M}\mathbf{p}$  via

$$\mathbb{Q}(\Phi(X) = z) = \sum_{i \in \mathcal{Y}} \mathbb{Q}(\Phi(X) = z | Y = i) \cdot \mathbb{Q}(Y = i) \quad \forall z \in \mathcal{Z}$$

and any of these systems can suit the purpose of finding  $\mathbf{p}$ .

The solution of the system  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- is Fisher-consistent by construction
- is asymptotically consistent:<sup>5</sup>

$$\|\lambda(\mathbf{B}) - \mathbf{p}^*\|_2 \leq \underbrace{\frac{2k(2 + \sqrt{2 \log \frac{2C}{\delta}})}{\sqrt{\lambda_2}}}_{\text{representation } \Phi} \left( \underbrace{\frac{\|\mathbf{p}^*\|_2}{\sqrt{|\mathbf{D}|}}}_{\text{shift \& volume}} + \underbrace{\frac{1}{\sqrt{|\mathbf{B}|}}}_{\text{volume}} \right) \quad \text{where } \begin{cases} k & \text{constant s.t. } \|\Phi(\mathbf{x})\|_2 \leq k \quad \forall \mathbf{x} \in \mathcal{X} \\ \lambda_2 & \text{second-smallest eigenvalue of } \mathbf{G} \\ \delta & \text{desired probability} \end{cases}$$

<sup>5</sup> Dussap, Blanchard, and Chérief-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

## Counter example: one-vs-rest quantification



Can we not just use binary (P)ACC with for each binary task in an OVR decomposition?

$$\hat{Q}(Y = i) = \frac{\hat{Q}(h(X) = 1) - \text{F}\hat{\text{P}}\text{R}_i}{\text{T}\hat{\text{P}}\text{R}_i - \text{F}\hat{\text{P}}\text{R}_i} \quad \forall i \in \mathcal{Y}$$

<sup>6</sup> Gövert, “Fisher-Konsistenz für Quantification-Algorithmen”, 2023, supervised by M. Bunse and S. Mücke.

<sup>7</sup> Donyavi, Serapião, and Batista, “MC-SQ: A Highly Accurate Ensemble for Multi-class Quantification”, 2023.

## Counter example: one-vs-rest quantification

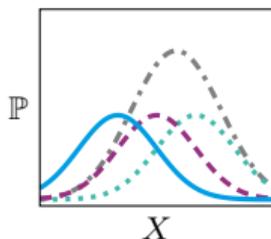


Can we not just use binary (P)ACC with for each binary task in an OVR decomposition?

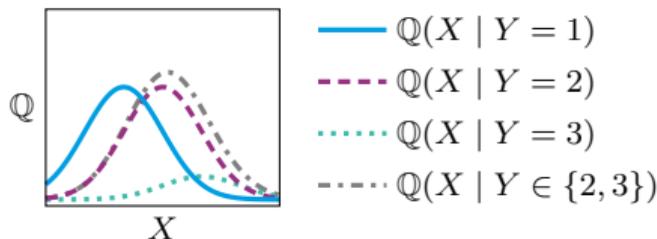
$$\hat{Q}(Y = i) = \frac{\hat{Q}(h(X) = 1) - \text{FPR}_i}{\text{TPR}_i - \text{FPR}_i} \quad \forall i \in \mathcal{Y}$$

**Assignment [2 min]:** What is the problem in the following situation?

training:  $\mathbb{P}(Y) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$



testing:  $\mathbb{Q}(Y) = (\frac{5}{11}, \frac{5}{11}, \frac{1}{11})$



<sup>6</sup> Gövert, “Fisher-Konsistenz für Quantification-Algorithmen”, 2023, supervised by M. Bunse and S. Mücke.

<sup>7</sup> Donyavi, Serapião, and Batista, “MC-SQ: A Highly Accurate Ensemble for Multi-class Quantification”, 2023.

## Counter example: one-vs-rest quantification

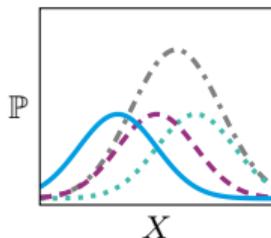


Can we not just use binary (P)ACC with for each binary task in an OVR decomposition?

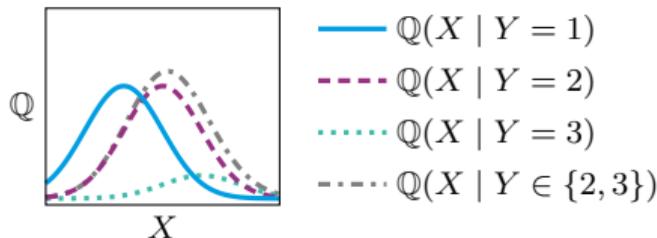
$$\hat{Q}(Y = i) = \frac{\hat{Q}(h(X) = 1) - \text{F}\hat{\text{P}}\text{R}_i}{\text{T}\hat{\text{P}}\text{R}_i - \text{F}\hat{\text{P}}\text{R}_i} \quad \forall i \in \mathcal{Y}$$

**Assignment [2 min]:** What is the problem in the following situation?

training:  $\mathbb{P}(Y) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$



testing:  $\mathbb{Q}(Y) = (\frac{5}{11}, \frac{5}{11}, \frac{1}{11})$



- PPS among  $C > 2$  classes leads to concept shift in OVR decompositions<sup>6,7</sup>
- hence, OVR quantification with Fisher-consistent binary quantifiers is not Fisher-consistent

<sup>6</sup> Gövert, "Fisher-Konsistenz für Quantification-Algorithmen", 2023, supervised by M. Bunse and S. Mücke.

<sup>7</sup> Donyavi, Serapião, and Batista, "MC-SQ: A Highly Accurate Ensemble for Multi-class Quantification", 2023.

# Synopsis



## We learned how to achieve Fisher consistency:

- define some data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$
- solve  $\mathbf{q} = \mathbf{M}\mathbf{p}$

## Next steps:

- discuss ways of solving  $\mathbf{q} = \mathbf{M}\mathbf{p}$
- define concrete representations  $\Phi$  (in addition to those of ACC and PACC)



1. Problem statement
2. Desirable properties of quantifiers
3. Binary quantifiers
4. Multi-class quantifiers

## 5. Numerical optimization

6. Loss functions & data representations
7. Beyond linear systems of equations

# Matrix inversion



**Goal:** solve  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- in other words, find  $\mathbf{p}$ , given  $\mathbf{q}$  and  $\mathbf{M}$
- $\mathbf{q}$  and  $\mathbf{M}$  are fully defined through  $\Phi$ ,  $\mathbf{B}$ , and  $\mathbf{D}$

<sup>8</sup> Mueller and Siltanen, *Linear and Nonlinear Inverse Problems with Practical Applications*, 2012.

# Matrix inversion



**Goal:** solve  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- in other words, find  $\mathbf{p}$ , given  $\mathbf{q}$  and  $\mathbf{M}$
- $\mathbf{q}$  and  $\mathbf{M}$  are fully defined through  $\Phi$ ,  $\mathbf{B}$ , and  $\mathbf{D}$

**Naive solution:** choose  $\hat{\mathbf{p}} = \mathbf{M}^{-1}\mathbf{q}$

- the inverse  $\mathbf{M}^{-1}$  is not guaranteed to exist ( $\mathbf{M}$  might not even be square)
- if  $\mathbf{M}^{-1}$  exists,  $\hat{\mathbf{p}}$  is not guaranteed to be in  $\Delta^{C-1}$  (an ad-hoc projection is necessary)

<sup>8</sup> Mueller and Siltanen, *Linear and Nonlinear Inverse Problems with Practical Applications*, 2012.



## Matrix inversion

**Goal:** solve  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- in other words, find  $\mathbf{p}$ , given  $\mathbf{q}$  and  $\mathbf{M}$
- $\mathbf{q}$  and  $\mathbf{M}$  are fully defined through  $\Phi$ ,  $\mathbf{B}$ , and  $\mathbf{D}$

**Naive solution:** choose  $\hat{\mathbf{p}} = \mathbf{M}^{-1}\mathbf{q}$

- the inverse  $\mathbf{M}^{-1}$  is not guaranteed to exist ( $\mathbf{M}$  might not even be square)
- if  $\mathbf{M}^{-1}$  exists,  $\hat{\mathbf{p}}$  is not guaranteed to be in  $\Delta^{C-1}$  (an ad-hoc projection is necessary)

**Naive improvement:** choose  $\hat{\mathbf{p}} = \mathbf{M}^\dagger\mathbf{q}$  with the Moore-Penrose pseudo-inverse  $\mathbf{M}^\dagger$

- $\mathbf{M}^\dagger$  always exists and  $\hat{\mathbf{p}}$  is a unique solution
- however,  $\hat{\mathbf{p}}$  is still not guaranteed to be in  $\Delta^{C-1}$
- $\hat{\mathbf{p}}$  is a minimum-norm least-squares solution<sup>8</sup> (while a minimum norm does not relate to quantification)

<sup>8</sup> Mueller and Siltanen, *Linear and Nonlinear Inverse Problems with Practical Applications*, 2012.

# Constrained optimization



**Proper solution:** constrain  $\hat{\mathbf{p}}$  to always be in  $\Delta^{C-1}$ , i.e.,

$$\hat{\mathbf{p}} = \arg \min_{\mathbf{p} \in \Delta^{C-1}} \ell(\mathbf{q}, \mathbf{M}\mathbf{p})$$

where  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  is a loss function, e.g., choose

$$\ell(\mathbf{q}, \mathbf{M}\mathbf{p}) = \|\mathbf{q} - \mathbf{M}\mathbf{p}\|_2^2 \quad (\text{least squares})$$

## Remarks:

- we will soon learn about other loss functions
- a straightforward implementation requires **constrained optimization algorithms**

# Implicit constraints



Can we use **unconstrained** optimization algorithms?

**Yes:**<sup>9</sup> use the soft-max operator  $\sigma : \mathbb{R}^{C-1} \rightarrow \Delta^{C-1}$  and optimize over log-odds  $\mathbf{l} \in \mathbb{R}^{C-1}$ , i.e.,

$$\hat{\mathbf{p}} = \sigma(\mathbf{l}^*)$$

$$\mathbf{l}^* = \arg \min_{\mathbf{l} \in \mathbb{R}^{C-1}} \ell(\mathbf{q}, \mathbf{M}\sigma(\mathbf{l}))$$

<sup>9</sup> Bunse, “On Multi-Class Extensions of Adjusted Classify and Count”, 2022.

# Implicit constraints



Can we use **unconstrained** optimization algorithms?

**Yes:**<sup>9</sup> use the soft-max operator  $\sigma : \mathbb{R}^{C-1} \rightarrow \Delta^{C-1}$  and optimize over log-odds  $\mathbf{l} \in \mathbb{R}^{C-1}$ , i.e.,

$$\hat{\mathbf{p}} = \sigma(\mathbf{l}^*)$$

$$\mathbf{l}^* = \arg \min_{\mathbf{l} \in \mathbb{R}^{C-1}} \ell(\mathbf{q}, \mathbf{M}\sigma(\mathbf{l}))$$

$$[\sigma(\mathbf{l})]_i = \begin{cases} \frac{1}{1 + \sum_{j=1}^{C-1} \exp(\mathbf{l}_j)} & \text{if } i = 1 \\ \frac{\exp(\mathbf{l}_{i-1})}{1 + \sum_{j=1}^{C-1} \exp(\mathbf{l}_j)} & \forall i \in \{2, 3, \dots, C\} \end{cases}$$

<sup>9</sup> Bunse, "On Multi-Class Extensions of Adjusted Classify and Count", 2022.

# Synopsis



## Components of a quantification algorithm:

- a data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$
- a loss function  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$
- an optimization algorithm (matrix inversion does not suffice!)

## Next steps:

- define concrete choices of  $\Phi$  and  $\ell$  (in addition to those of ACC and PACC)
- discuss algorithms beyond solutions of  $\mathbf{q} = \mathbf{M}\mathbf{p}$



1. Problem statement
2. Desirable properties of quantifiers
3. Binary quantifiers
4. Multi-class quantifiers
5. Numerical optimization

## **6. Loss functions & data representations**

7. Beyond linear systems of equations

## (Probabilistic) Adjusted Classify & Count



Loss function:

$$\ell(\mathbf{q}, \mathbf{Mp}) = \|\mathbf{q} - \mathbf{Mp}\|_2^2 \quad (\text{least squares})$$

Representation:<sup>10,11</sup>

$$\Phi(\mathbf{x}) = \begin{cases} \mathbb{1}_{h(\mathbf{x})} \in \{0, 1\}^C & \text{ACC (a one-hot encoding of } h(\mathbf{x})) \\ s(\mathbf{x}) \in \Delta^{C-1} & \text{PACC} \end{cases}$$

where

- $h : \mathcal{X} \rightarrow \mathcal{Y}$  a “hard” classifier such that  $h(\mathbf{x}) = \hat{y}$
- $s : \mathcal{X} \rightarrow \Delta^{C-1}$  a “soft” classifier such that  $s(\mathbf{x}) \approx \mathbb{P}(Y | \mathbf{x})$

<sup>10</sup> Firat, “Unified Framework for Quantification”, 2016.

<sup>11</sup> Bunse, “Unification of Algorithms for Quantification and Unfolding”, 2022.

## Distribution matching: HDx and HDy



Loss function:

$$\ell(\mathbf{q}, \mathbf{M}_p) = \frac{1}{d} \sum_{i=1}^d \text{HD}(\mathbf{q}_{i\bullet}, \mathbf{M}_{i\bullet p}) \quad \text{where} \quad \text{HD}(\mathbf{a}, \mathbf{b}) = \sqrt{\sum_i^b \left( \sqrt{\mathbf{a}_i} - \sqrt{\mathbf{b}_i} \right)^2}$$

Representation:

$$\Phi(\mathbf{x}) = \begin{cases} \left( \mathbb{1}_{b_1(\mathbf{x}_1)}, \mathbb{1}_{b_2(\mathbf{x}_2)}, \dots, \mathbb{1}_{b_d(\mathbf{x}_d)} \right) \in \{0, 1\}^{Bd} & \text{HDx} \\ \left( \mathbb{1}_{b_1([s(\mathbf{x})]_1)}, \mathbb{1}_{b_2([s(\mathbf{x})]_2)}, \dots, \mathbb{1}_{b_d([s(\mathbf{x})]_d)} \right) \in \{0, 1\}^{BC} & \text{HDy} \end{cases}$$

where  $b_i : \mathbb{R} \rightarrow \{1, 2, \dots, B\}$  is a binning of the  $i$ -th feature (or class probability)

**Problem:** HD is not twice differentiable  $\Rightarrow$  prefer HD<sup>2</sup> instead.<sup>12</sup>

<sup>12</sup> Bunse, “qunifold: Composable Quantification and Unfolding Methods in Python”, 2023.

# Kernel mean matching: EDx, EDy, and others



Representation:

$$[\Phi(\mathbf{x})]_i = \frac{1}{|D_i|} \sum_{\mathbf{x}' \in D_i} K(\mathbf{x}, \mathbf{x}')$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is some kernel between data items, e.g.,

- $\|\mathbf{x} - \mathbf{x}'\|_2$  (Euclidean distance; EDx<sup>13</sup>)
- $\sum_{i=1}^{C-1} \left| \sum_{j=1}^i [s(\mathbf{x})]_j - [s(\mathbf{x}')]_j \right|$  (Earth Mover's Distance; EDy<sup>14</sup>)
- $\exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|_2}{(2\sigma)^2}\right)$  (Gaussian kernel)

<sup>13</sup> Kawakubo, Plessis, and Sugiyama, “Computationally Efficient Class-Prior Estimation under Class Balance Change Using Energy Distance”, 2016.

<sup>14</sup> Castaño et al., “Matching Distributions Algorithms Based on the Earth Mover's Distance for Ordinal Quantification”, 2022.

<sup>15</sup> Dussap, Blanchard, and Chérief-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Kernel mean matching: EDx, EDy, and others



**Representation:**

$$[\Phi(\mathbf{x})]_i = \frac{1}{|D_i|} \sum_{\mathbf{x}' \in D_i} K(\mathbf{x}, \mathbf{x}')$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is some kernel between data items, e.g.,

- $\|\mathbf{x} - \mathbf{x}'\|_2$  (Euclidean distance; EDx<sup>13</sup>)
- $\sum_{i=1}^{C-1} \left| \sum_{j=1}^i [s(\mathbf{x})]_j - [s(\mathbf{x}')]_j \right|$  (Earth Mover's Distance; EDy<sup>14</sup>)
- $\exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|_2}{(2\sigma)^2}\right)$  (Gaussian kernel)

**Problem:** a high computational cost  $\Rightarrow$  use a random Fourier approximation<sup>15</sup>

<sup>13</sup> Kawakubo, Plessis, and Sugiyama, “Computationally Efficient Class-Prior Estimation under Class Balance Change Using Energy Distance”, 2016.

<sup>14</sup> Castaño et al., “Matching Distributions Algorithms Based on the Earth Mover's Distance for Ordinal Quantification”, 2022.

<sup>15</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Kernel mean matching: EDx, EDy, and others



**Representation:**

$$[\Phi(\mathbf{x})]_i = \frac{1}{|D_i|} \sum_{\mathbf{x}' \in D_i} K(\mathbf{x}, \mathbf{x}')$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is some kernel between data items, e.g.,

- $\|\mathbf{x} - \mathbf{x}'\|_2$  (Euclidean distance; EDx<sup>13</sup>)
- $\sum_{i=1}^{C-1} \left| \sum_{j=1}^i [s(\mathbf{x})]_j - [s(\mathbf{x}')]_j \right|$  (Earth Mover's Distance; EDy<sup>14</sup>)
- $\exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|_2}{(2\sigma)^2}\right)$  (Gaussian kernel)

**Problem:** a high computational cost  $\Rightarrow$  use a random Fourier approximation<sup>15</sup>

**Loss function (EDx and EDy):**  $\ell(\mathbf{q}, \mathbf{M}\mathbf{p}) = 2\mathbf{p}^T \mathbf{q} - \mathbf{p}^T \mathbf{M}\mathbf{p}$  (any other loss is possible)

<sup>13</sup> Kawakubo, Plessis, and Sugiyama, “Computationally Efficient Class-Prior Estimation under Class Balance Change Using Energy Distance”, 2016.

<sup>14</sup> Castaño et al., “Matching Distributions Algorithms Based on the Earth Mover’s Distance for Ordinal Quantification”, 2022.

<sup>15</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Regularization



So far, we've been very concerned about consistency. But what if

- the data volume is small?
- we've strong assumptions about how the predictions should look like?

<sup>16</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

# Regularization



So far, we've been very concerned about consistency. But what if

- the data volume is small?
- we've strong assumptions about how the predictions should look like?

## Regularization:

$$\ell'(\mathbf{q}, \mathbf{M}\mathbf{p}) = \ell(\mathbf{q}, \mathbf{M}\mathbf{p}) + \tau \cdot r(\mathbf{p})$$

- $\tau \geq 0$  is the regularization impact (i.e., a hyper-parameter that needs to be optimized)
- $r : \Delta^{C-1} \rightarrow \mathbb{R}$  is a regularization term that penalizes any deviation from our assumptions

<sup>16</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

# Regularization



So far, we've been very concerned about consistency. But what if

- the data volume is small?
- we've strong assumptions about how the predictions should look like?

## Regularization:

$$\ell'(\mathbf{q}, \mathbf{M}\mathbf{p}) = \ell(\mathbf{q}, \mathbf{M}\mathbf{p}) + \tau \cdot r(\mathbf{p})$$

- $\tau \geq 0$  is the regularization impact (i.e., a hyper-parameter that needs to be optimized)
- $r : \Delta^{C-1} \rightarrow \mathbb{R}$  is a regularization term that penalizes any deviation from our assumptions

## Tikhonov regularization:<sup>16</sup>

$$r(\mathbf{p}) = \frac{1}{2}(\mathbf{C}\mathbf{p})^2 = \begin{cases} \frac{1}{2} \sum_{i=2}^{C-1} (-\mathbf{p}_{i-1} + 2\mathbf{p}_i - \mathbf{p}_{i+1})^2 & \text{ordinal quantification} \\ \frac{1}{2} \sum_{i=1}^{C-1} (\mathbf{p}_i - \mathbf{p}_{i+1})^2 & \text{deviation from } \mathbf{p}_i = \frac{1}{C} \end{cases}$$

<sup>16</sup> Bunse et al., "Regularization-based Methods for Ordinal Quantification", 2024.

# Synopsis



**Many quantification algorithms are combinations of:**

- a data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$
- a loss function  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$
- an optimization algorithm

We have still omitted many methods from this family (ReadMe, PDF, unfolding methods, ...)

# Synopsis



**Many quantification algorithms are combinations of:**

- a data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$
- a loss function  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$
- an optimization algorithm

We have still omitted many methods from this family (ReadMe, PDF, unfolding methods, ...)

**We can alter these algorithms by:**

- re-combining their  $\Phi$  and  $\ell$
- approximating their representations
- adding regularization

**Next step:** discuss algorithms beyond solutions of  $\mathbf{q} = \mathbf{M}\mathbf{p}$



1. Problem statement
2. Desirable properties of quantifiers
3. Binary quantifiers
4. Multi-class quantifiers
5. Numerical optimization
6. Loss functions & data representations

## **7. Beyond linear systems of equations**

# Expectation maximization



## Preliminaries:

$$Q(\mathbf{x} | y) \stackrel{\text{Bayes}}{=} \frac{Q(y | \mathbf{x}) \cdot Q(\mathbf{x})}{Q(y)} \stackrel{\text{PPS}}{=} P(\mathbf{x} | y) \stackrel{\text{Bayes}}{=} \frac{P(y | \mathbf{x}) \cdot P(\mathbf{x})}{P(y)} \quad \forall (\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$$

<sup>17</sup> Saerens, Latinne, and Decaestecker, "Adjusting the Outputs of a Classifier to New a Priori Probabilities: A Simple Procedure", 2002.

# Expectation maximization



## Preliminaries:

$$Q(\mathbf{x} | y) \stackrel{\text{Bayes}}{=} \frac{Q(y | \mathbf{x}) \cdot Q(\mathbf{x})}{Q(y)} \stackrel{\text{PPS}}{=} P(\mathbf{x} | y) \stackrel{\text{Bayes}}{=} \frac{P(y | \mathbf{x}) \cdot P(\mathbf{x})}{P(y)} \quad \forall (\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$$

$$\Rightarrow Q(y | \mathbf{x}) = \frac{\frac{Q(y)}{P(y)} \cdot P(y | \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \frac{Q(y')}{P(y')} \cdot P(y' | \mathbf{x})} \quad (\text{the original purpose of this method})$$

<sup>17</sup> Saerens, Latinne, and Decaestecker, "Adjusting the Outputs of a Classifier to New a Priori Probabilities: A Simple Procedure", 2002.

# Expectation maximization



## Preliminaries:

$$\begin{aligned} \mathbb{Q}(\mathbf{x} | y) &\stackrel{\text{Bayes}}{=} \frac{\mathbb{Q}(y | \mathbf{x}) \cdot \mathbb{Q}(\mathbf{x})}{\mathbb{Q}(y)} \stackrel{\text{PPS}}{=} \mathbb{P}(\mathbf{x} | y) \stackrel{\text{Bayes}}{=} \frac{\mathbb{P}(y | \mathbf{x}) \cdot \mathbb{P}(\mathbf{x})}{\mathbb{P}(y)} \quad \forall (\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y} \\ \Rightarrow \mathbb{Q}(y | \mathbf{x}) &= \frac{\frac{\mathbb{Q}(y)}{\mathbb{P}(y)} \cdot \mathbb{P}(y | \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \frac{\mathbb{Q}(y')}{\mathbb{P}(y')} \cdot \mathbb{P}(y' | \mathbf{x})} \quad (\text{the original purpose of this method}) \end{aligned}$$

**SLD / EMQ:**<sup>17</sup> repeat the E-step and the M-step until convergence.

**Initialize:**  $\mathbf{p}^{(0)} \leftarrow \hat{\mathbb{P}}(Y)$   
 $\hat{\mathbb{Q}}^{(0)}(y | \mathbf{x}) \leftarrow \hat{\mathbb{P}}(y | \mathbf{x}) \quad \forall y \in \mathcal{Y}, \mathbf{x} \in \mathcal{B}$

**E-step:**  $[\mathbf{p}^{(k)}]_i \leftarrow \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} \hat{\mathbb{Q}}^{(k-1)}(Y = i | \mathbf{x})$

**M-step:**  $\hat{\mathbb{Q}}^{(k)}(y | \mathbf{x}) \leftarrow \frac{\frac{[\mathbf{p}^{(k)}]_y}{\hat{\mathbb{P}}(y)} \cdot \hat{\mathbb{P}}(y | \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \frac{[\mathbf{p}^{(k)}]_{y'}}{\hat{\mathbb{P}}(y')} \cdot \hat{\mathbb{P}}(y' | \mathbf{x})} \quad \forall y \in \mathcal{Y}, \mathbf{x} \in \mathcal{B}$

<sup>17</sup> Saerens, Latinne, and Decaestecker, "Adjusting the Outputs of a Classifier to New a Priori Probabilities: A Simple Procedure", 2002.

# Expectation maximization



**Properties:** SLD / EMQ is...

- Fisher consistent<sup>18</sup>
- equivalent to the maximum likelihood quantifier<sup>19</sup> (which is to be presented next)
- and it maintains per-example contributions  $\hat{\mathbb{Q}}(Y = i | \mathbf{x})$

<sup>18</sup> Tasche, “Fisher consistency for prior probability shift”, 2017.

<sup>19</sup> Alexandari, Kundaje, and Shrikumar, “Maximum Likelihood with Bias-Corrected Calibration is Hard-To-Beat at Label Shift Adaptation”, 2020.

# Maximum likelihood



**Likelihood principle:**

$$\begin{aligned}\mathcal{L}(\mathbf{p} \mid B) &= Q(B \mid \mathbf{p}) \\ &= \prod_{\mathbf{x} \in B} Q(\mathbf{x} \mid \mathbf{p}) \\ &\stackrel{\text{PPS}}{=} \prod_{\mathbf{x} \in B} \sum_{y \in \mathcal{Y}} \mathbb{P}(\mathbf{x} \mid y) \cdot \mathbf{p}_y\end{aligned}$$

# Maximum likelihood



Likelihood principle:

$$\begin{aligned}\mathcal{L}(\mathbf{p} \mid B) &= \mathbb{Q}(B \mid \mathbf{p}) \\ &= \prod_{\mathbf{x} \in B} \mathbb{Q}(\mathbf{x} \mid \mathbf{p}) \\ &\stackrel{\text{PPS}}{=} \prod_{\mathbf{x} \in B} \sum_{y \in \mathcal{Y}} \mathbb{P}(\mathbf{x} \mid y) \cdot \mathbf{p}_y\end{aligned}$$

$$\begin{aligned}\Rightarrow -\log \mathcal{L}(\mathbf{p} \mid B) &= -\sum_{\mathbf{x} \in B} \log \sum_{y \in \mathcal{Y}} \mathbb{P}(\mathbf{x} \mid y) \cdot \mathbf{p}_y \\ &\propto -\sum_{\mathbf{x} \in B} \log \sum_{y \in \mathcal{Y}} \frac{\mathbb{P}(y \mid \mathbf{x})}{\mathbb{P}(y)} \cdot \mathbf{p}_y\end{aligned}$$

# Maximum likelihood



Likelihood principle:

$$\begin{aligned}\mathcal{L}(\mathbf{p} \mid B) &= \mathbb{Q}(B \mid \mathbf{p}) \\ &= \prod_{\mathbf{x} \in B} \mathbb{Q}(\mathbf{x} \mid \mathbf{p}) \\ &\stackrel{\text{PPS}}{=} \prod_{\mathbf{x} \in B} \sum_{y \in \mathcal{Y}} \mathbb{P}(\mathbf{x} \mid y) \cdot \mathbf{p}_y\end{aligned}$$

$$\begin{aligned}\Rightarrow -\log \mathcal{L}(\mathbf{p} \mid B) &= -\sum_{\mathbf{x} \in B} \log \sum_{y \in \mathcal{Y}} \mathbb{P}(\mathbf{x} \mid y) \cdot \mathbf{p}_y \\ &\propto -\sum_{\mathbf{x} \in B} \log \sum_{y \in \mathcal{Y}} \frac{\mathbb{P}(y \mid \mathbf{x})}{\mathbb{P}(y)} \cdot \mathbf{p}_y\end{aligned}$$

Therefore, choose  $\hat{\mathbf{p}} = \arg \min_{\mathbf{p} \in \Delta^{C-1}} -\sum_{\mathbf{x} \in B} \log \sum_{y \in \mathcal{Y}} \frac{\hat{\mathbb{P}}(y \mid \mathbf{x})}{\hat{\mathbb{P}}(y)} \cdot \mathbf{p}_y$

## Continuous representations



**KDEy:**<sup>20</sup> represent all probabilities through kernel density estimates (KDEs), i.e.,

$$\hat{\mathbb{Q}}(\mathbf{x}) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}' \in \mathcal{B}} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \text{and} \quad \hat{\mathbb{Q}}(\mathbf{x} | y) = \frac{1}{|\mathcal{D}_y|} \sum_{\mathbf{x}' \in \mathcal{D}_y} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \forall \mathbf{x} \in \mathcal{X}$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel function (e.g., a Gaussian kernel with some bandwidth)

<sup>20</sup> Moreo, González, and Coz, “Kernel Density Estimation for Multiclass Quantification”, 2024.

## Continuous representations



**KDEy:**<sup>20</sup> represent all probabilities through kernel density estimates (KDEs), i.e.,

$$\hat{Q}(\mathbf{x}) = \frac{1}{|B|} \sum_{\mathbf{x}' \in B} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \text{and} \quad \hat{Q}(\mathbf{x} | y) = \frac{1}{|D_y|} \sum_{\mathbf{x}' \in D_y} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \forall \mathbf{x} \in \mathcal{X}$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel function (e.g., a Gaussian kernel with some bandwidth)

**Remark:** this is different from KMM, where  $[\Phi(\mathbf{x})]_i = \frac{1}{|D_i|} \sum_{\mathbf{x}' \in D_i} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \forall \mathbf{x} \in D \cup B$

<sup>20</sup> Moreo, González, and Coz, “Kernel Density Estimation for Multiclass Quantification”, 2024.

## Continuous representations



**KDEy:**<sup>20</sup> represent all probabilities through kernel density estimates (KDEs), i.e.,

$$\hat{Q}(\mathbf{x}) = \frac{1}{|B|} \sum_{\mathbf{x}' \in B} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \text{and} \quad \hat{Q}(\mathbf{x} | y) = \frac{1}{|D_y|} \sum_{\mathbf{x}' \in D_y} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \forall \mathbf{x} \in \mathcal{X}$$

where  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel function (e.g., a Gaussian kernel with some bandwidth)

**Remark:** this is different from KMM, where  $[\Phi(\mathbf{x})]_i = \frac{1}{|D_i|} \sum_{\mathbf{x}' \in D_i} K(s(\mathbf{x}), s(\mathbf{x}')) \quad \forall \mathbf{x} \in D \cup B$

**Optimization task:** is determined by the choice of loss function.

- **Losses with closed-form solutions** lead to specific tasks (e.g., Cauchy-Schwarz)
- **Negative log-likelihood** leads to the maximum likelihood estimator (with a KDE representation)
- **MC-sampled losses** lead to  $q = \text{Mp}$  tasks

<sup>20</sup> Moreo, González, and Coz, “Kernel Density Estimation for Multiclass Quantification”, 2024.

## Symmetric learning



So far, we have assumed a training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ , but what if we have

$$D' = \left\{ (B_i, \mathbf{p}_i) \in \cup_{m=1}^{\infty} \mathcal{X}^m \times \Delta^{C-1} \right\}_{i=1}^n$$

<sup>21</sup> Pérez-Mon et al., “Quantification using Permutation-Invariant Networks based on Histograms”, 2024.

## Symmetric learning



So far, we have assumed a training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ , but what if we have

$$D' = \left\{ (B_i, \mathbf{p}_i) \in \cup_{m=1}^{\infty} \mathcal{X}^m \times \Delta^{C-1} \right\}_{i=1}^n$$

**Requirements:** the representations of the  $B_i$  need to be...

<sup>21</sup> Pérez-Mon et al., “Quantification using Permutation-Invariant Networks based on Histograms”, 2024.

## Symmetric learning



So far, we have assumed a training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ , but what if we have

$$D' = \left\{ (B_i, \mathbf{p}_i) \in \cup_{m=1}^{\infty} \mathcal{X}^m \times \Delta^{C-1} \right\}_{i=1}^n$$

**Requirements:** the representations of the  $B_i$  need to be...

- variable-size
- permutation-invariant

<sup>21</sup> Pérez-Mon et al., “Quantification using Permutation-Invariant Networks based on Histograms”, 2024.

## Symmetric learning



So far, we have assumed a training set  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ , but what if we have

$$D' = \left\{ (B_i, \mathbf{p}_i) \in \cup_{m=1}^{\infty} \mathcal{X}^m \times \Delta^{C-1} \right\}_{i=1}^n$$

**Requirements:** the representations of the  $B_i$  need to be...

- variable-size
- permutation-invariant

**HistNetQ:**<sup>21</sup>

$$\ell(\theta) = \frac{1}{|D'|} \sum_{(B, \mathbf{p}) \in D'} \text{RAE}(\lambda_{\theta}(B), \mathbf{p})$$

where  $\lambda_{\theta}$  a neural network with a differentiable histogram layer

<sup>21</sup> Pérez-Mon et al., “Quantification using Permutation-Invariant Networks based on Histograms”, 2024.

# Ensembling



**Idea:** compute a central tendency (mean, median) of multiple predictions.

- multiple classifiers within different quantifiers (MC-MQ) or within duplicates of the same (MC-SQ)<sup>22</sup>

<sup>22</sup> Donyavi, Serapião, and Batista, “MC-SQ and MC-MQ: Ensembles for Multi-class Quantification”, 2024.

<sup>23</sup> Pérez-Gállego, Quevedo, and Coz, “Using ensembles for problems with characterizable changes in data distribution: A case study on quantification”, 2017.

<sup>24</sup> Pérez-Gállego et al., “Dynamic ensemble selection for quantification tasks”, 2019.

<sup>25</sup> Janssen, “Ensembles für Quantification durch Konkatenieren von Quantifier-Modellen”, 2024, supervised by M. Bunse and S. Buschjäger.

# Ensembling



**Idea:** compute a central tendency (mean, median) of multiple predictions.

- multiple classifiers within different quantifiers (MC-MQ) or within duplicates of the same (MC-SQ)<sup>22</sup>
- for each member, use all data or use different subsamples<sup>23</sup>

<sup>22</sup> Donyavi, Serapião, and Batista, “MC-SQ and MC-MQ: Ensembles for Multi-class Quantification”, 2024.

<sup>23</sup> Pérez-Gállego, Quevedo, and Coz, “Using ensembles for problems with characterizable changes in data distribution: A case study on quantification”, 2017.

<sup>24</sup> Pérez-Gállego et al., “Dynamic ensemble selection for quantification tasks”, 2019.

<sup>25</sup> Janssen, “Ensembles für Quantification durch Konkatenieren von Quantifier-Modellen”, 2024, supervised by M. Bunse and S. Buschjäger.

# Ensembling



**Idea:** compute a central tendency (mean, median) of multiple predictions.

- multiple classifiers within different quantifiers (MC-MQ) or within duplicates of the same (MC-SQ)<sup>22</sup>
- for each member, use all data or use different subsamples<sup>23</sup>
- maintain all members, select a subset at training time, or select a subset at prediction time<sup>24</sup>

<sup>22</sup> Donyavi, Serapião, and Batista, “MC-SQ and MC-MQ: Ensembles for Multi-class Quantification”, 2024.

<sup>23</sup> Pérez-Gállego, Quevedo, and Coz, “Using ensembles for problems with characterizable changes in data distribution: A case study on quantification”, 2017.

<sup>24</sup> Pérez-Gállego et al., “Dynamic ensemble selection for quantification tasks”, 2019.

<sup>25</sup> Janssen, “Ensembles für Quantification durch Konkatenieren von Quantifier-Modellen”, 2024, supervised by M. Bunse and S. Buschjäger.

# Ensembling



**Idea:** compute a central tendency (mean, median) of multiple predictions.

- multiple classifiers within different quantifiers (MC-MQ) or within duplicates of the same (MC-SQ)<sup>22</sup>
- for each member, use all data or use different subsamples<sup>23</sup>
- maintain all members, select a subset at training time, or select a subset at prediction time<sup>24</sup>
- concatenate  $\Phi(\mathbf{x}) = (\Phi_1(\mathbf{x}), \Phi_2(\mathbf{x}), \dots, \Phi_E(\mathbf{x}))$  and minimize the loss once<sup>25</sup>

**Open issue:** under which circumstances are ensembles *provably* better than single models?

<sup>22</sup> Donyavi, Serapião, and Batista, “MC-SQ and MC-MQ: Ensembles for Multi-class Quantification”, 2024.

<sup>23</sup> Pérez-Gállego, Quevedo, and Coz, “Using ensembles for problems with characterizable changes in data distribution: A case study on quantification”, 2017.

<sup>24</sup> Pérez-Gállego et al., “Dynamic ensemble selection for quantification tasks”, 2019.

<sup>25</sup> Janssen, “Ensembles für Quantification durch Konkatenieren von Quantifier-Modellen”, 2024, supervised by M. Bunse and S. Buschjäger.



**Conclusion:  
supervised methods for quantification**

## Conclusion: supervised methods for quantification



**Goal:** under PPS, find a quantifier  $\lambda : \cup_{m=1}^{\infty} \mathcal{X}^m \rightarrow \Delta^{C-1}$  that is

- Fisher-consistent
- has low estimation error // what about other settings than PPS?

## Conclusion: supervised methods for quantification



**Goal:** under PPS, find a quantifier  $\lambda : \cup_{m=1}^{\infty} \mathcal{X}^m \rightarrow \Delta^{C-1}$  that is

- Fisher-consistent
- has low estimation error // what about other settings than PPS?

**Linear systems of equations:** most algorithms can be expressed as solutions of some  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- choose a data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$  // what makes a good  $\Phi$ ?
- choose a loss function  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  // what makes a good  $\ell$ ?
- choose an optimization algorithm (constrained or soft-max)

## Conclusion: supervised methods for quantification



**Goal:** under PPS, find a quantifier  $\lambda : \cup_{m=1}^{\infty} \mathcal{X}^m \rightarrow \Delta^{C-1}$  that is

- Fisher-consistent
- has low estimation error // what about other settings than PPS?

**Linear systems of equations:** most algorithms can be expressed as solutions of some  $\mathbf{q} = \mathbf{M}\mathbf{p}$

- choose a data representation  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$  // what makes a good  $\Phi$ ?
- choose a loss function  $\ell : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  // what makes a good  $\ell$ ?
- choose an optimization algorithm (constrained or soft-max)

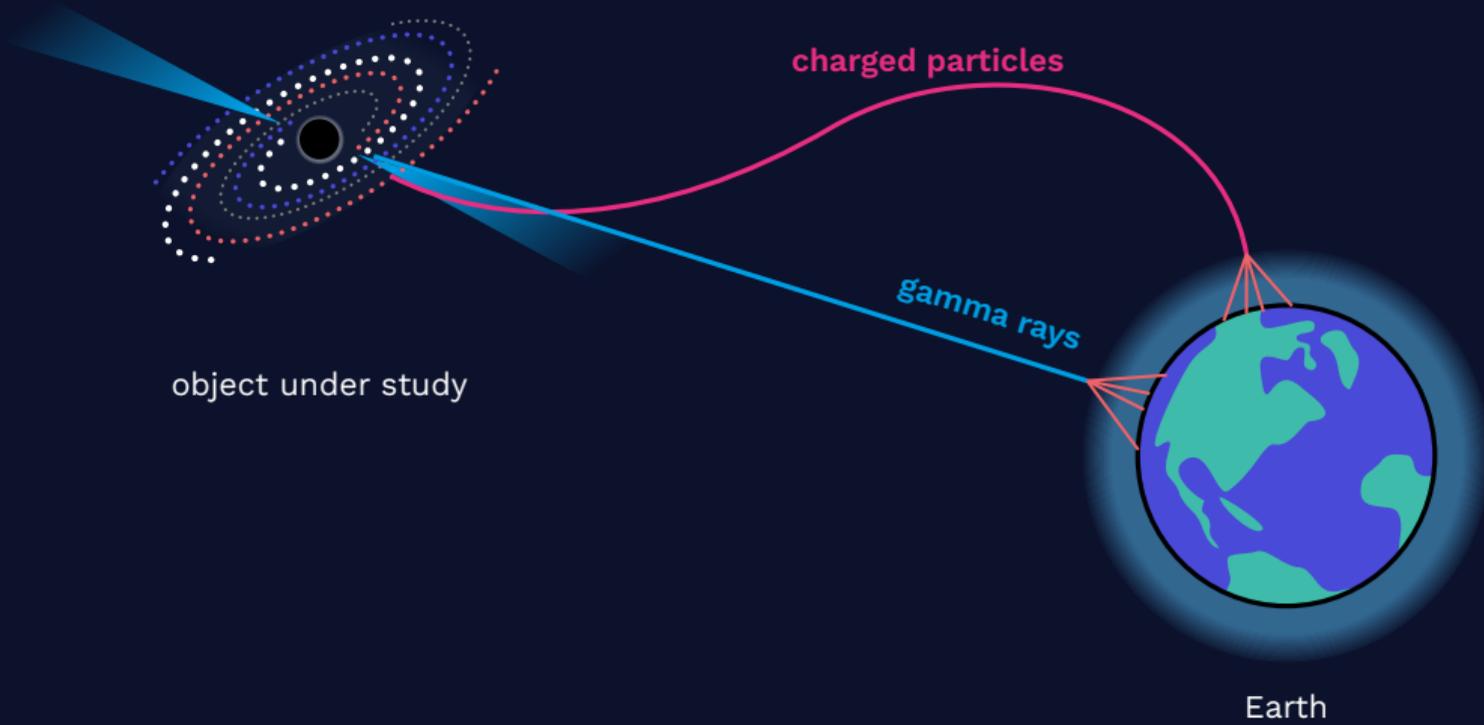
**Notable other methods:**

- SLD / EMQ
- maximum likelihood // can we bound its error?
- continuous representations
- symmetric learning
- ensembles // are they provably better?

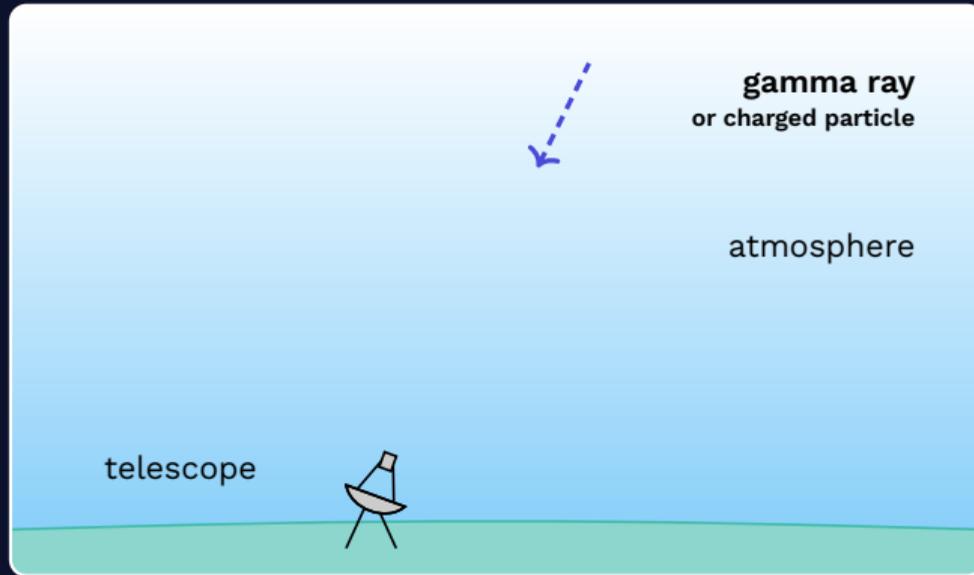


# Advanced topics for experimental physics

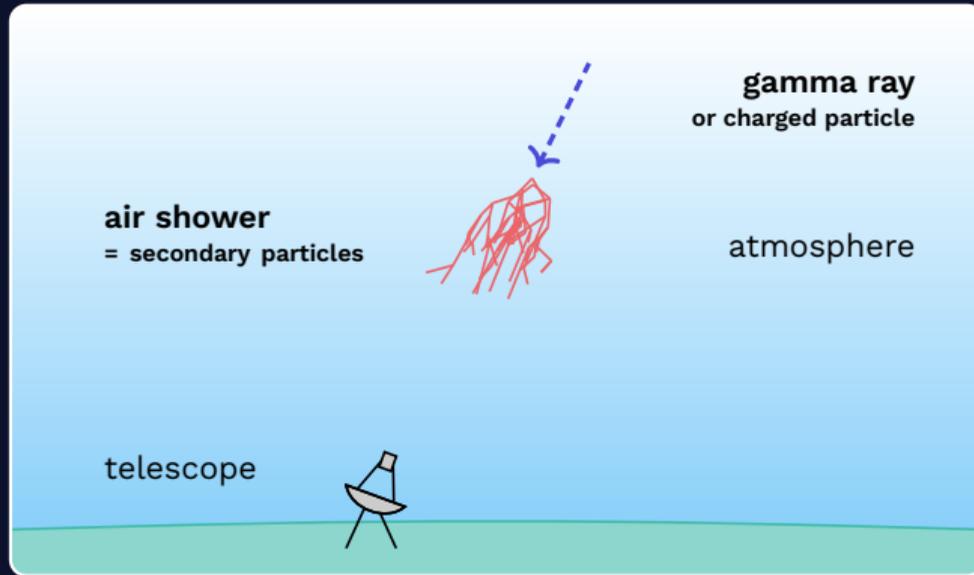
# Example: astro-particle physics



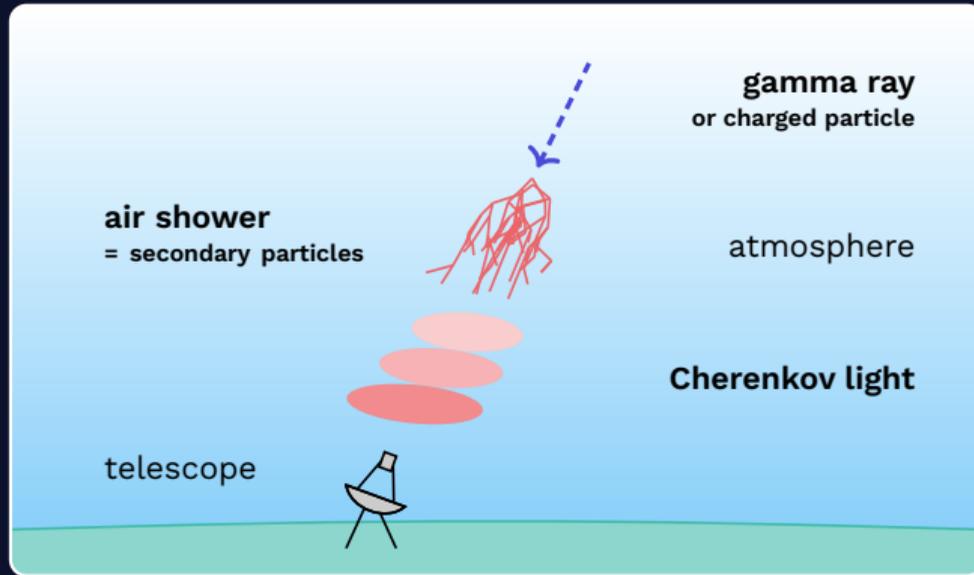
## Example: astro-particle physics



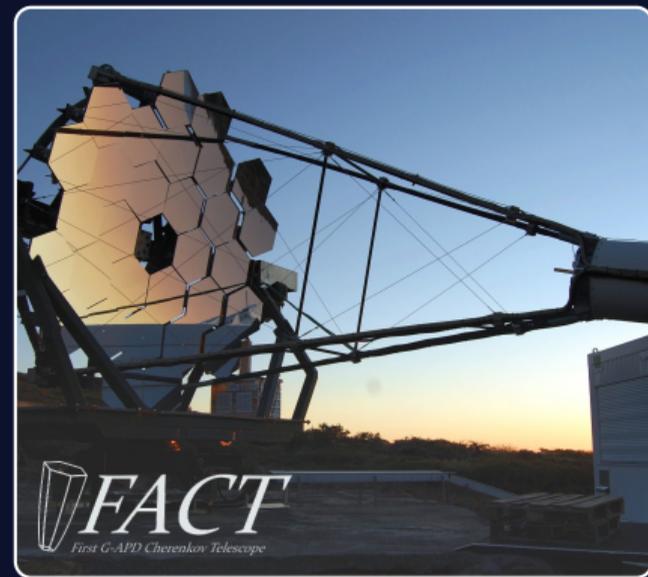
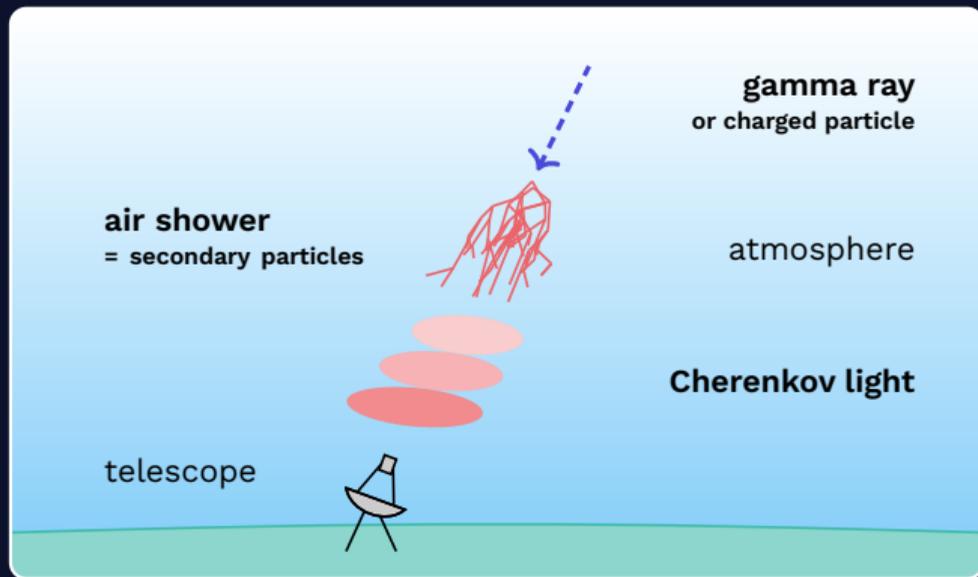
# Example: astro-particle physics



## Example: astro-particle physics



## Example: astro-particle physics

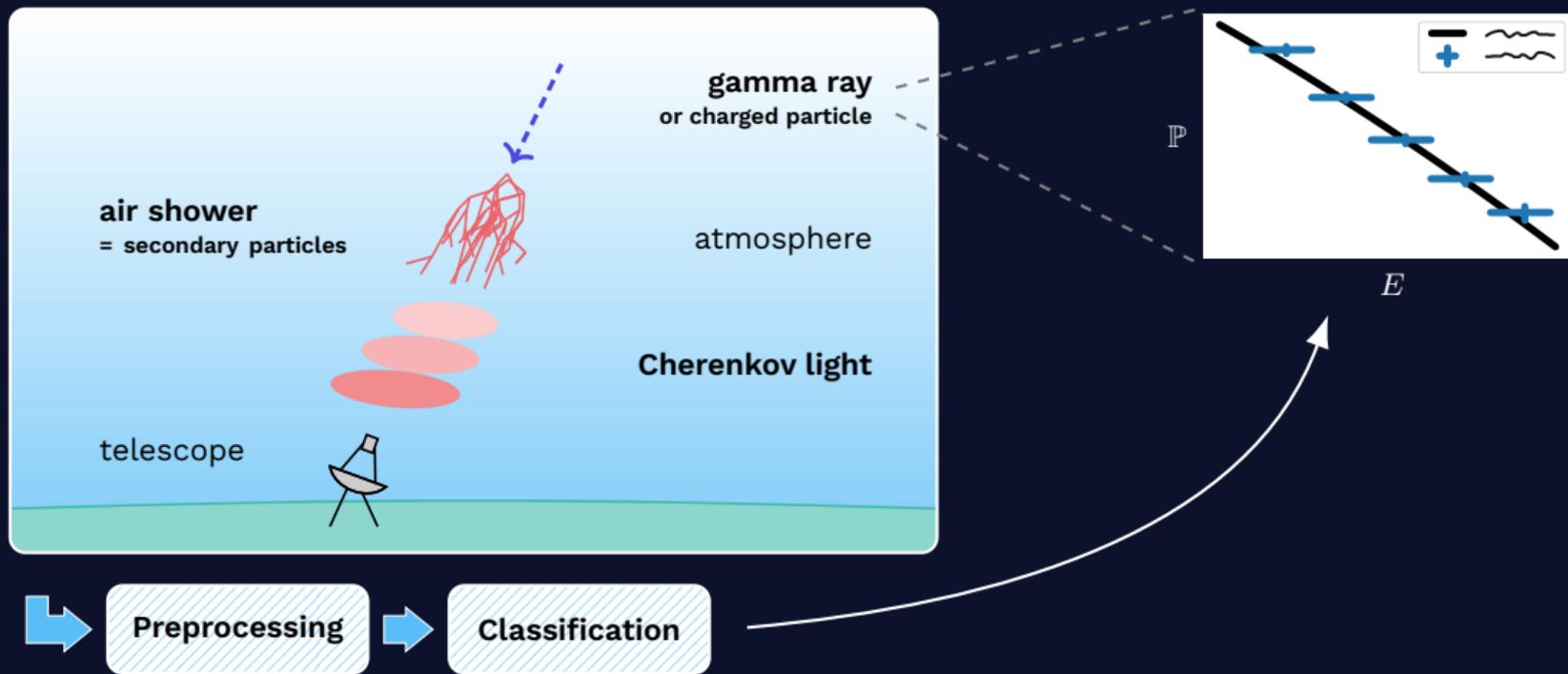


**Preprocessing**

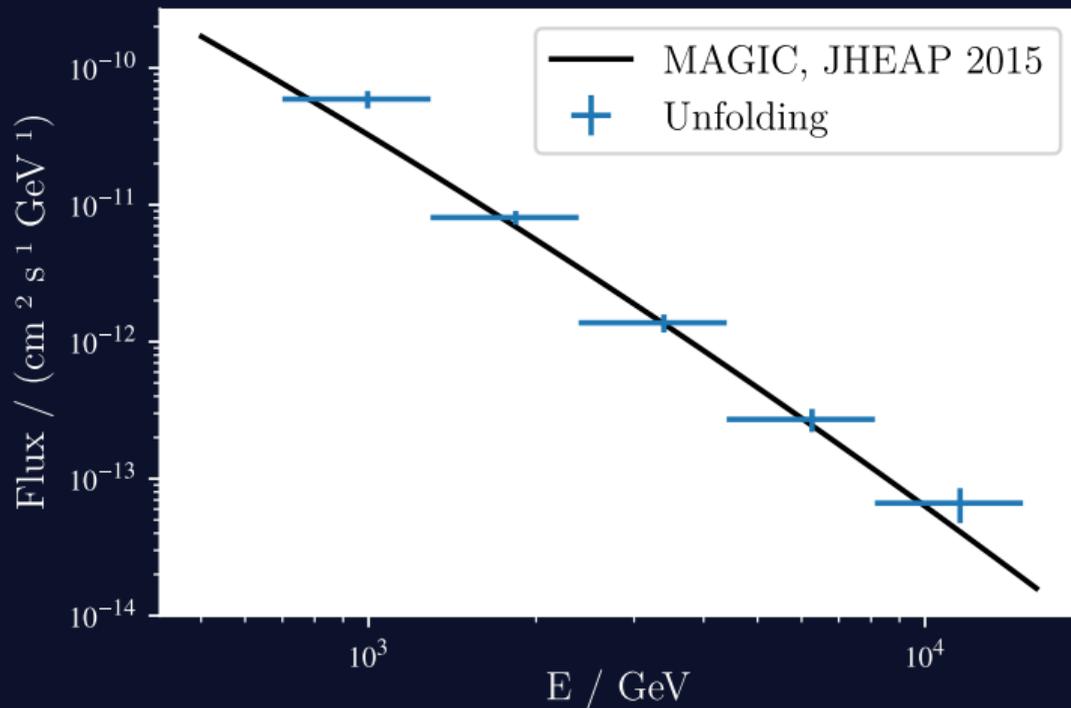


**Classification**

# Example: astro-particle physics



## Example: astro-particle physics



# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \quad \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)
- background:  $Q(\mathbf{x}) = Q(\mathbf{x}, \emptyset) + \sum_{y=1}^C Q(\mathbf{x}, y)$  (PPS with a noise class<sup>29</sup>)

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)
- background:  $Q(\mathbf{x}) = Q(\mathbf{x}, \emptyset) + \sum_{y=1}^C Q(\mathbf{x}, y)$  (PPS with a noise class<sup>29</sup>)
- class-conditional selection bias:  $Q(\mathbf{x} \in B \mid y_i) \neq Q(\mathbf{x} \in B \mid y_j) \exists i \neq j$

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)
- background:  $Q(\mathbf{x}) = Q(\mathbf{x}, \emptyset) + \sum_{y=1}^C Q(\mathbf{x}, y)$  (PPS with a noise class<sup>29</sup>)
- class-conditional selection bias:  $Q(\mathbf{x} \in B | y_i) \neq Q(\mathbf{x} \in B | y_j) \exists i \neq j$
- changing environment:  $Q(\mathbf{x}, y) = \sum_{e \in \mathcal{E}} Q(\mathbf{x}, y, e)$

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)
- background:  $Q(\mathbf{x}) = Q(\mathbf{x}, \emptyset) + \sum_{y=1}^C Q(\mathbf{x}, y)$  (PPS with a noise class<sup>29</sup>)
- class-conditional selection bias:  $Q(\mathbf{x} \in B | y_i) \neq Q(\mathbf{x} \in B | y_j) \exists i \neq j$
- changing environment:  $Q(\mathbf{x}, y) = \sum_{e \in \mathcal{E}} Q(\mathbf{x}, y, e)$
- concept shift:  $Q(\mathbf{x} | y) \neq \mathbb{P}(\mathbf{x} | y)$  (in addition to PPS)

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.

# Advanced topics for experimental physics



## Issues to be resolved:

- ordinality:  $y_i \prec y_{i+1} \forall i \in \mathcal{Y}$  (to be covered through regularization for ordinal plausibility<sup>28</sup>)
- background:  $Q(\mathbf{x}) = Q(\mathbf{x}, \emptyset) + \sum_{y=1}^C Q(\mathbf{x}, y)$  (PPS with a noise class<sup>29</sup>)
- class-conditional selection bias:  $Q(\mathbf{x} \in B | y_i) \neq Q(\mathbf{x} \in B | y_j) \exists i \neq j$
- changing environment:  $Q(\mathbf{x}, y) = \sum_{e \in \mathcal{E}} Q(\mathbf{x}, y, e)$
- concept shift:  $Q(\mathbf{x} | y) \neq \mathbb{P}(\mathbf{x} | y)$  (in addition to PPS)
- inspect contributions of individual data items  $\mathbf{x} \in B$  to  $\lambda(B)$  (data selection, human in the loop)

<sup>28</sup> Bunse et al., “Regularization-based Methods for Ordinal Quantification”, 2024.

<sup>29</sup> Dussap, Blanchard, and Chérif-Abdellatif, “Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching”, 2023.