



Using Stratified Sampling to Improve LIME Image Explanations

Muhammad Rashid¹, Elvio G. Amparore¹, Enrico Ferrari², Damiano Verda²

Abstract

PROBLEM

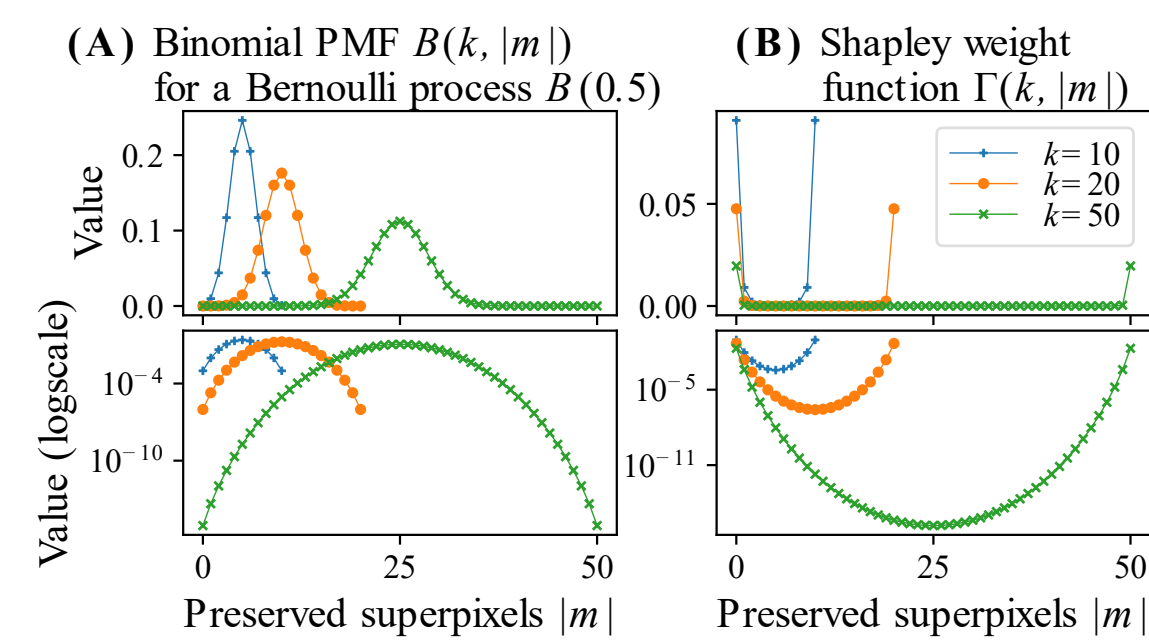
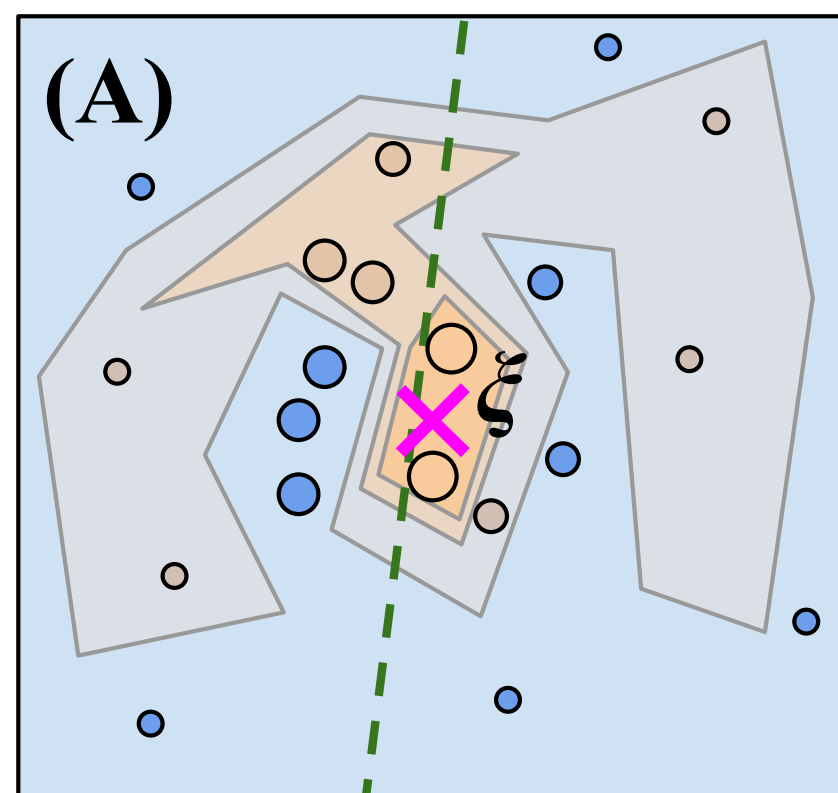
- LIME Image employs Monte Carlo sampling to generate synthetic neighborhoods
- The synthetic neighborhood aims to resemble a point cloud around the explained instance ξ , with perturbed samples varying in distance from ξ ; some closer, other further away.
- However, the use of the Bernoulli distribution B with coefficient $\frac{1}{2}$ concentrates the probability mass at the distribution center, allocating close-to-zero probability to the tails.
- Consequently, the synthetic neighborhood tends to look like a hypersphere positioned halfway between the original input and the fully-masked image.
- This can lead to a significant under-representation of the neighborhood, as very few synthetic samples closely resemble the original image.

PROPOSAL

- Integrate a stratified sampling approach into LIME equations to overcome the undersampling issues.
- Samples are drawn from the entire sampling space of the possible perturbed inputs, divided by strata.
- Add weight factors to samples, based on the frequency of the strata they belong to, to keep the estimator unbiased.

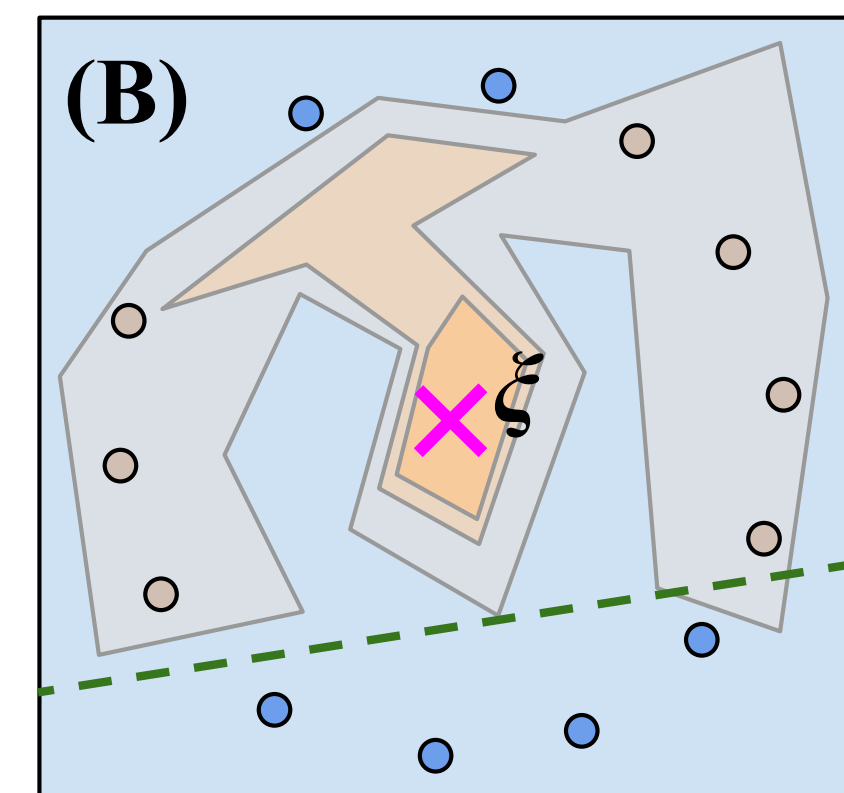
The problem with unbiased Bernoulli distribution

Ideally, the synthetic neighborhood $N(\xi)$ should provide a “good enough” coverage of the variations around ξ .



Samples at the tails of the Bernoulli distribution $B(0.5)$ are more rare than samples at the center of the distribution.

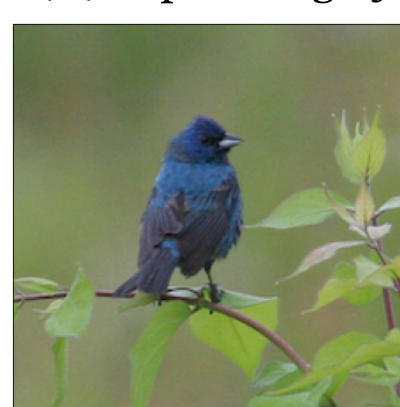
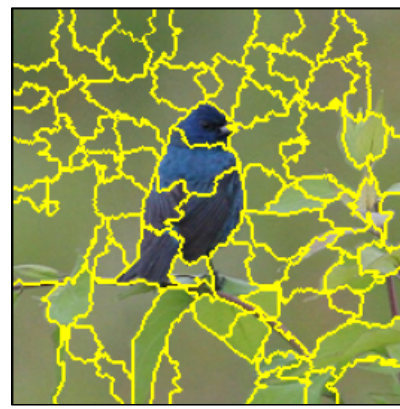
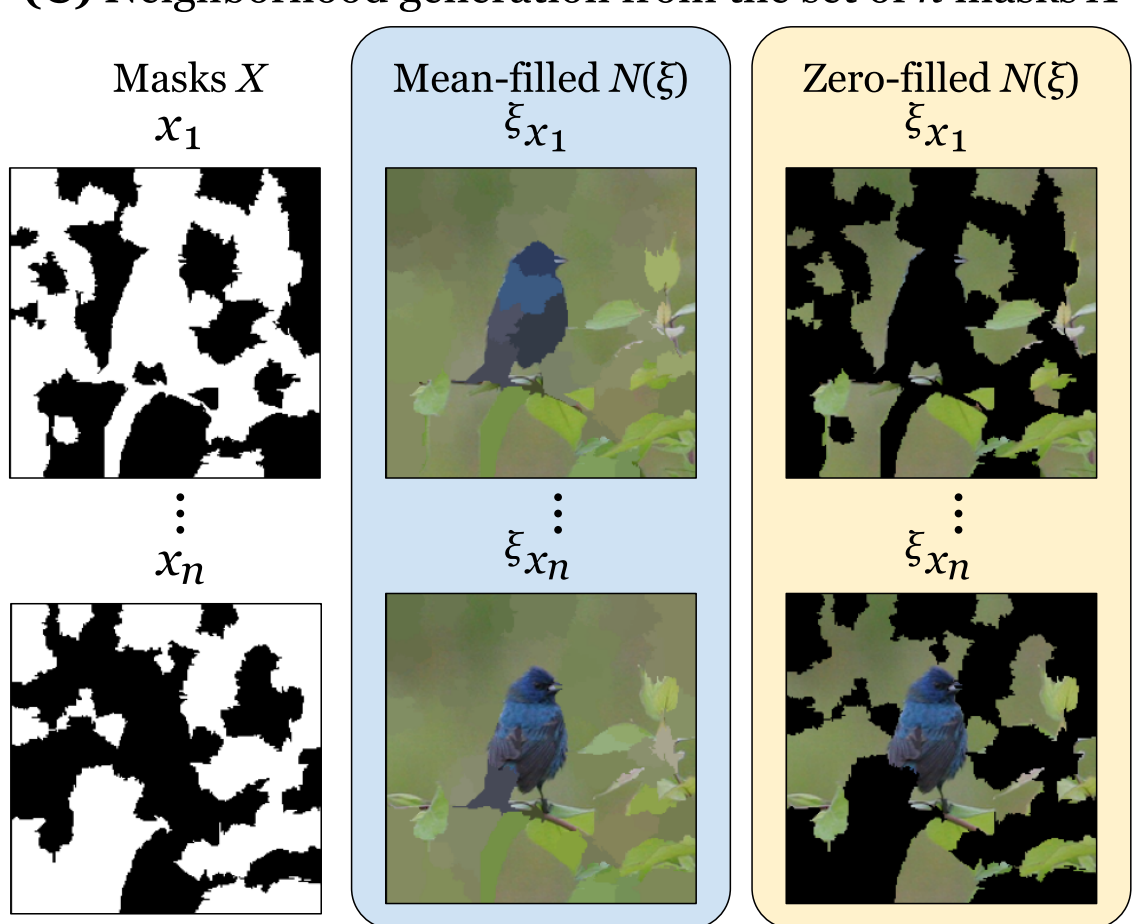
In practice the synthetic neighborhood $N(\xi)$ sampled by LIME Image looks like an hypersphere, with ver few to no samples close to the input image ξ .



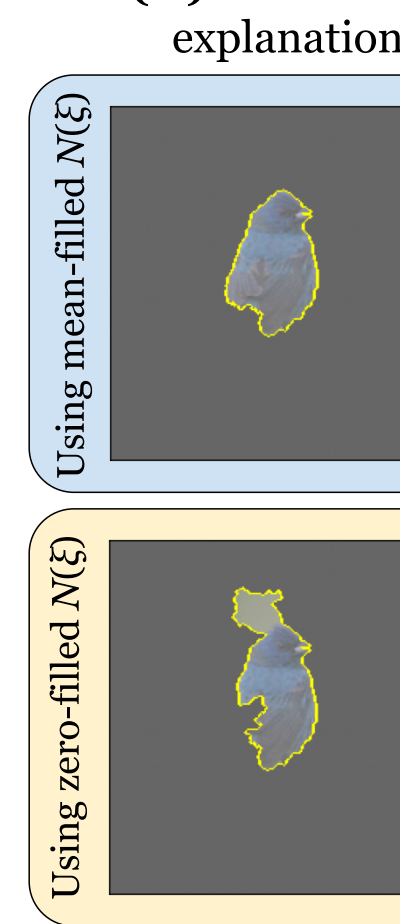
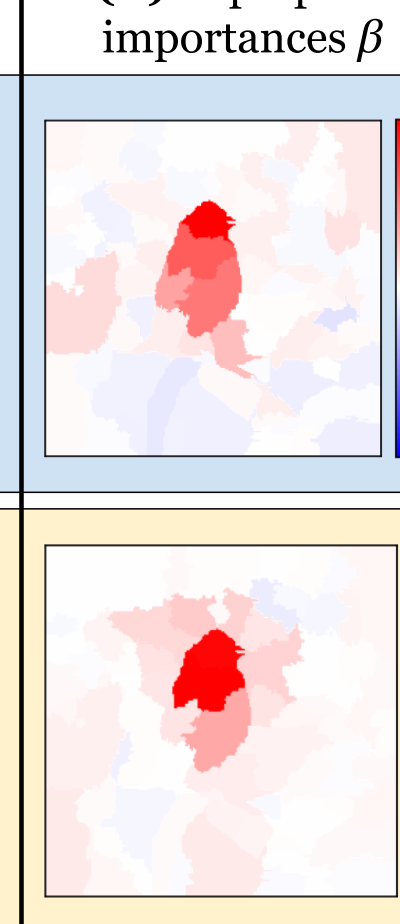
How standard LIME Image works

NOTATION

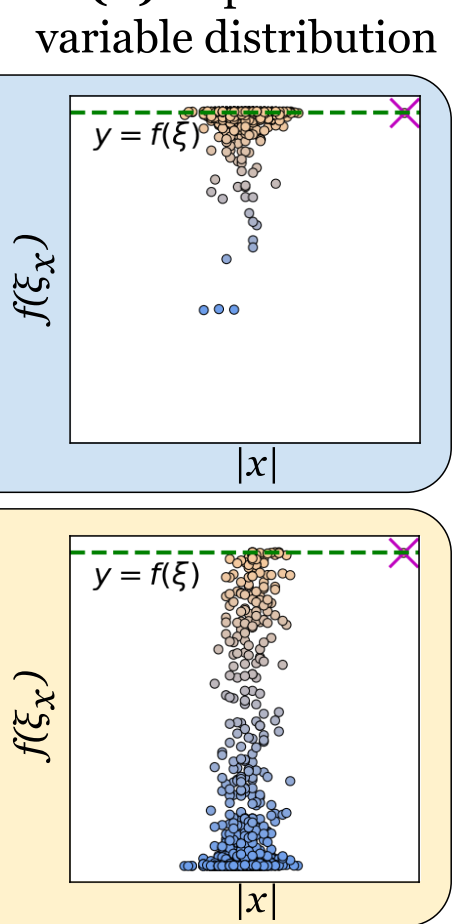
- Initial image $\xi \rightarrow$ divided into k superpixels
- Superpixel masking: $x \in \{0, 1\}^k$ generates a perturbed image ξ_x
- In LIME Image, masks are sampled using an unbiased Monte Carlo strategy: $x[i] \sim B(0.5)$, $1 \leq i \leq k$ where $B(p)$ is a Bernoulli-distributed random variable having probability $p = 0.5$
- $N(\xi)$ = synthetic neighborhood of n perturbed images
- Dependent variables: $Y = \{f(\xi_x) \mid \xi_x \in N(\xi)\}$

(A) Input image ξ (B) Superpixels $k=84$ (C) Neighborhood generation from the set of n masks X 

(D) LIME Image explanations

(E) Superpixel importances β 

(F) Dependent variable distribution



Dependent variable undersampling

Evaluation metrics

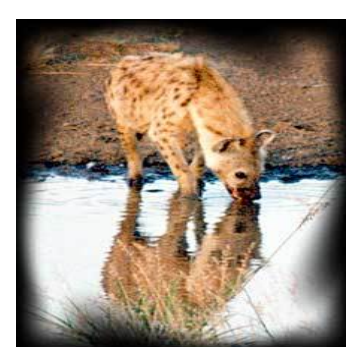
- Coefficient of Variation of the explanation β
 $CV(\beta) = \frac{\sigma_\beta}{\mu_\beta}$
- Range Coverage of the values of the dependent variable Y in the synthetic neighborhood.
 $RC(Y) = \frac{IQR_{1-99}(Y)}{f(\xi)}$

Observations

- Almost no sample is close to ξ .
- Samples drawn from $B(p)$ have ~50% of the superpixels masked.
- Under-representation of the local behaviour of the black-box model f .

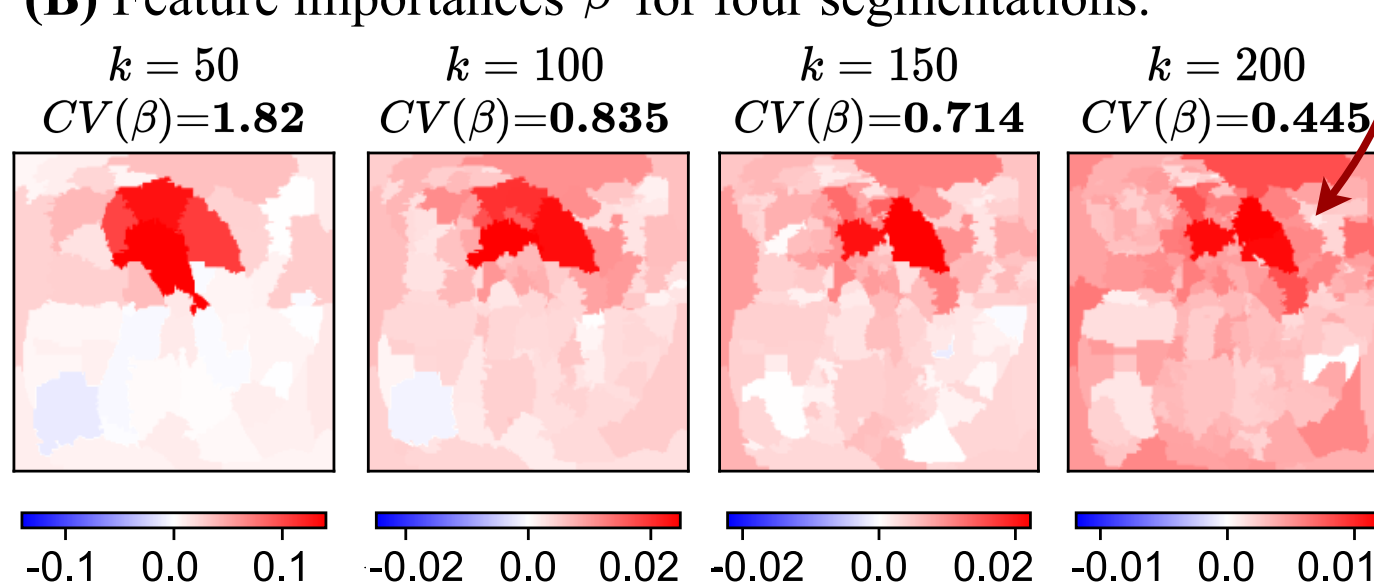
Stratified sampling

- Consider a partitioning stratified on the number of masked superpixels.
- $\mathcal{X}^{(i)}$ = set of possible masks having $|x| = i$
- Stratum i size is known a priori:
 $|\mathcal{X}^{(i)}| = \binom{k}{i}$, $0 \leq i \leq k$
- Oversampled probability
 $Prob\{x \in \mathcal{X}^{(i)} \mid x \in \hat{\mathcal{X}}\} = \frac{1}{k+1}$
- Adjustment factors
 $adj(i) = \frac{Prob\{x \in \mathcal{X}^{(i)} \mid x \in \mathcal{X}\}}{Prob\{x \in \mathcal{X}^{(i)} \mid x \in \hat{\mathcal{X}}\}} = \frac{(k+1)\binom{k}{i}}{2^k}$

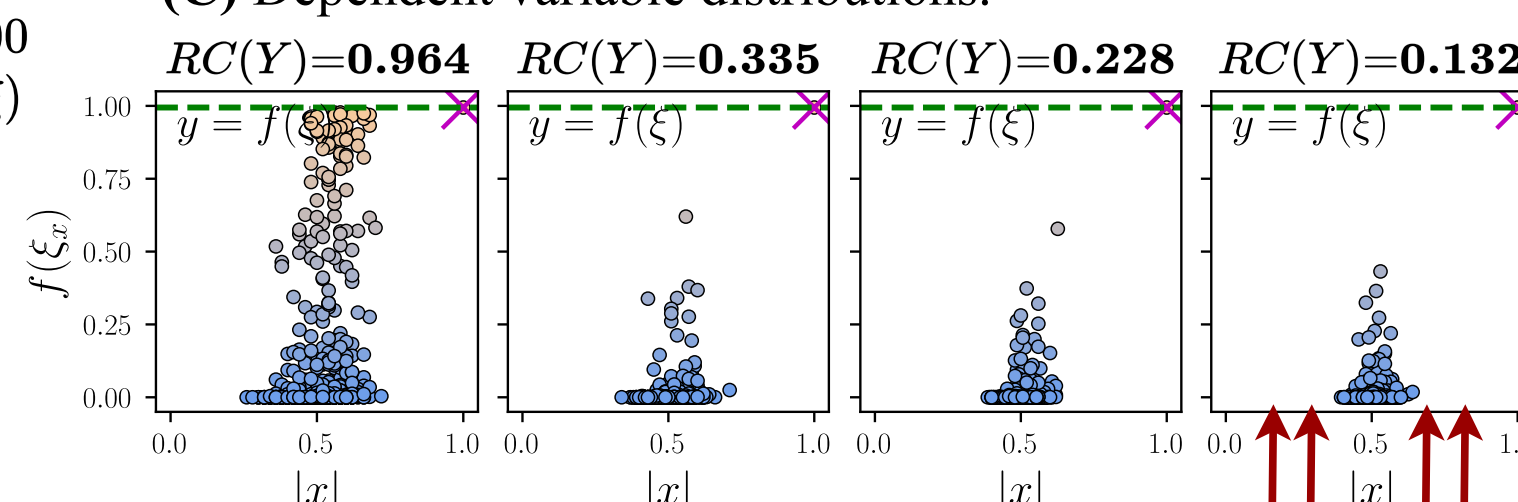


Input image ξ
Class: *hyena*
Probability: 99.46%
num_samples $n = 1000$
Using mean-filled $N(\xi)$

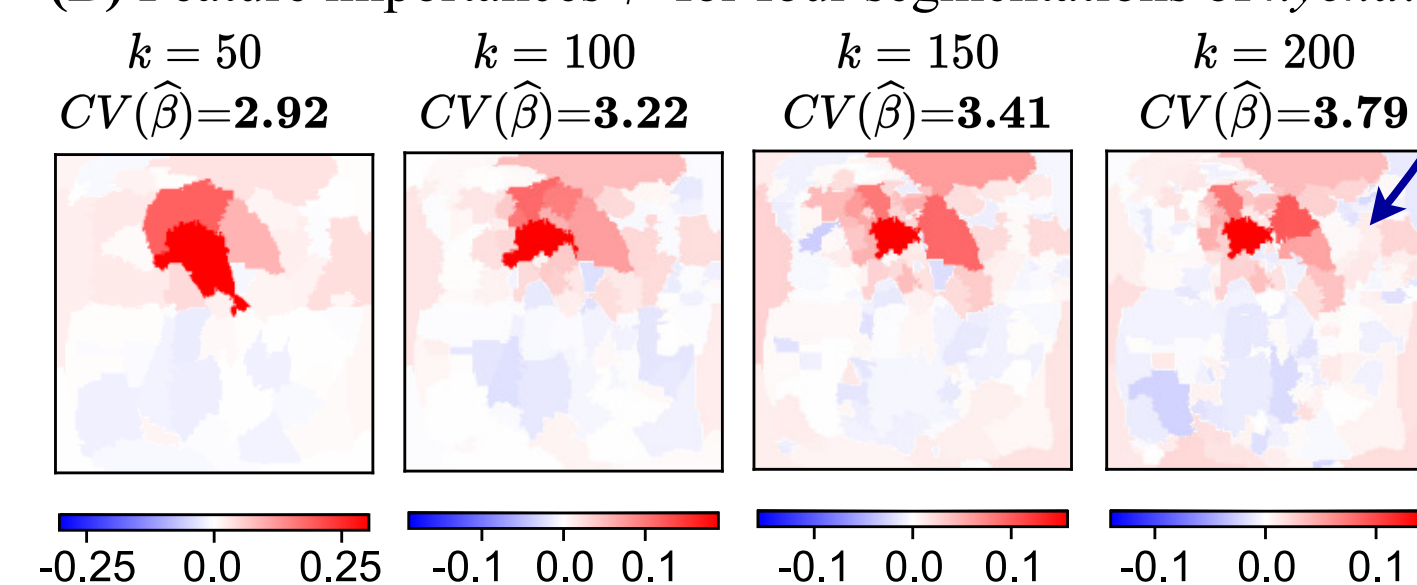
Monte Carlo sampling (default)

(B) Feature importances β for four segmentations.

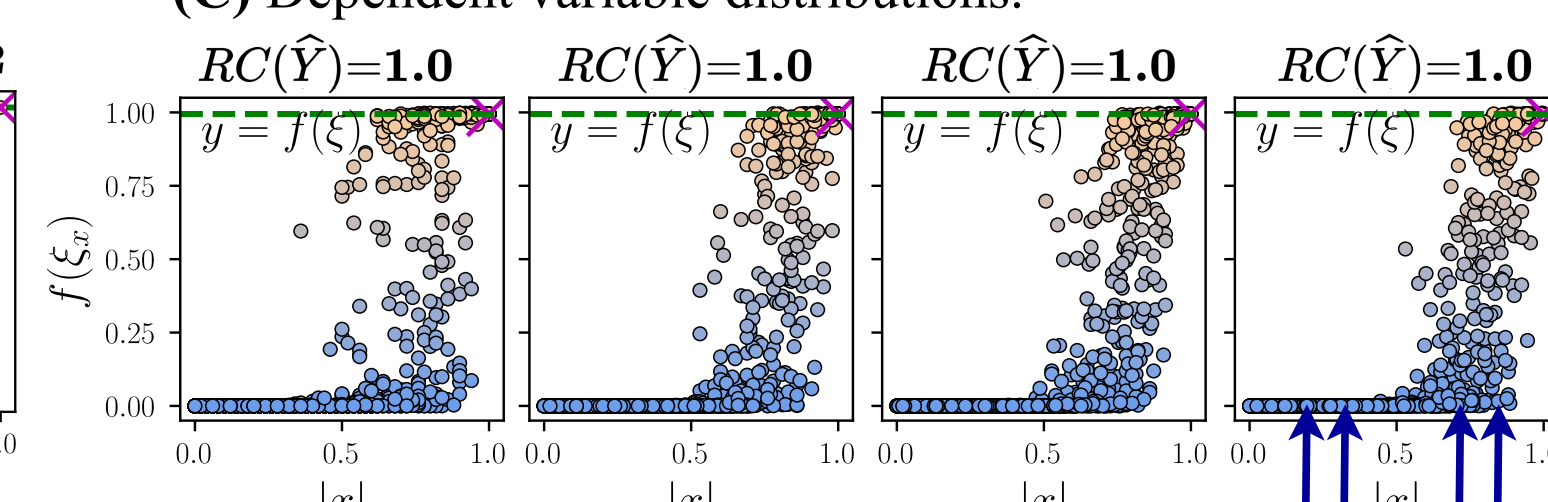
(C) Dependent variable distributions.



Stratified sampling

(B) Feature importances $\hat{\beta}$ for four segmentations of *hyena*.

(C) Dependent variable distributions.



Monte Carlo vs Stratified Sampling

Monte Carlo sampling

- LIME Image uses a simple linear homoscedastic model
 $Y = X \cdot \beta + \epsilon$
- The explanation coefficients β results from
 $\beta = (X^T W X)^{-1} X^T W Y$

Stratified sampling

- β coefficients may vary by stratum.
- We adopt instead a mixture model
 $\hat{Y}^{(i)} = \hat{X}^{(i)} \cdot \hat{\beta}^{(i)} + \hat{\epsilon}^{(i)}$
- The explanation coefficients β results from
 $\hat{\beta} = (\hat{X}^T \hat{W} \hat{X})^{-1} \hat{X}^T \hat{W} \hat{Y}$
- \hat{W} accounts for the adjustment factors that correct the bias introduced by oversampling the distribution tails.

Impact of stratified sampling in LIME Image

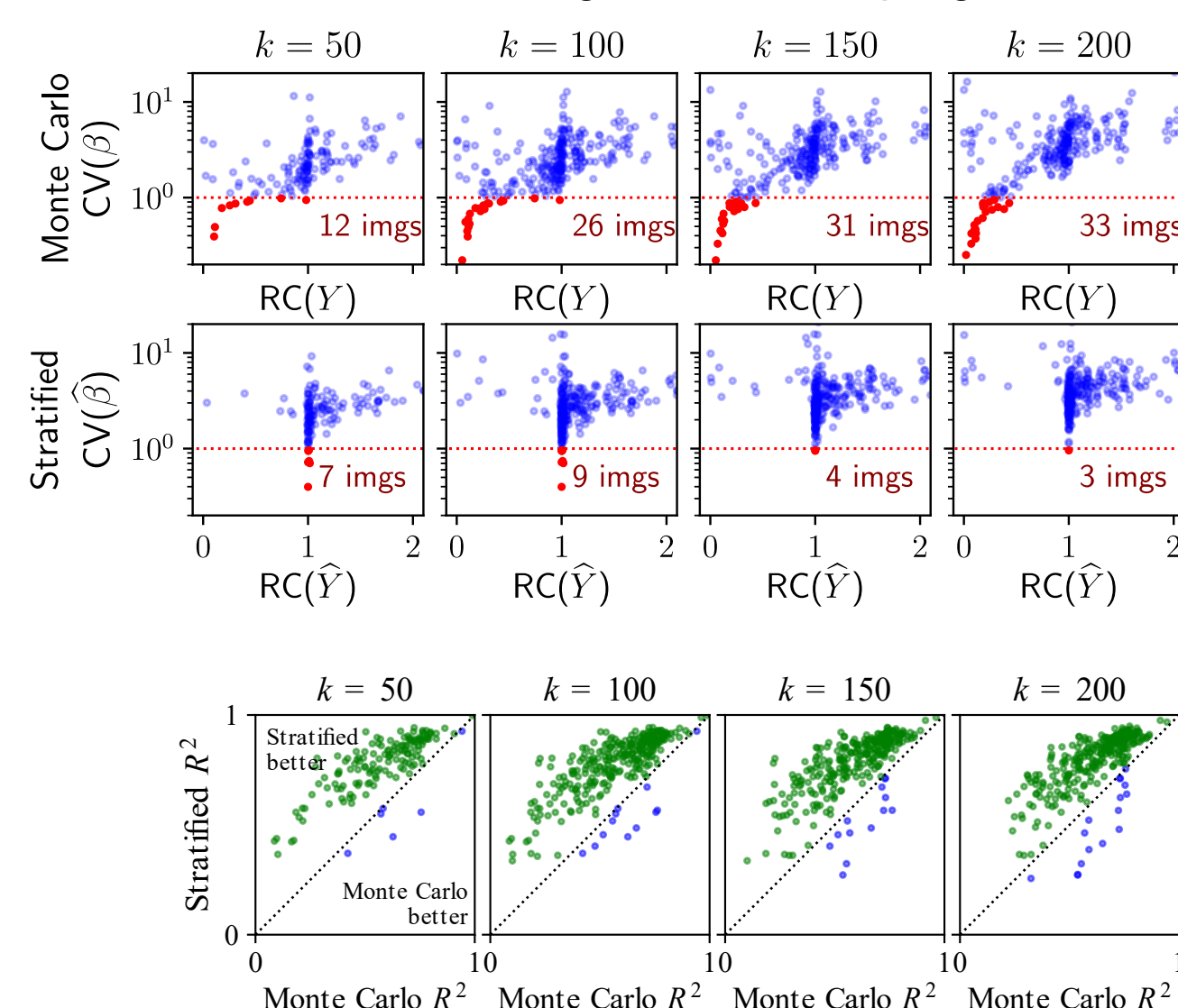
- Case (A):** The mean and variance of $\hat{\beta}^{(i)}$ are independent from the strata.
- weighted regression model is not needed.
- Monte Carlo and stratified sampling should behave similarly.
- Unlikely to happen using complex black-box machine learning models.

Case (B):

- The mean and variance of $\hat{\beta}^{(i)}$ varies by stratum.
- weighted regression model is highly advisable (DuMouchel and Duncan 1983).
- Monte Carlo will perform badly, stratified sampling is relevant.
- Common scenario for complex models and/or large num. of superpixels.

Evaluation & conclusions

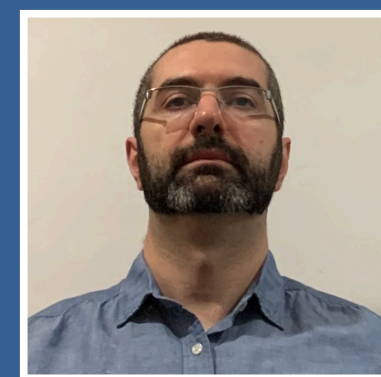
- About 1 image out of 5 in the ImageNet Object Localization Dataset suffers from severe undersampling using the default Monte Carlo sampling of LIME-Image
- Misbehaviours are corrected using stratified sampling



Conclusions

- Reformulation of LIME Image sampling strategy (not restricted to image data) for stratified sampling.
- Drawing lessons from the Shapley theory.
- Empirical evaluation shows that Monte Carlo undersampling is not rare, and stratified sampling provides practical improvements, at no additional cost.

Image	Monte Carlo sampling				Stratified sampling			
	(A) $k=50$ $CV(\beta) \mid RC(Y)$	(B) $k=200$ $CV(\beta) \mid RC(Y)$	(C) $k=50$ $CV(\hat{\beta}) \mid RC(\hat{Y})$	(D) $k=200$ $CV(\hat{\beta}) \mid RC(\hat{Y})$	(A) $k=50$ $CV(\beta) \mid RC(Y)$	(B) $k=200$ $CV(\beta) \mid RC(Y)$	(C) $k=50$ $CV(\hat{\beta}) \mid RC(\hat{Y})$	(D) $k=200$ $CV(\hat{\beta}) \mid RC(\hat{Y})$
orange	0.387 0.0833	0.178 0.0217	2.12 1.0	3.54 1.01	0.387 0.0833	0.178 0.0217	2.12 1.0	3.54 1.01
wardrobe	0.135 0.0182	0.164 0.00908	1.72 1.0	3.04 0.993	0.135 0.0182	0.164 0.00908	1.72 1.0	3.04 0.993
milk can	0.967 0.357	0.571 0.145	3.04 1.14	4.71 1.14	0.967 0.357	0.571 0.145	3.04 1.14	4.71 1.14
lynx	1.57 1.19	0.464 0.137	3.05 1.71	4.79 1.58	1.57 1.19	0.464 0.137	3.05 1.71	4.79 1.58
ringneck snake	1.10 0.365	0.221 0.0238	3.45 1.3	5.11 1.19	1.10 0.365	0.221 0.0238	3.45 1.3	5.11 1.19
chickadee	4.05 0.999	5.73 0.996	3.74 1.0	4.05 1.0	4.05 0.999	5.73 0.996	3.74 1.0	4.05 1.0
polecat	7.07 1.86	5.76 1.51	6.54 1.79	5.76 1.46	7.07 1.86	5.76 1.51	6.54 1.79	5.76 1.46



Elvio G. Amparore¹
elvio.g.amparore@unito.it



Muhammad Rashid¹
muhammad.rashid@unito.it



Enrico Ferrari²
enrico.ferrari@rulex.ai



Damiano Verda²
damiano.verda@rulex.ai

Source code:

https://github.com/rashidrao-pk/lime_stratified
<https://github.com/rashidrao-pk/lime-stratified-examples>

