



Memorial Sloan Kettering  
Cancer Center

# IMPACT annotator

November 5, 2018

Pierre Guilmin

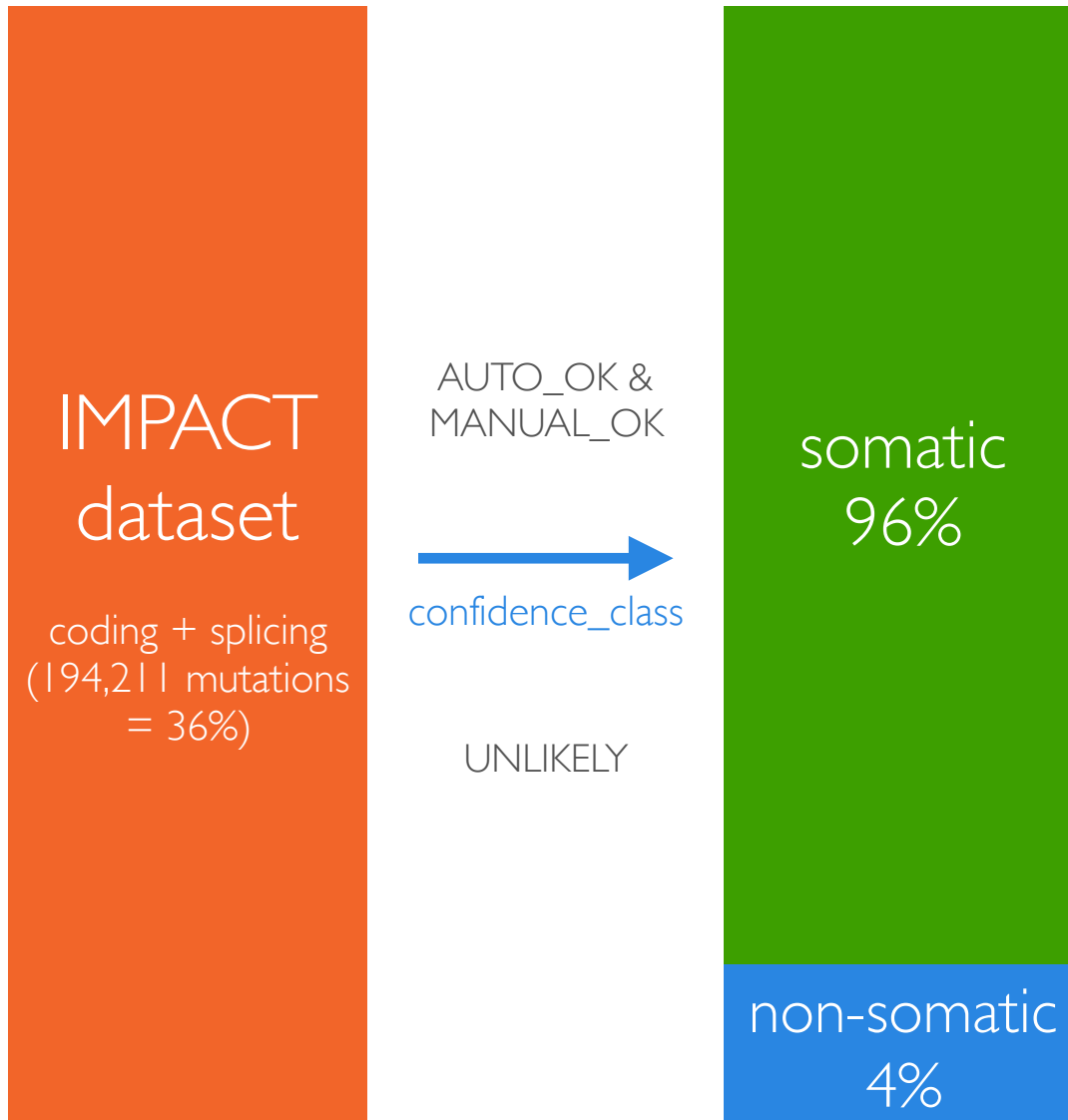
Elsa Bernard



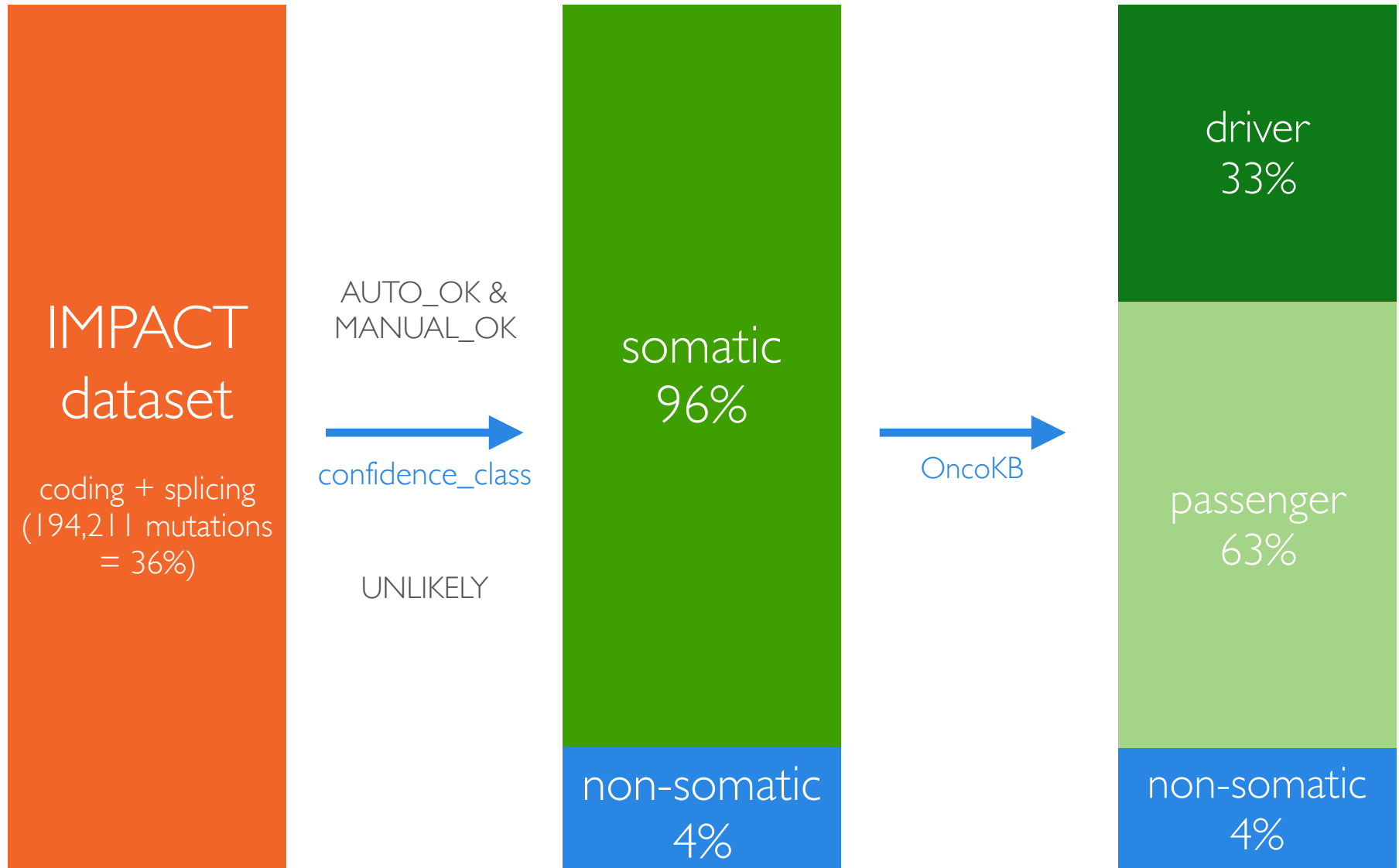
## IMPACT dataset

coding + splicing  
(194,211 mutations  
= 36%)

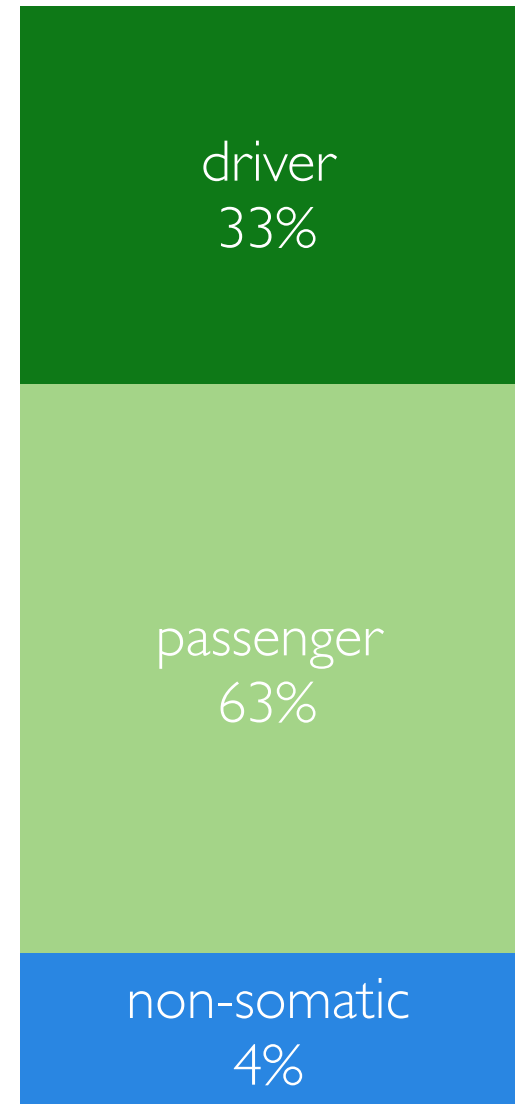
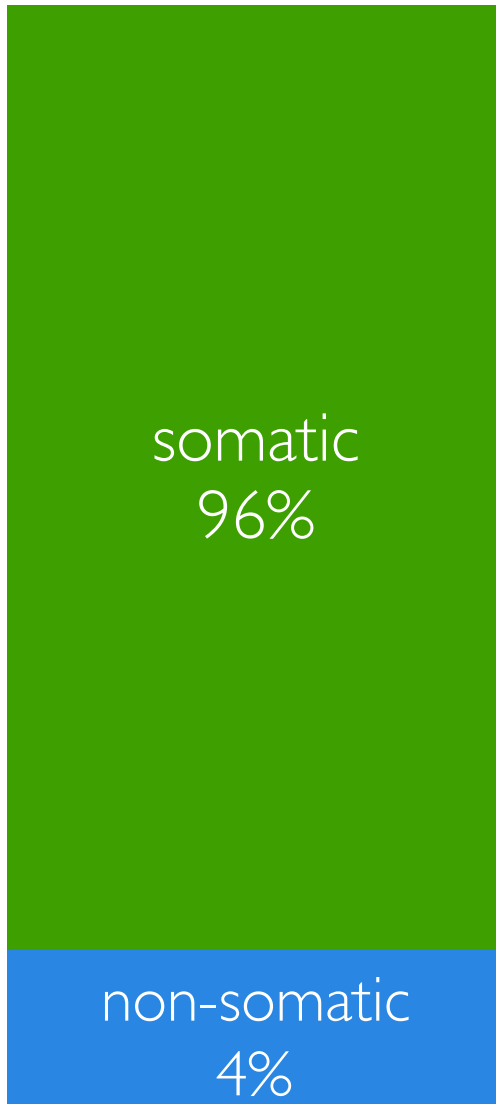
# Different classes



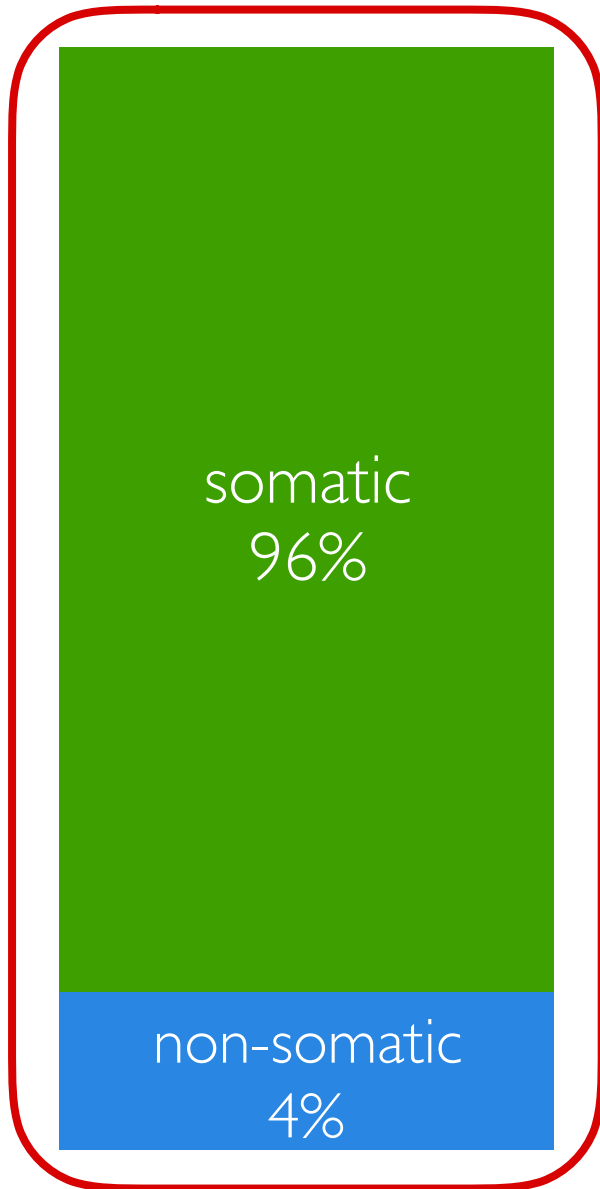
# Different classes



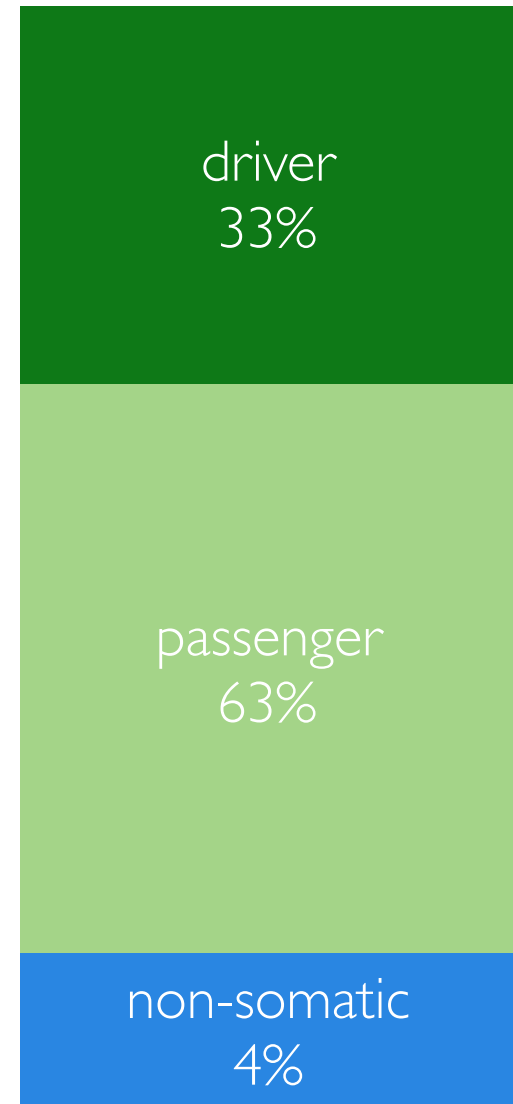
# Supervised learning



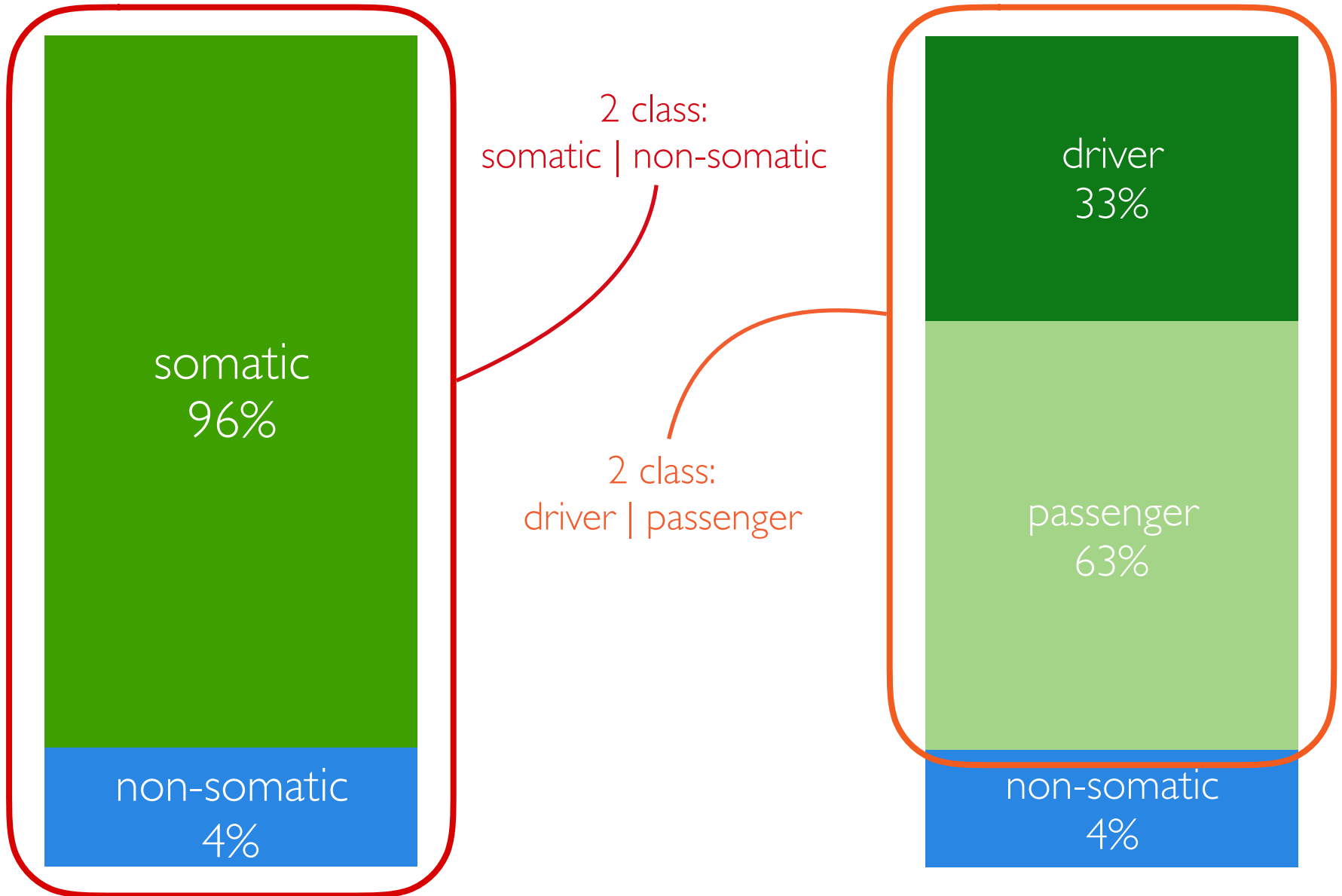
# Supervised learning



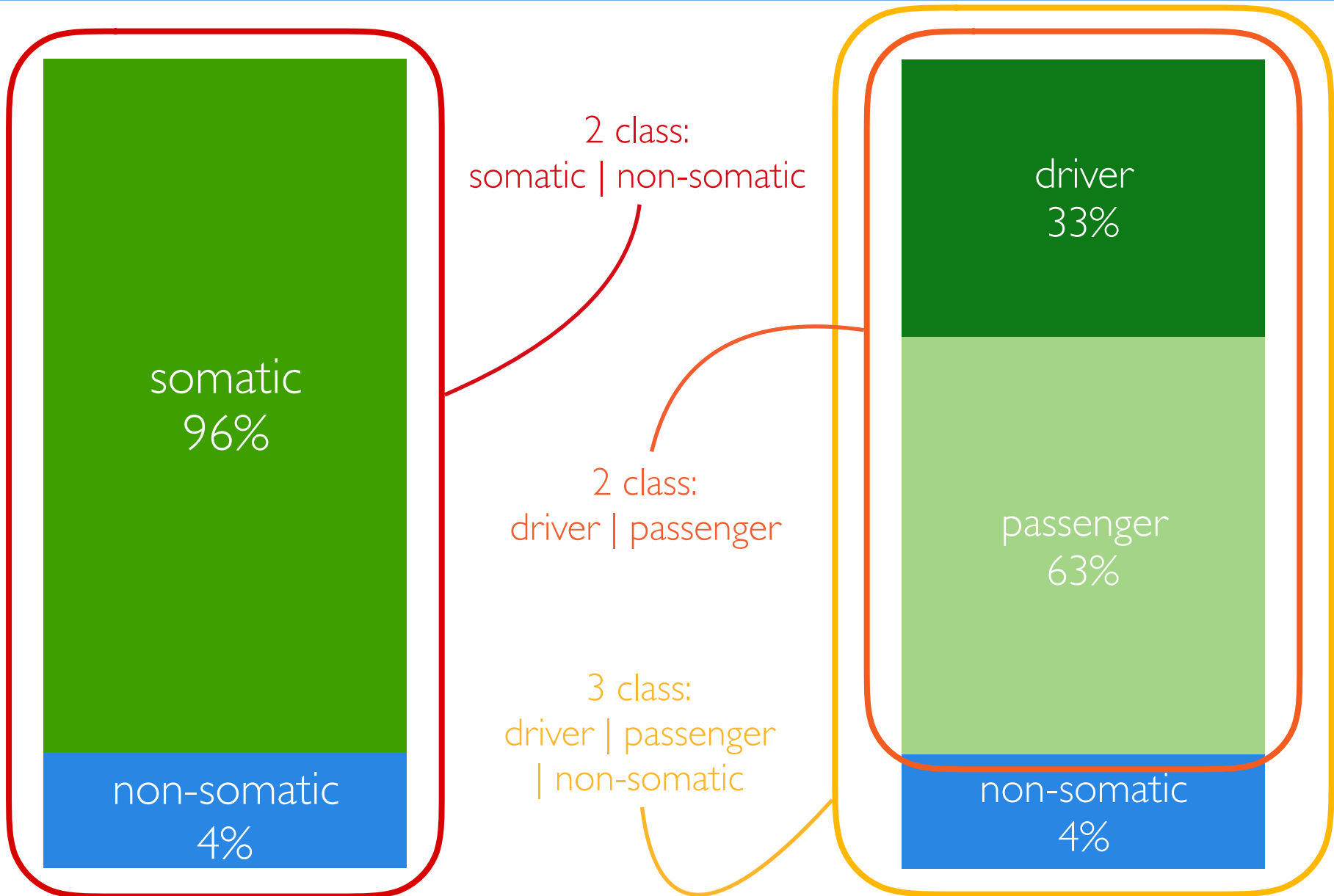
2 class:  
somatic | non-somatic



# Supervised learning

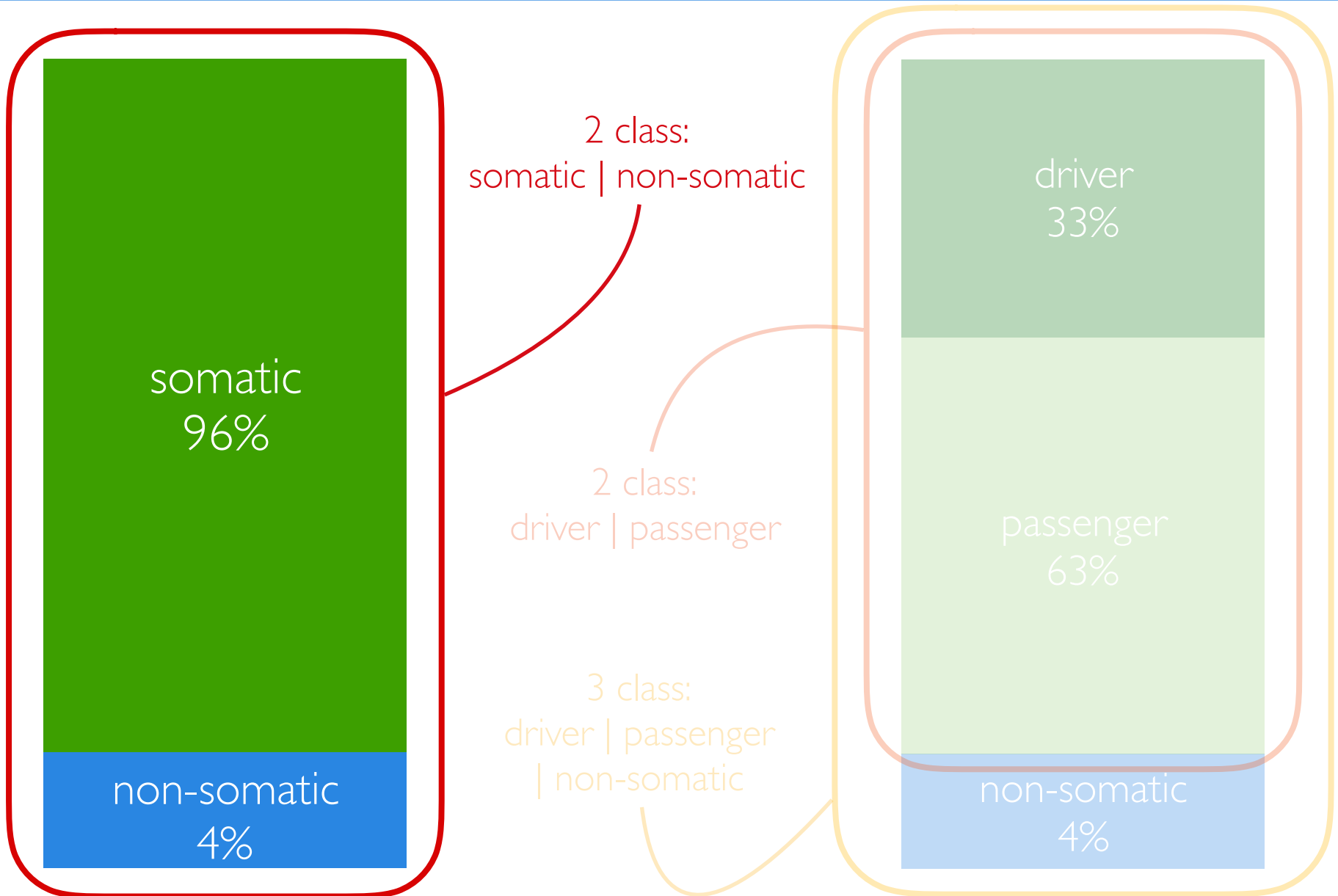


# Supervised learning

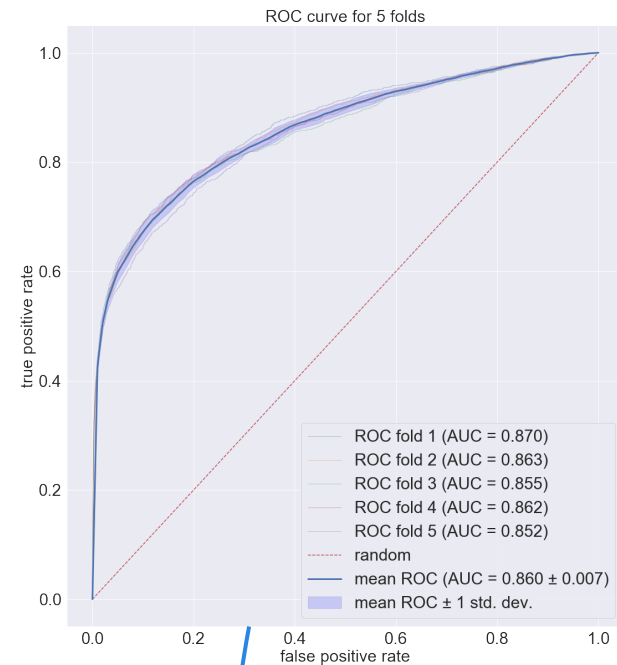




# Supervised learning

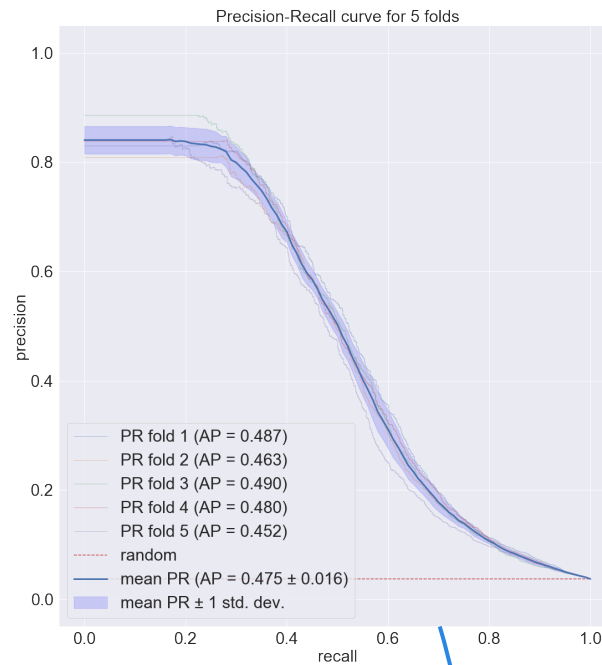


# Main metrics used



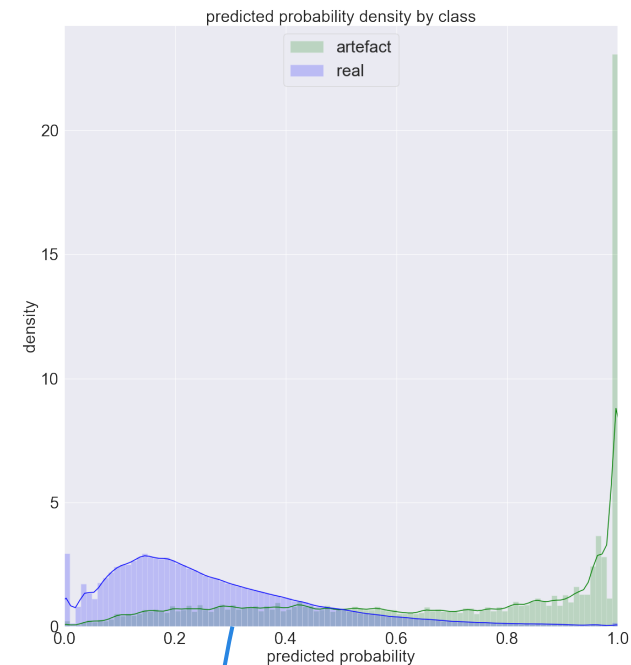
ROC AUC

area under the  
ROC curve



average precision

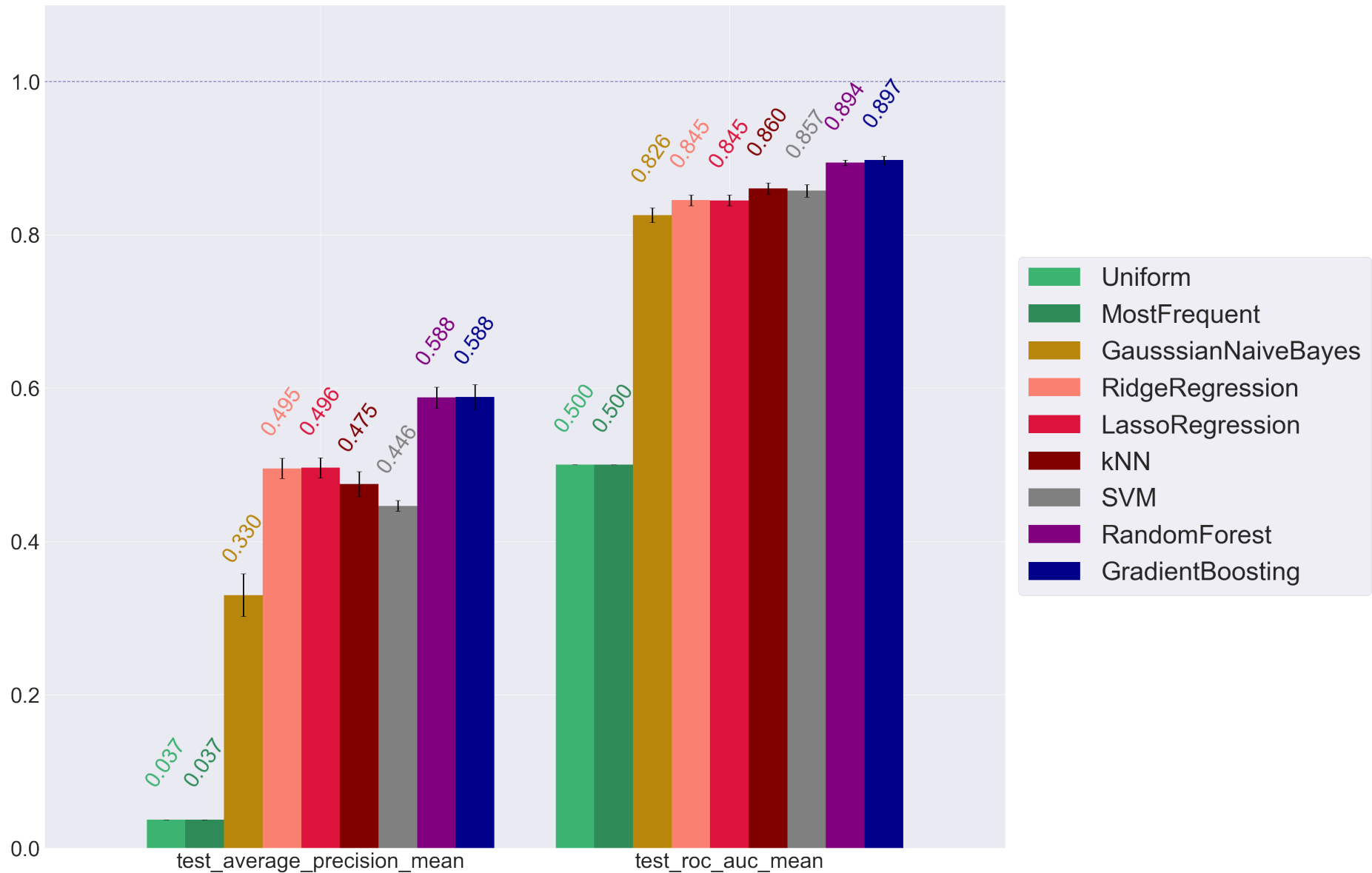
area under the precision-  
recall curve



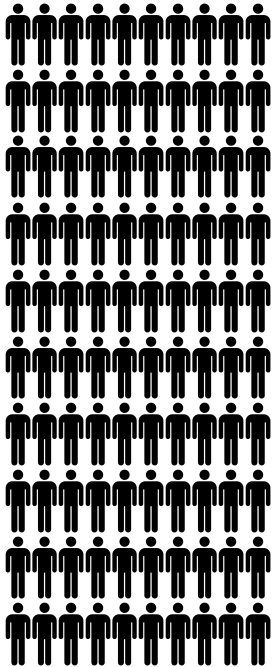
Probability distribution

probability predicted by the  
classifier for each class

# Algorithm comparison



100 patients



# IMPACT annotator

100 patients

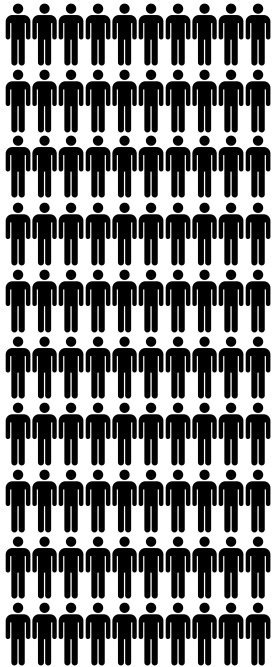
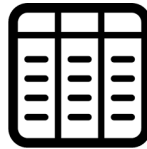
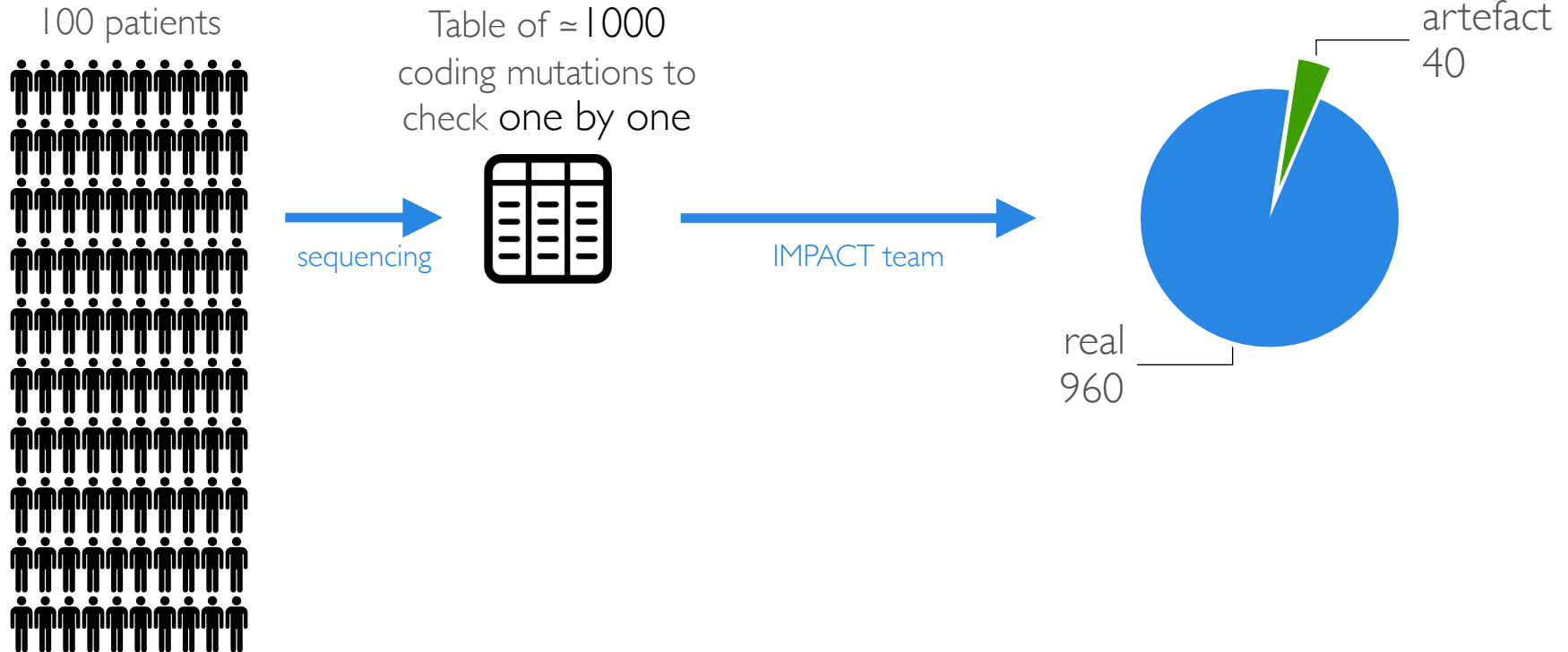


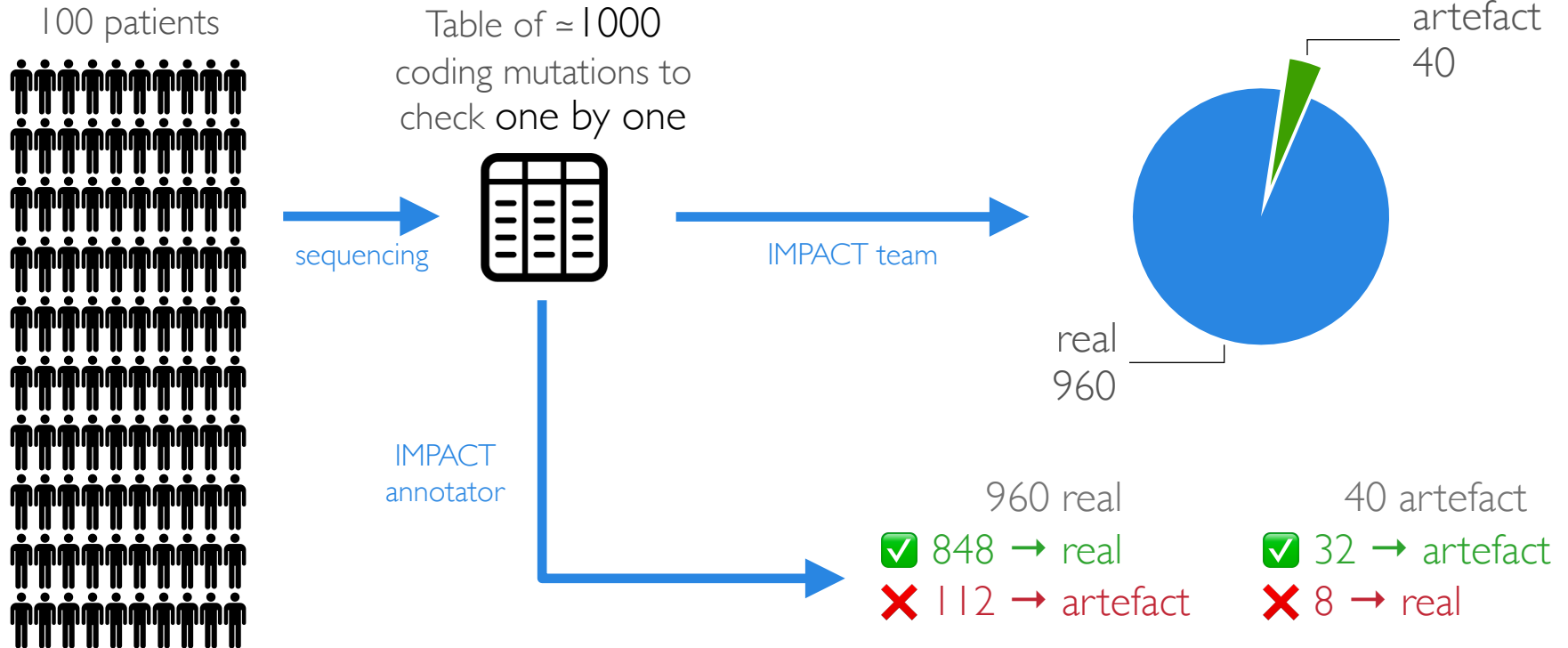
Table of  $\approx 1000$   
coding mutations to  
check one by one



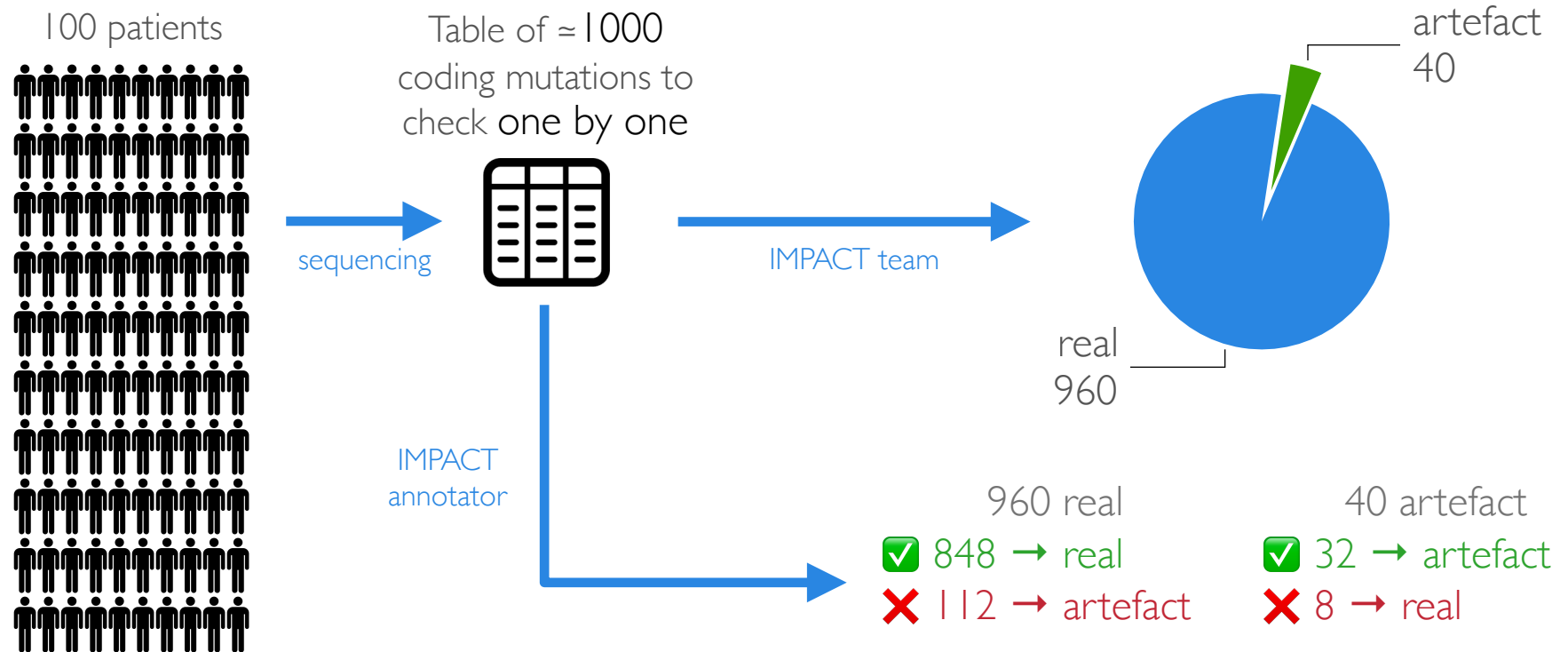
# IMPACT annotator



# IMPACT annotator



# IMPACT annotator



## IMPACT team



1000/1000 mutations to check one by one

✗ 0/40 artefacts considered as real

## IMPACT annotator

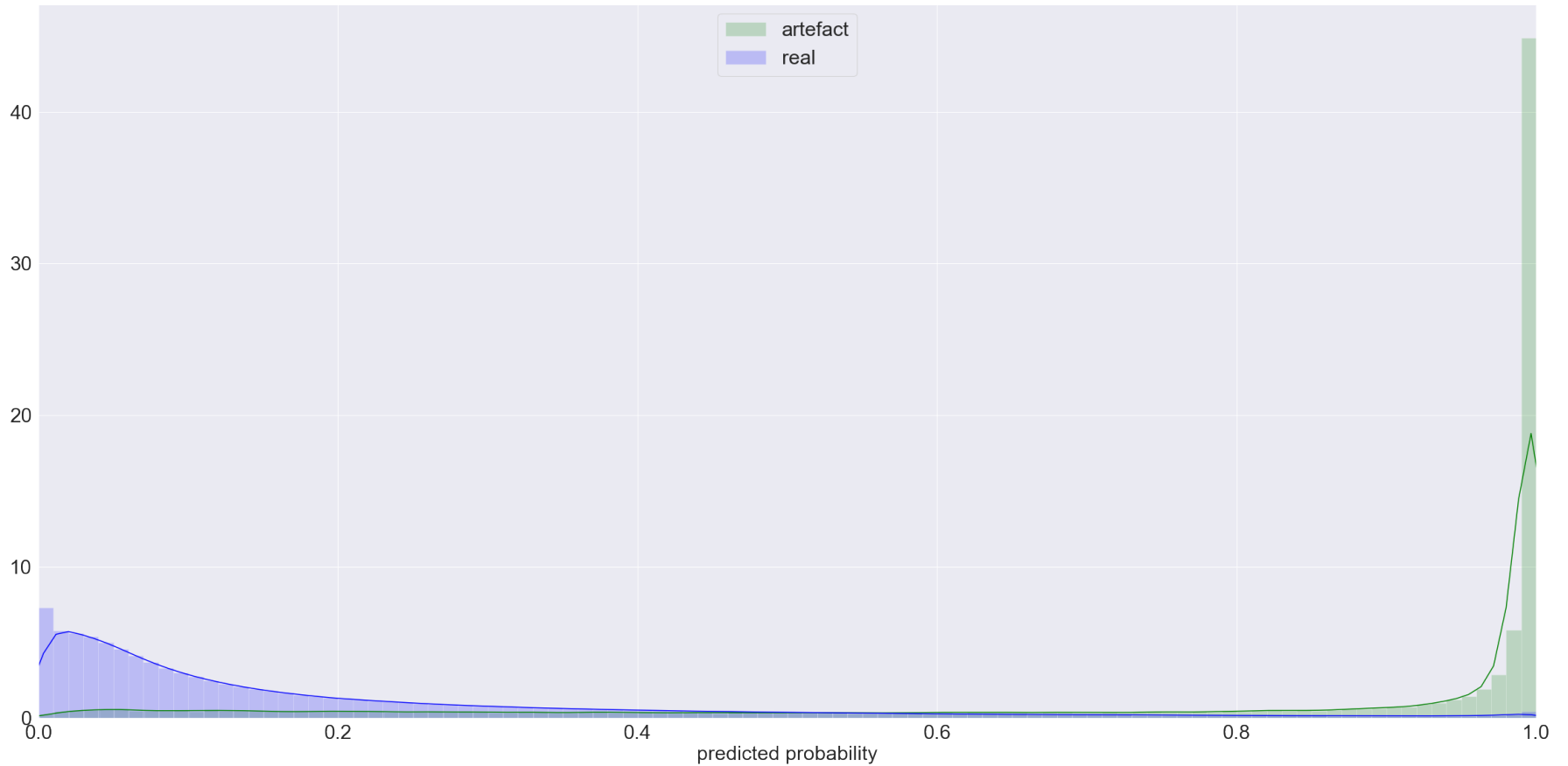


144/1000 mutations to check one by one

✗ 8/40 artefacts considered as real



# Threshold choice



# Threshold choice

threshold = 0.5

 144 mutations to check

960 real

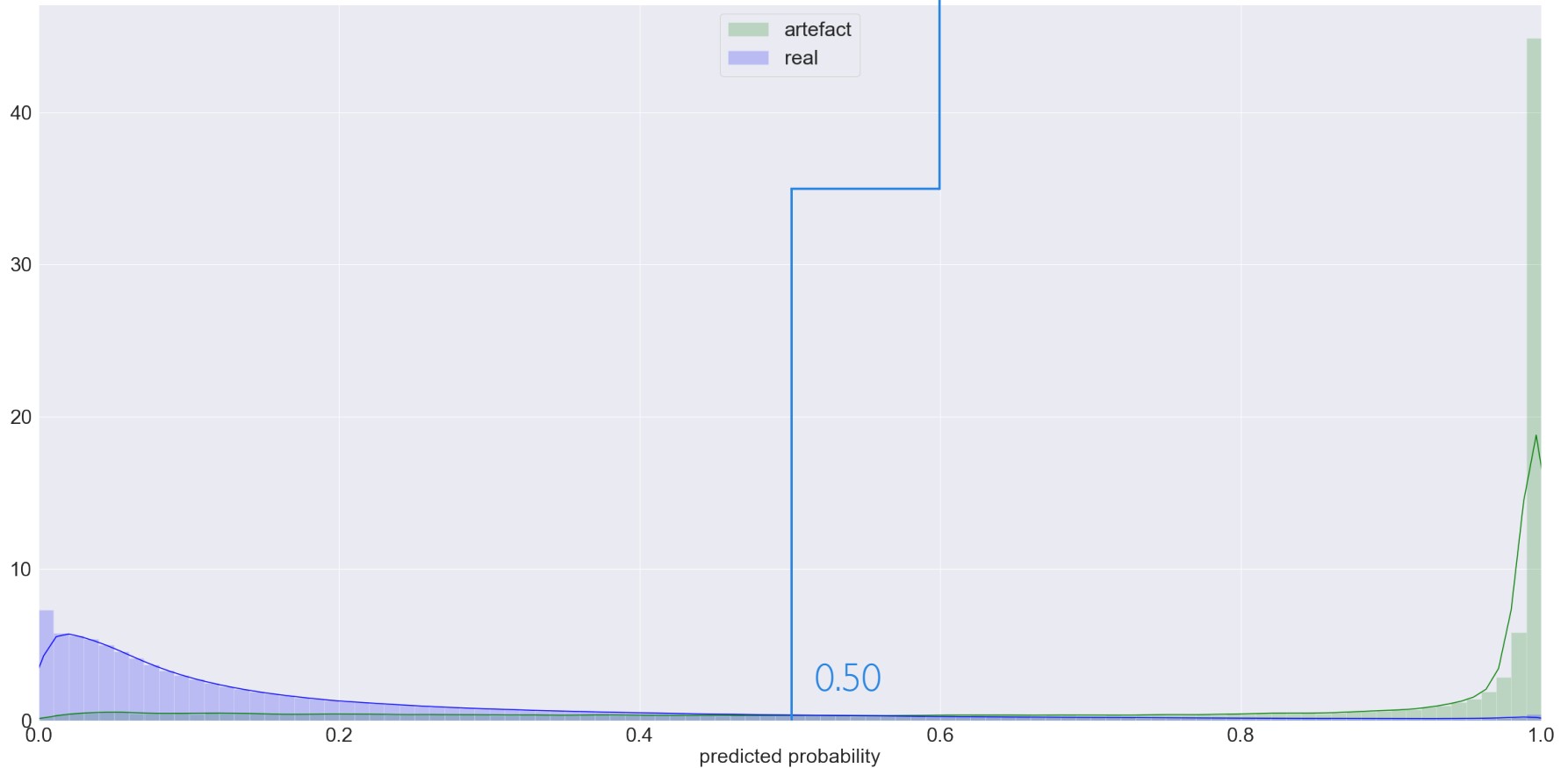
40 artefact

✓ 848 → real

✓ 32 → artefact

✗ 112 → artefact

✗ 8 → real



# Threshold choice

threshold = 0.25

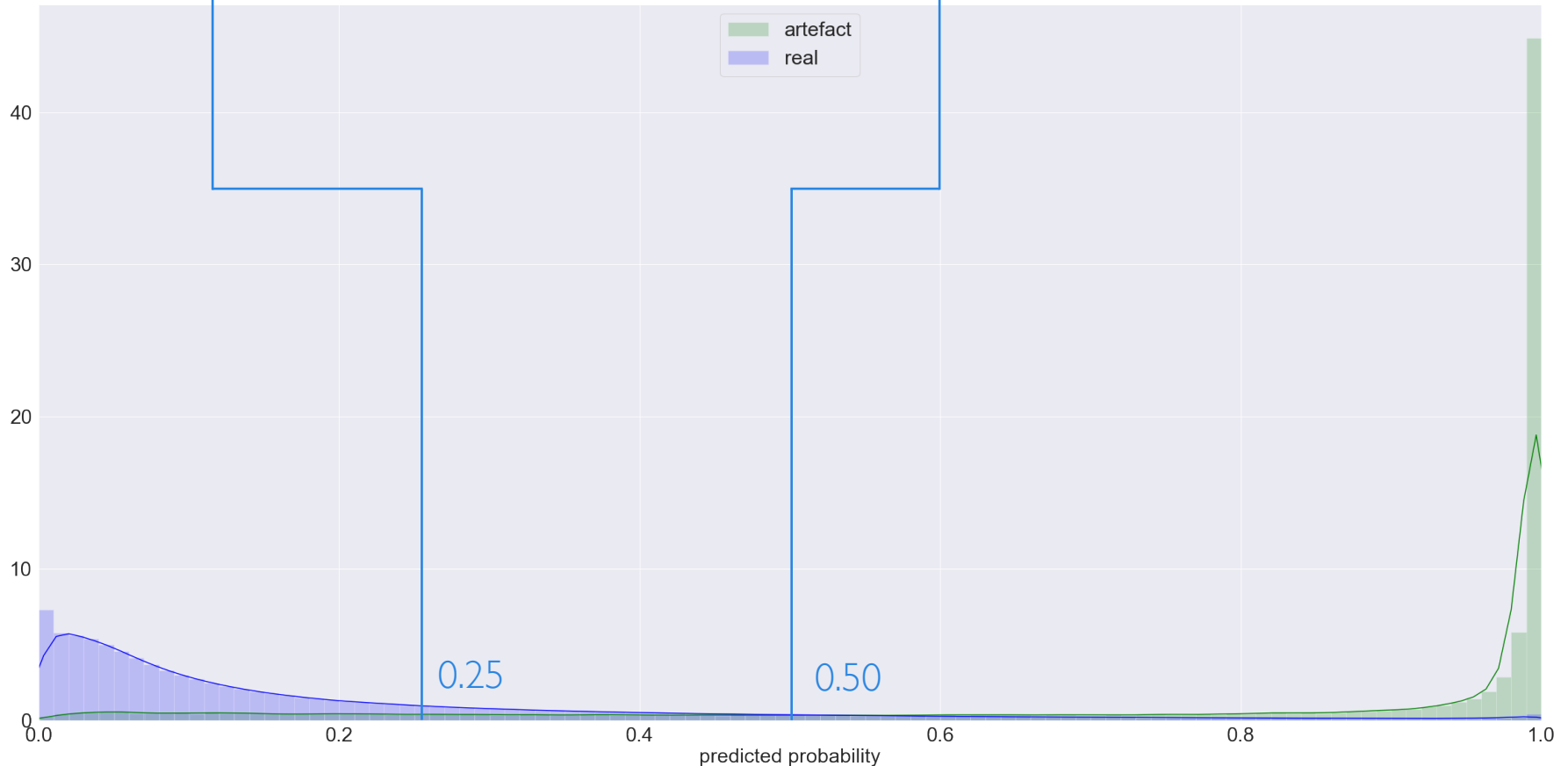
 296

960 real	40 artefact
✓ 699	✓ 35
✗ 261	✗ 5

threshold = 0.5

 144 mutations to check

960 real	40 artefact
✓ 848 → real	✓ 32 → artefact
✗ 112 → artefact	✗ 8 → real



# Threshold choice

threshold = 0.25

296

960 real	40 artefact
✓ 699	✓ 35
✗ 261	✗ 5

threshold = 0.5

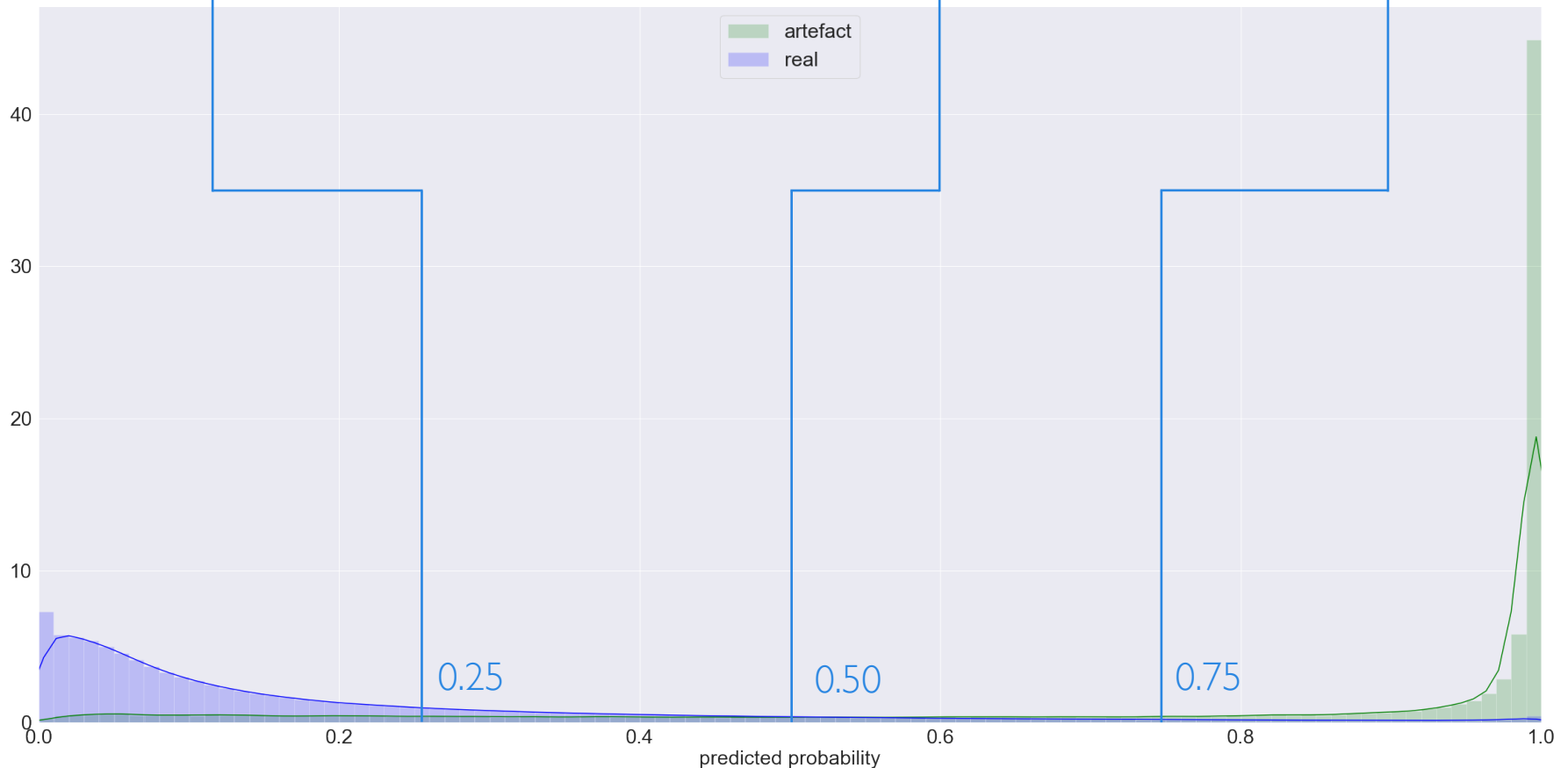
144 mutations to check

960 real	40 artefact
✓ 848 → real	✓ 32 → artefact
✗ 112 → artefact	✗ 8 → real

threshold = 0.75

72

960 real	40 artefact
✓ 916	✓ 28
✗ 44	✗ 12



# The features group



## NGS (x 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

# The features group

## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

## genome (× 3)

- Chromosome
- Hugo Symbol
- Variant Class

# The features group



## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

## genome (× 3)

- Chromosome
- Hugo Symbol
- Variant Class

## somatic driver (× 4)

- COSMIC count
- is a hotspot, is a 3d hotspot
- OncoKB oncogenic

# The features group

## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

## AF (× 12)

- in dbSNP
- gnomAD total AF and AF by population (AFR, AMR, ASJ, ...)

## genome (× 3)

- Chromosome
- Hugo Symbol
- Variant Class

## somatic driver (× 4)

- COSMIC count
- is a hotspot, is a 3d hotspot
- OncoKB oncogenic



# The features group



## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

## AF (× 12)

- in dbSNP
- gnomAD total AF and AF by population (AFR, AMR, ASJ, ...)

## genome (× 3)

- Chromosome
- Hugo Symbol
- Variant Class

## freq (× 1)

- frequency in normals

## somatic driver (× 4)

- COSMIC count
- is a hotspot, is a 3d hotspot
- OncoKB oncogenic

# The features group

## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor +/- strand count
- sample\_coverage

## AF (× 12)

- in dbSNP
- gnomAD total AF and AF by population (AFR, AMR, ASJ, ...)

## genome (× 3)

- Chromosome
- Hugo Symbol
- Variant Class

## freq (× 1)

- frequency in normals

## somatic driver (× 4)

- COSMIC count
- is a hotspot, is a 3d hotspot
- OncoKB oncogenic

## Consequence (× 6)

- gene type
- mutation consequence (stopgain, frameshift, ...)
- VEP\_IMPACT, VEP\_CLIN\_SIG
- SIFT and PolyPhen class

# The features group

## AF (× 12)

- in dbSNP
- gnomAD total AF and AF by population (AFR, AMR, ASJ, ...)

## NGS (× 11)

- tumor & normal depth, vaf and count
- tumor positive and negative count
- sample\_coverage

## freq (× 1)

- frequency in normals

## genome (× 3)

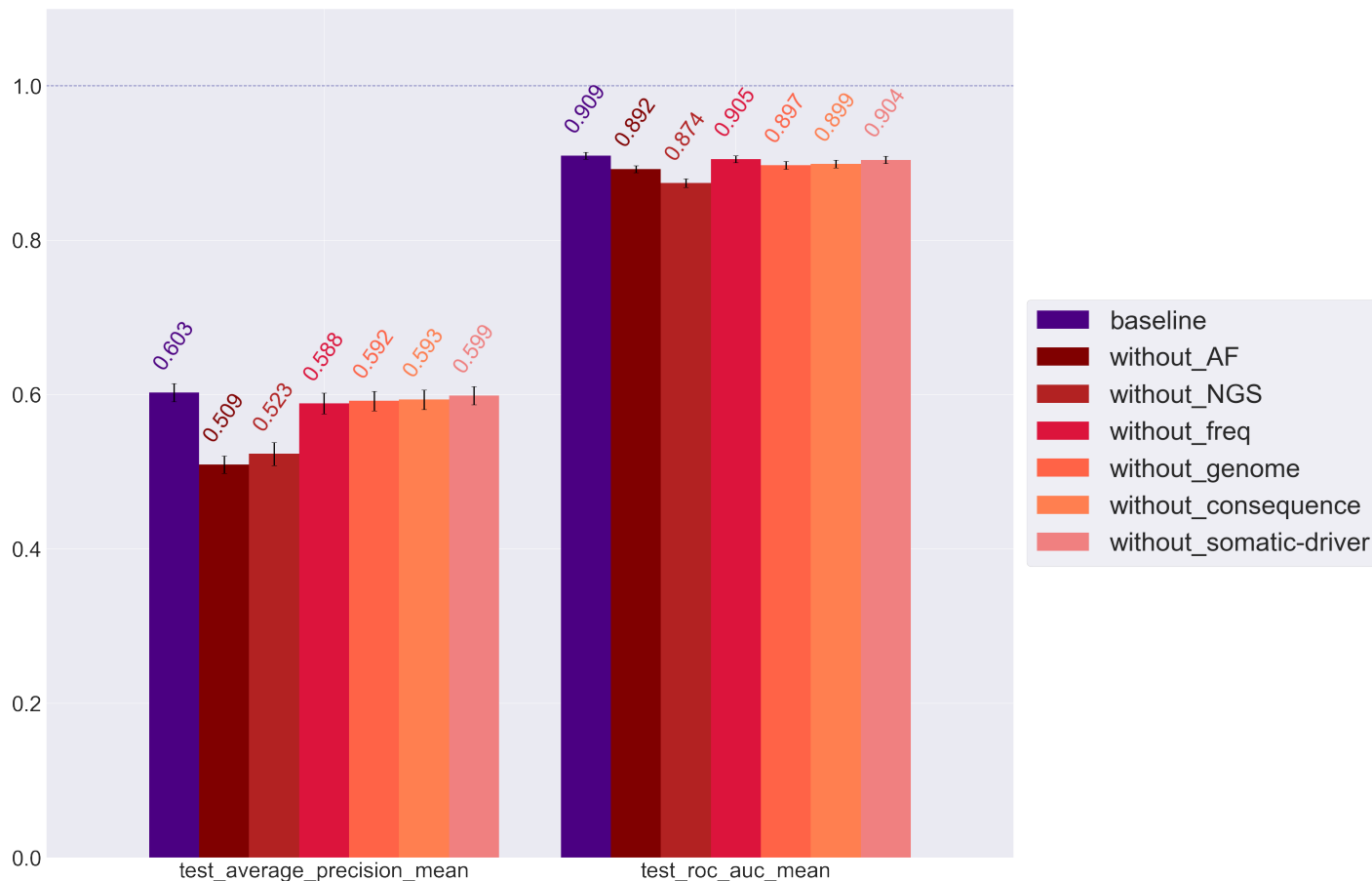
- Chromosome
- Hugo Symbol
- Variant Class

## Consequence (× 6)

- gene type
- Mutation consequence (stopgain, frameshift, ...)
- VEP\_IMPACT, VEP\_CLIN\_SIG
- SIFT and PolyPhen class

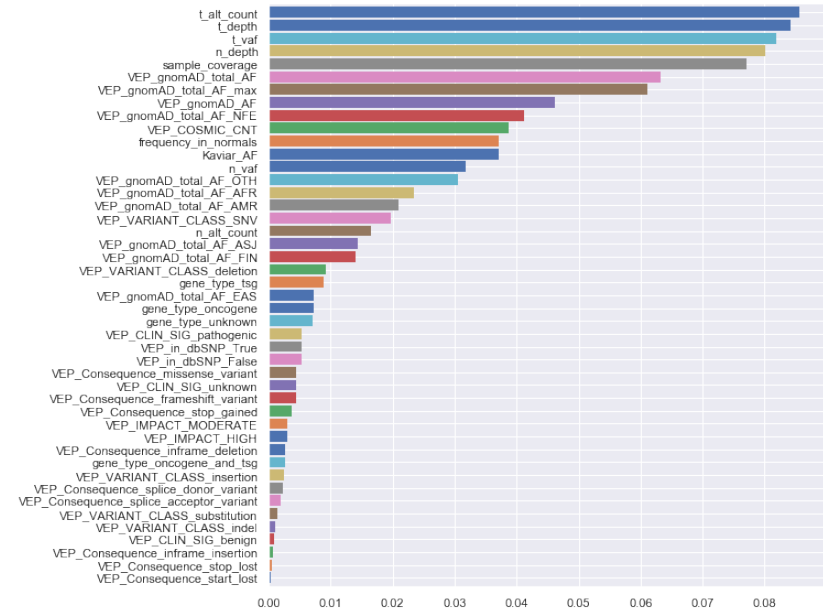
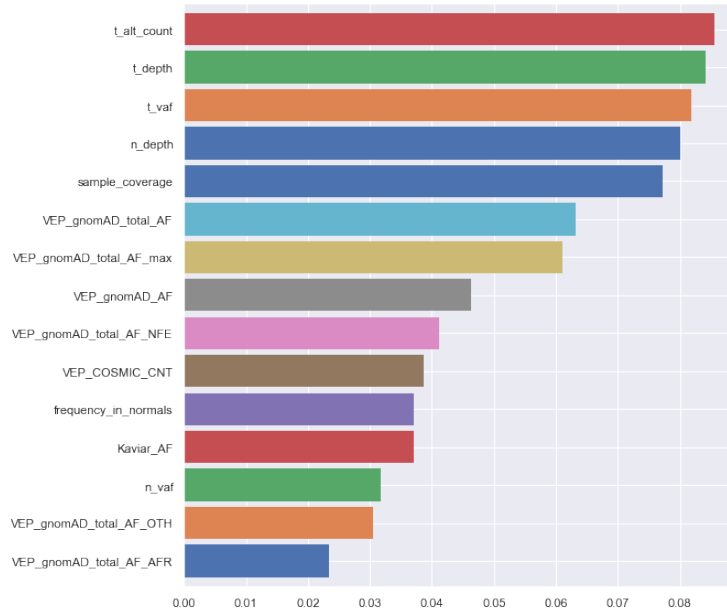
## somatic driver (× 4)

- COSMIC count
- is a hotspot, is a 3d hotspot
- OncoKB oncogenic



# Individual features importance

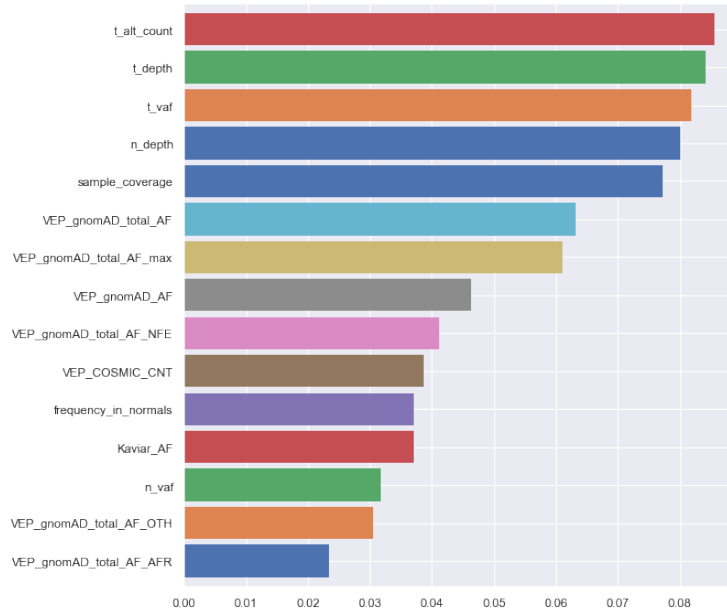
Random  
Forest  
(1000 trees)



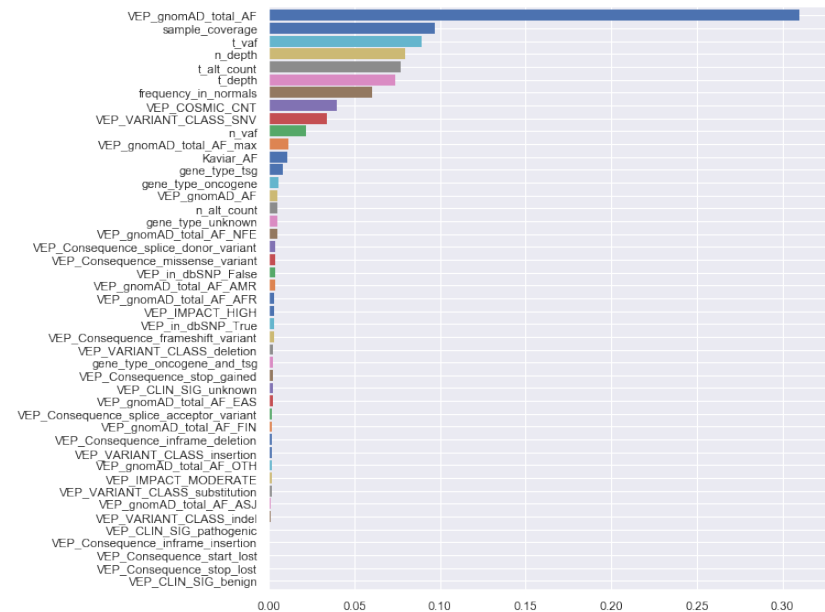
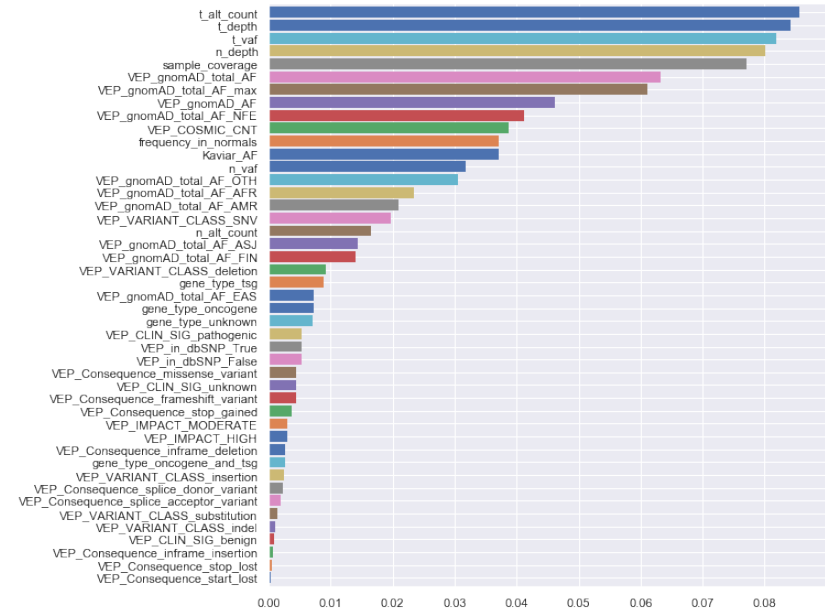
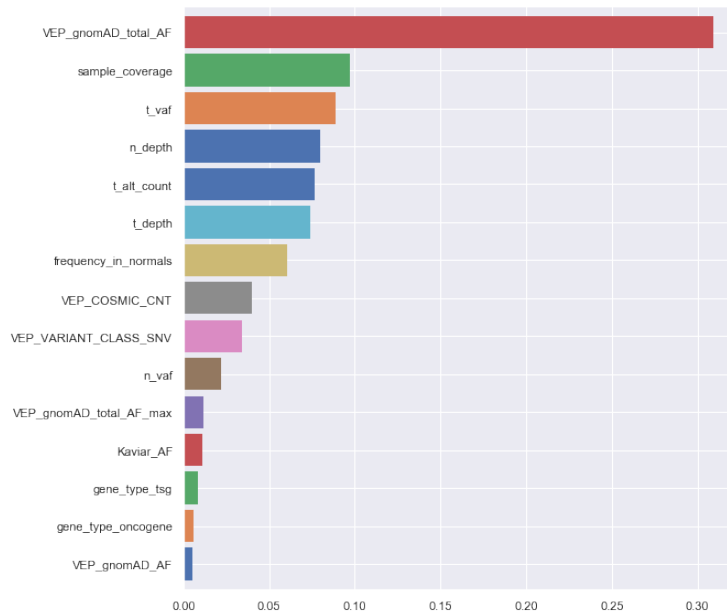
# Individual features importance



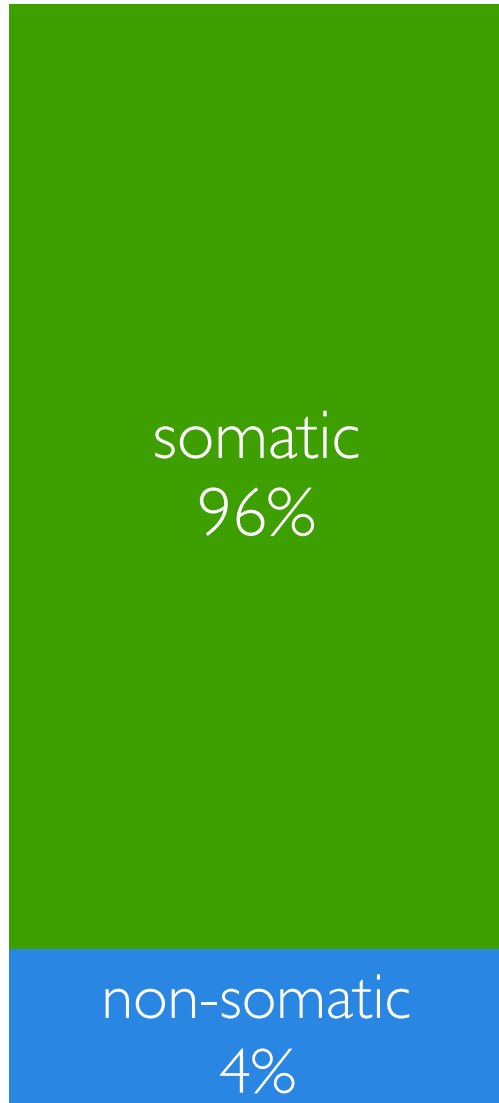
Random  
Forest  
(1000 trees)



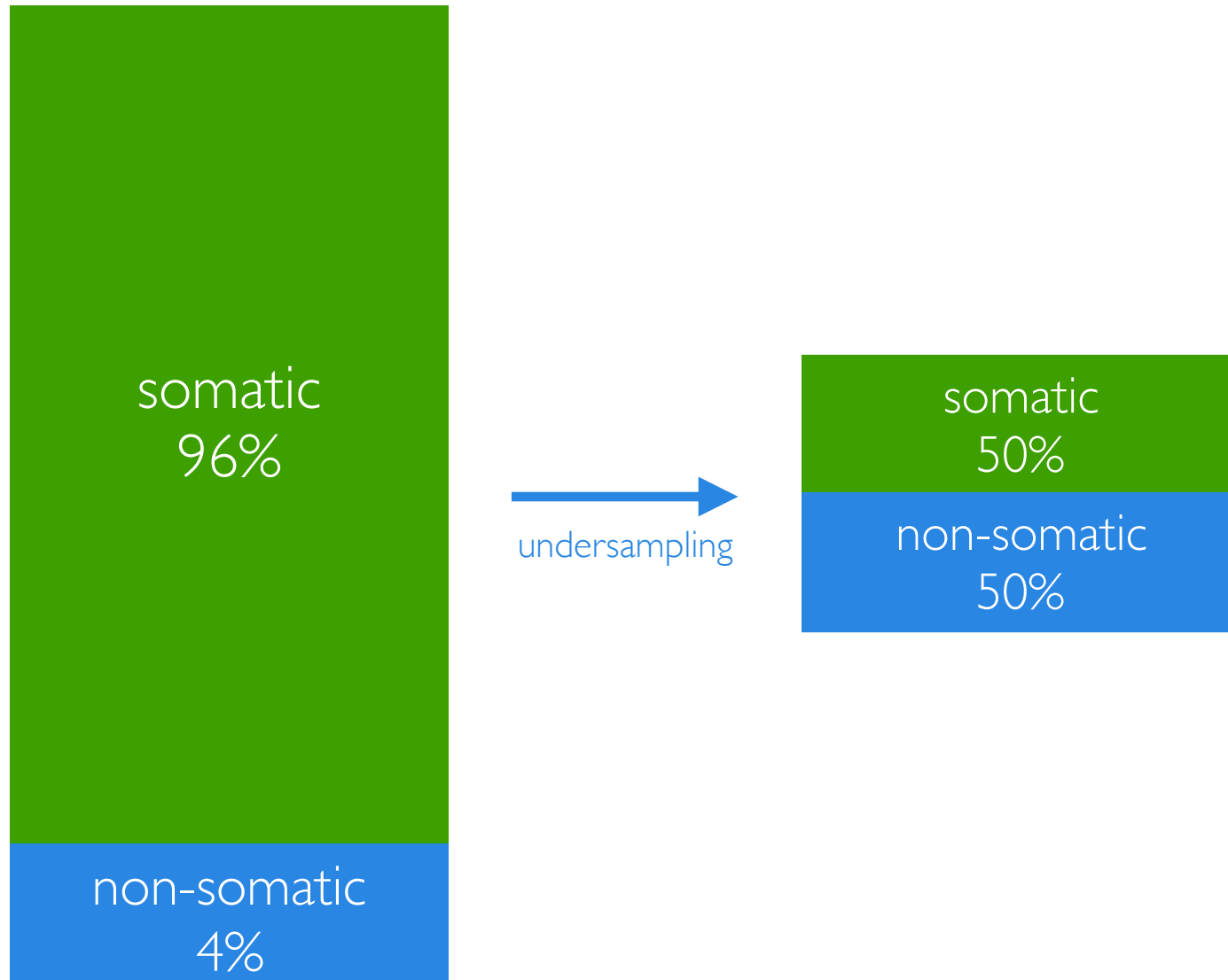
Gradient  
Boosting  
(1000 trees)



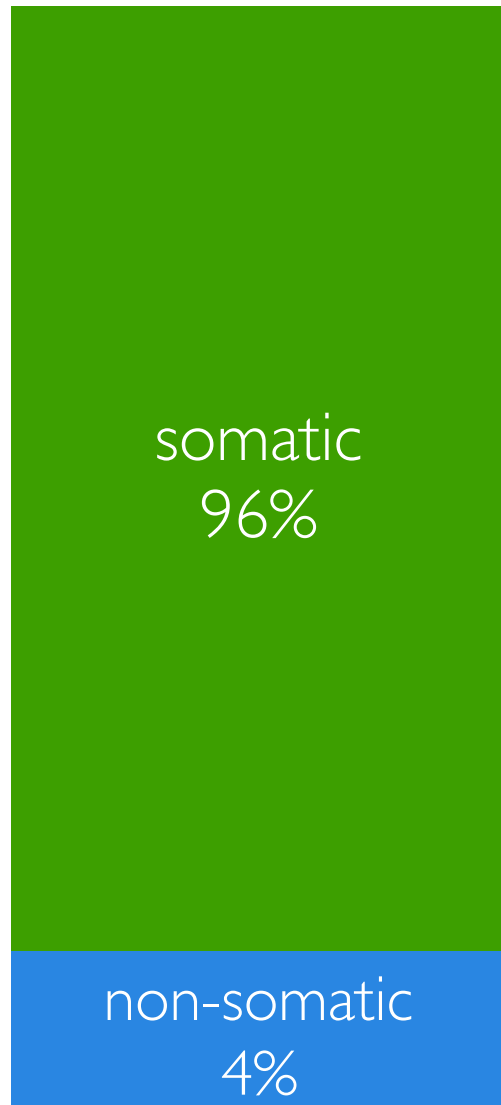
# Imbalanced learning



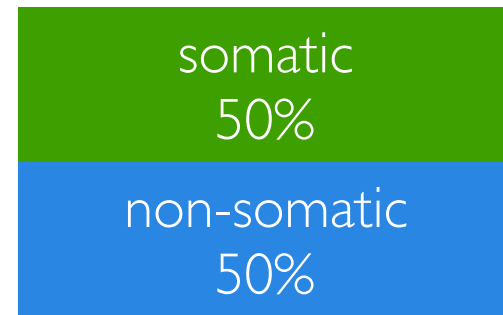
# Imbalanced learning



# Imbalanced learning



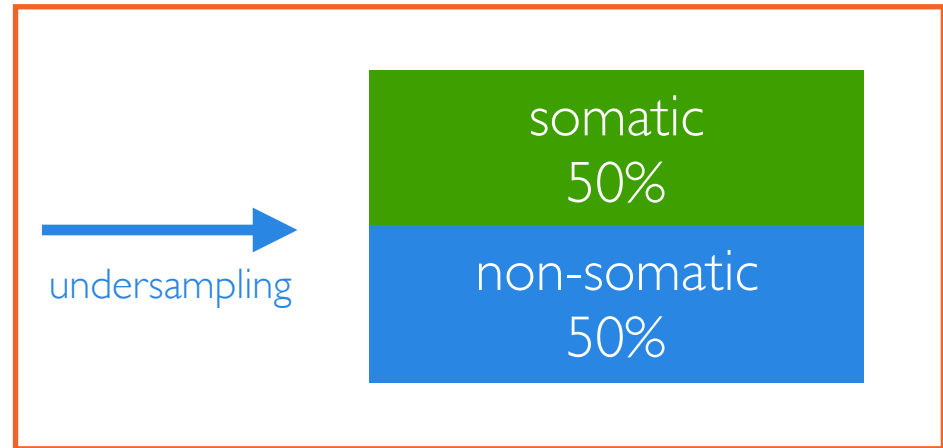
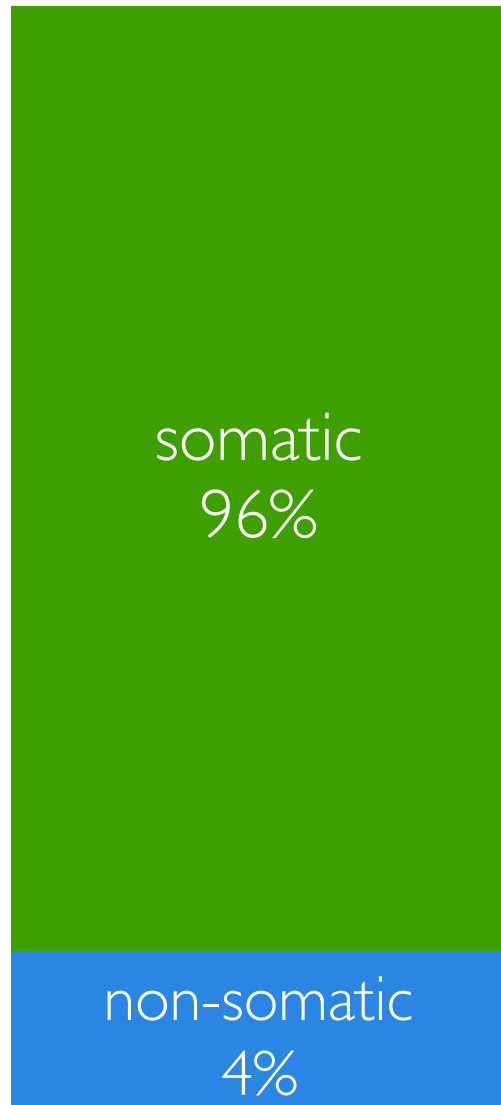
→  
undersampling



within the model



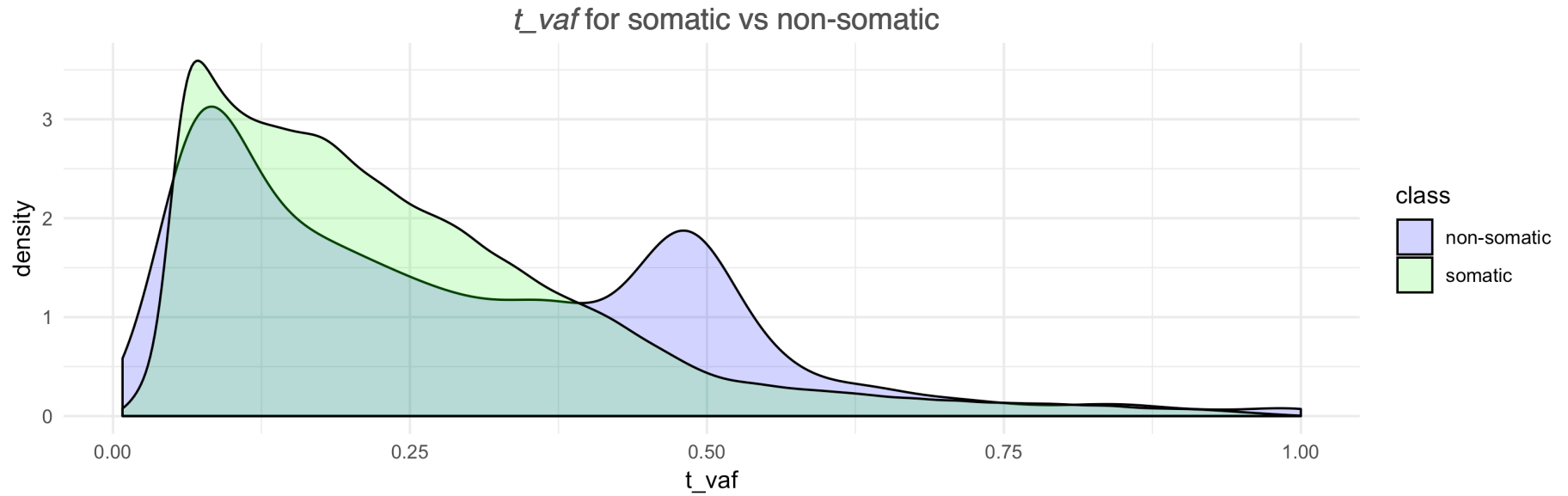
# Imbalanced learning



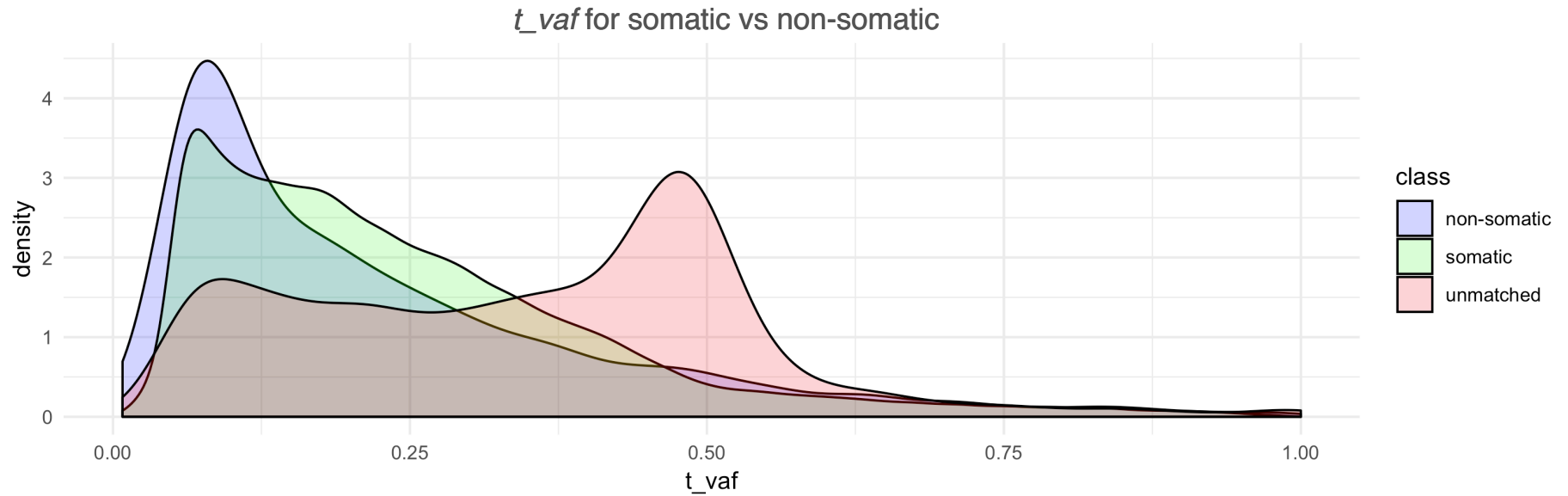
within the model

So the dataset used to train the model is only 8% of the whole dataset (7,199 real and 7,199 artefact mutations)

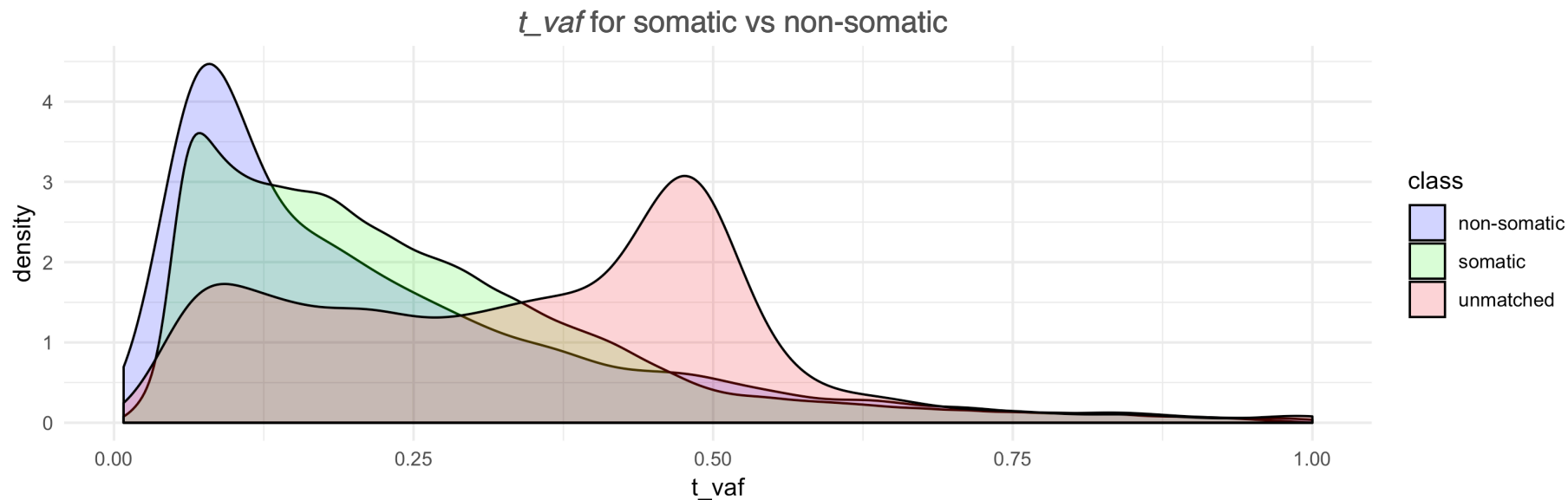
# Removing the unmatched samples



# Removing the unmatched samples



# Removing the unmatched samples



Without the unmatched artefacts (and missing rows):  
4,477 artefacts (instead of 7,199)

# The sorted dataset phenomenon



IMPACT dataset

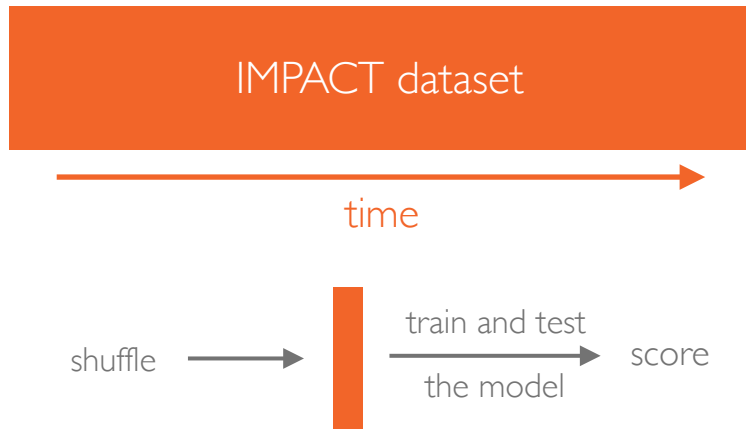


time

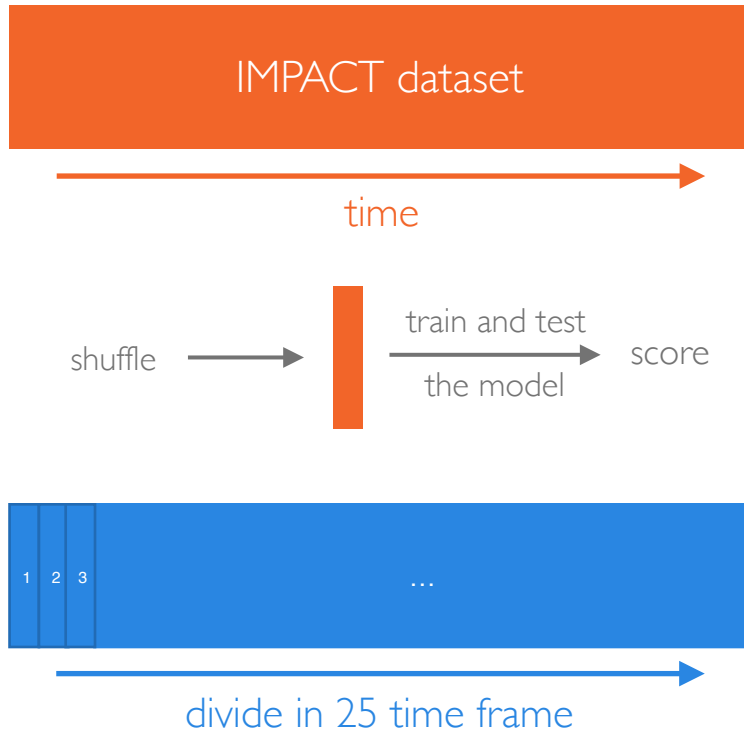
# The sorted dataset phenomenon



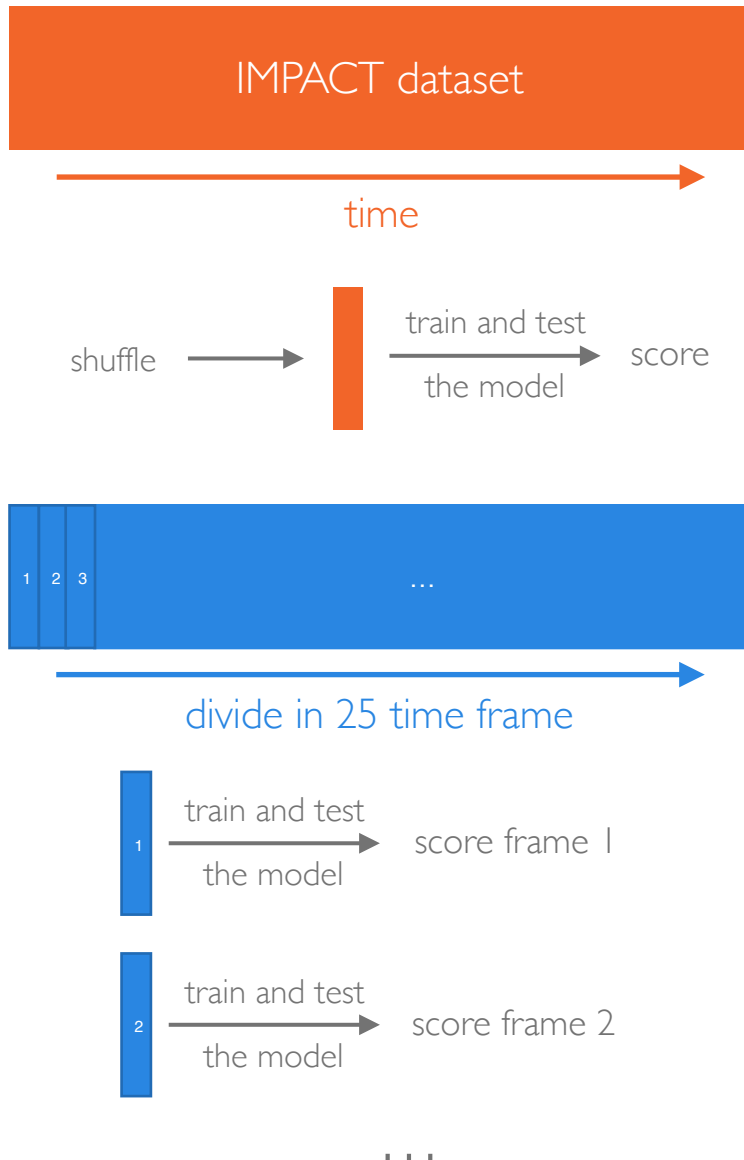
1884



# The sorted dataset phenomenon



# The sorted dataset phenomenon

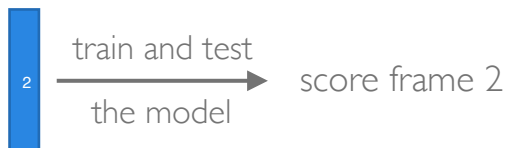
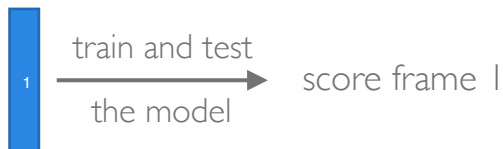




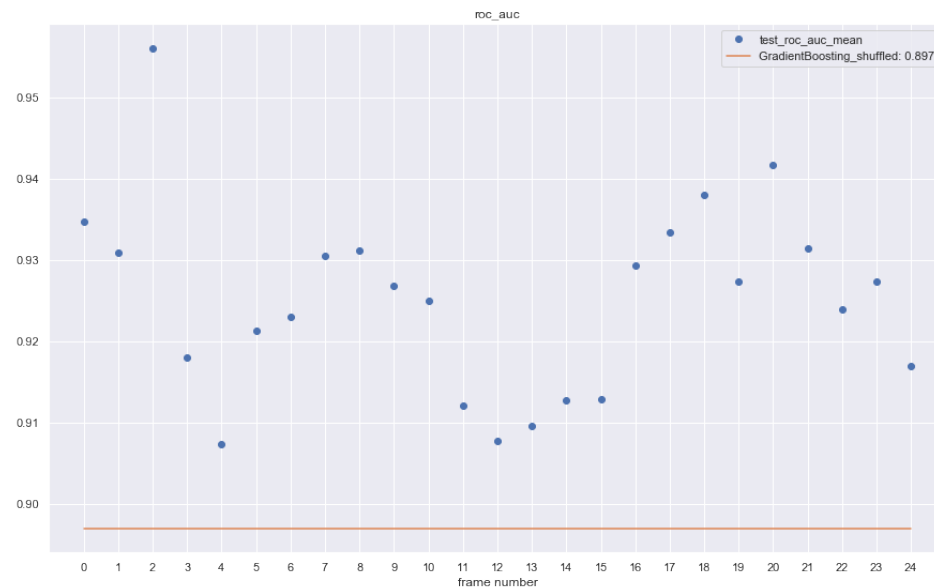
# The sorted dataset phenomenon



divide in 25 time frame



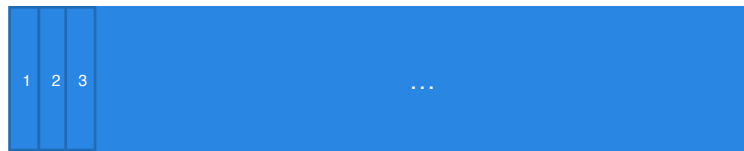
...



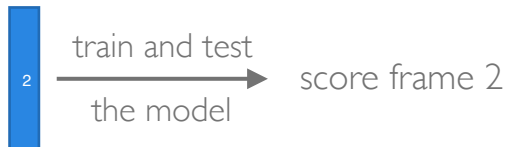
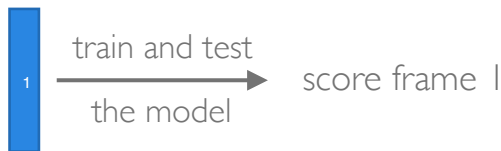
# The sorted dataset phenomenon



time



divide in 25 time frame



...

