Black Box Model

SMILE

# From Machine Learning to Generative AI: Advancing Statistical Model-Agnostic Interpretability with Local Explanations (SMILE)

# Table of Content

What I am going to discuss

**Highlights on eXplainable Artificial Intelligence (XAI)**

Some basic information.

**LIME Explainability Procedure**

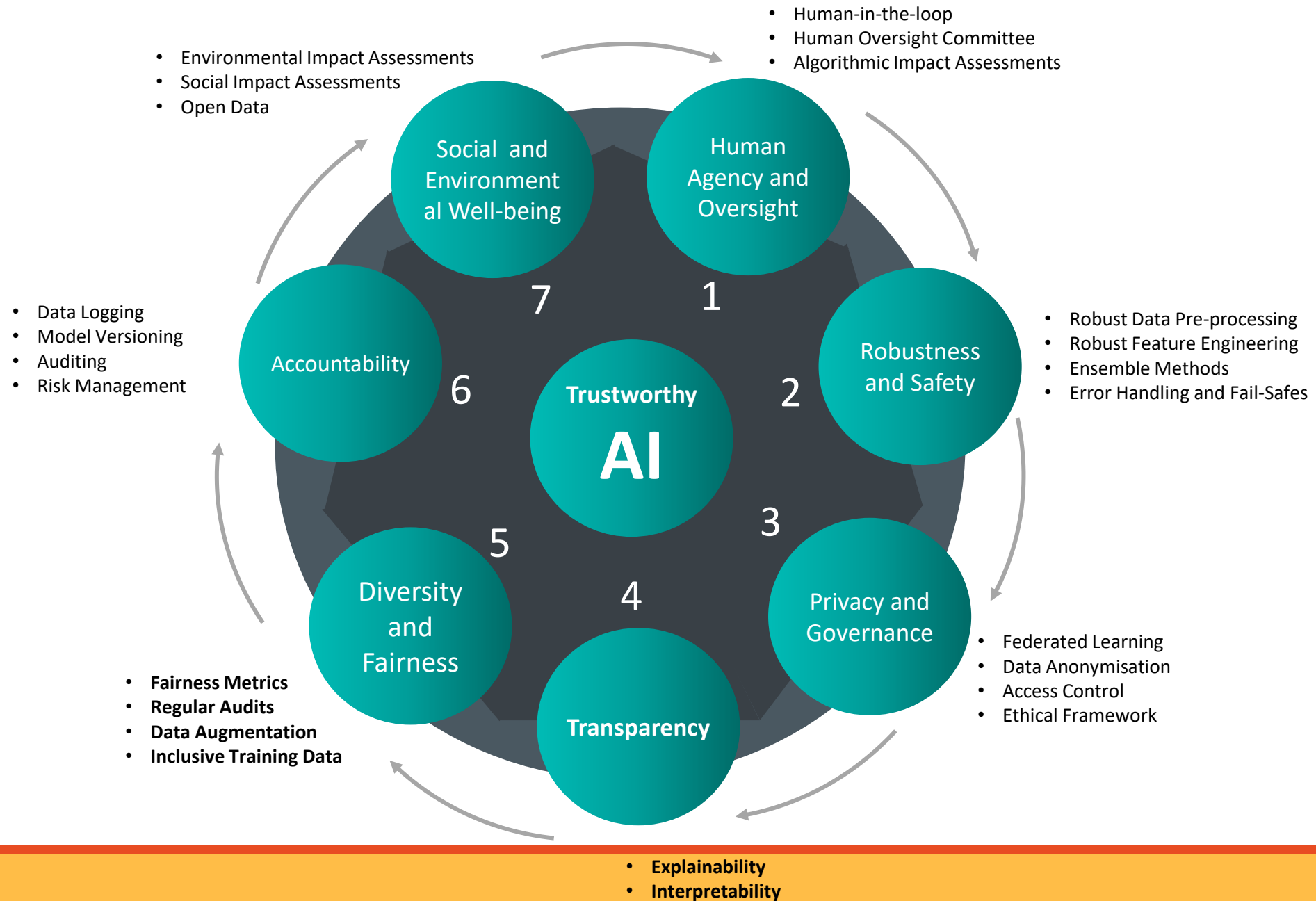Discussing how LIME works on a simple example.

**SMILE IDEA**

Discussing the idea and providing some experiments.

**Expanding SMILE for GenAI**

Discussing how SMILE has been expanded for GenAI (LLMs, MLLMs, etc).

# EU AI ACT



- Human-in-the-loop
- Human Oversight Committee
- Algorithmic Impact Assessments

- Environmental Impact Assessments
- Social Impact Assessments
- Open Data

Social and Environmental Well-being

Human Agency and Oversight

- Robust Data Pre-processing
- Robust Feature Engineering
- Ensemble Methods
- Error Handling and Fail-Safes

- Data Logging
- Model Versioning
- Auditing
- Risk Management

Accountability

7

1

6

**Trustworthy AI**

2

Robustness and Safety

5

3

4

Diversity and Fairness

Privacy and Governance

**Transparency**

- Federated Learning
- Data Anonymisation
- Access Control
- Ethical Framework

- **Fairness Metrics**
- **Regular Audits**
- **Data Augmentation**
- **Inclusive Training Data**

- **Explainability**
- **Interpretability**

3

# Key Questions

**What** explanation methods are applied, and are they appropriate for the model and user needs?

**Who** are the target users of explanations, and how are their needs and understanding levels considered?

**Where** in the system lifecycle is explainability integrated (e.g., model design, deployment)?

**Which** model components are explained, and what type of explanations are provided

**When** is explainability assessed or updated (e.g., continuously, post-deployment)?

# Key Stakeholders

## Data Scientist

- Technical report
- XAI metrics



- Mitigation strategies
- Improvement recommendations

## Business Owner

- Business implications
- Potential risks
- Opportunities
- Technical overview



## Regulators

- Compliance shown
- XAI evidence
- Mitigation details



## Consumer



**Examples:**
- User-friendly infographics.
- Online user feedback forms.
- Demographic-based testing.
- Personalized & non-biased product recommendations.
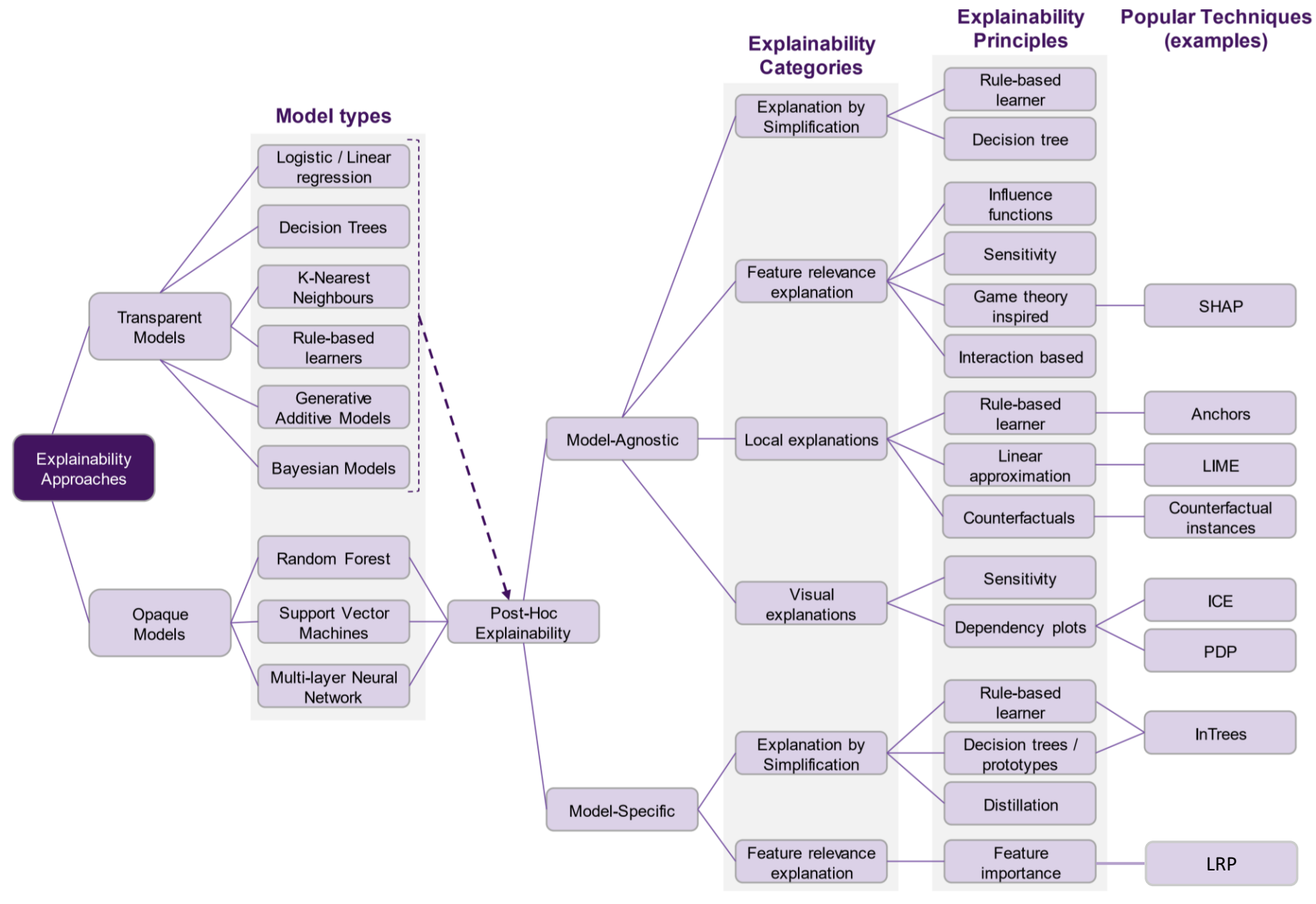
**Fairness Report:**
- Non-technical focus
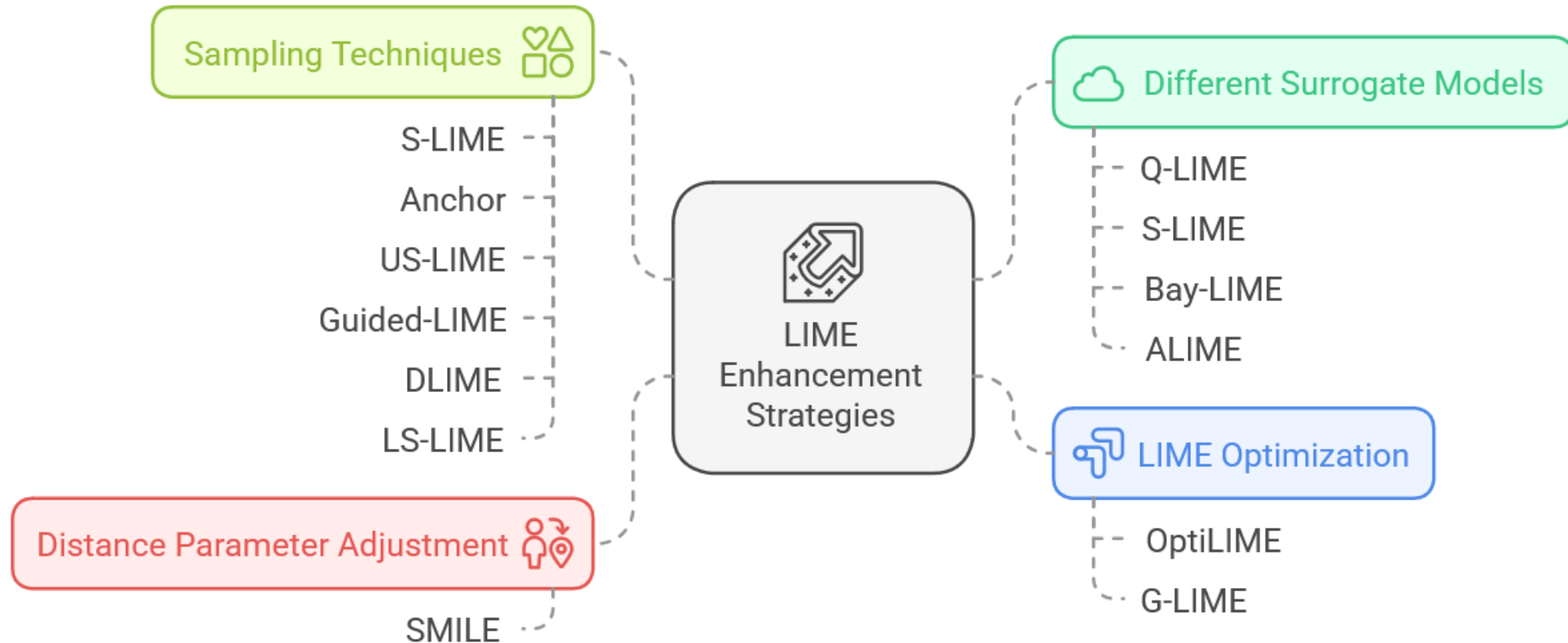- Model effects
- XAI steps
- Feedback avenues

## XAI Advocacy Groups

- Bias transparency
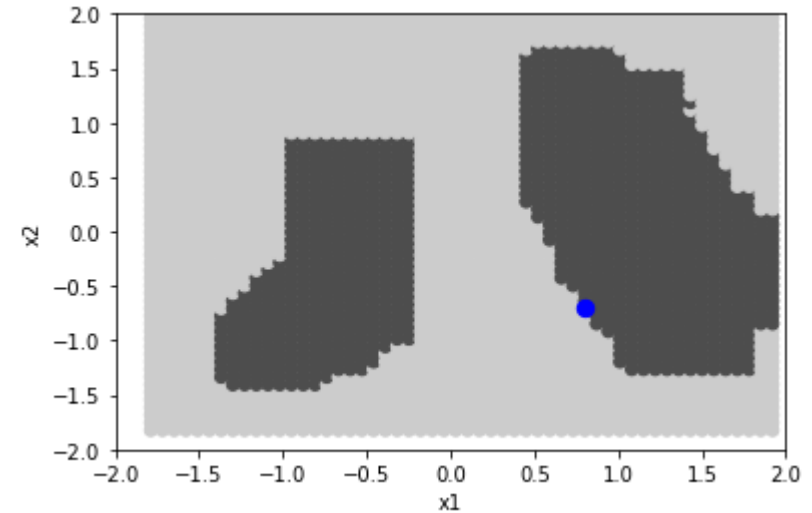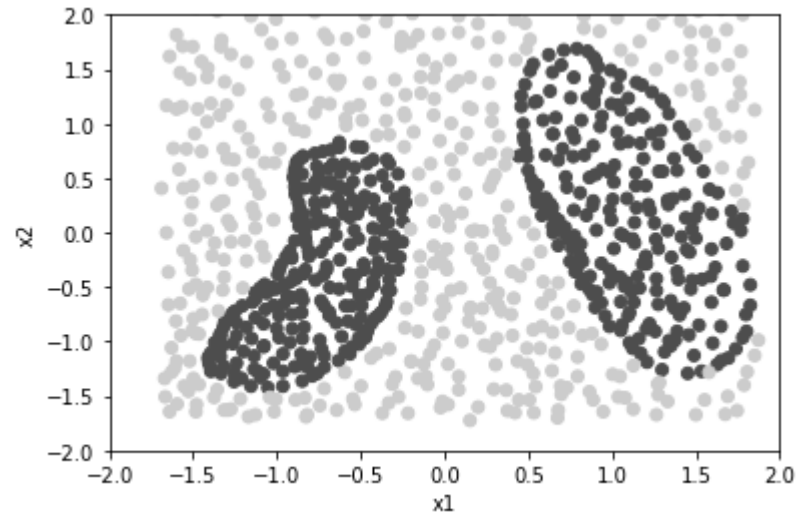- Societal impacts
- Addressing steps
- Real-world evidence

Belle, V., & Papantonis, I. (2021). Principles and practice of explainable machine learning. *Frontiers in big Data*, 39.

# LIME-based Approaches



Sampling Techniques
- S-LIME
- Anchor
- US-LIME
- Guided-LIME
- DLIME
- LS-LIME

LIME Enhancement Strategies

Different Surrogate Models
- Q-LIME
- S-LIME
- Bay-LIME
- ALIME

LIME Optimization
- OptiLIME
- G-LIME

Distance Parameter Adjustment
- SMILE

# LIME

Discussing LIME Algorithm in a simple tabular example

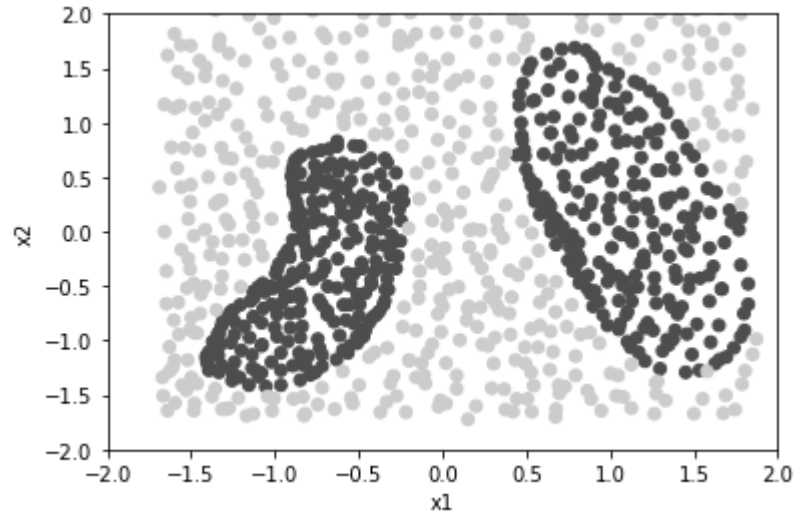# LIME: Local Interpretable Model-agnostic Explanations

Let's discuss a simple example of LIME on a tabular hypothetical dataset

# LIME: Local Interpretable Model-agnostic Explanations

Let's discuss a simple example of LIME on a tabular hypothetical dataset
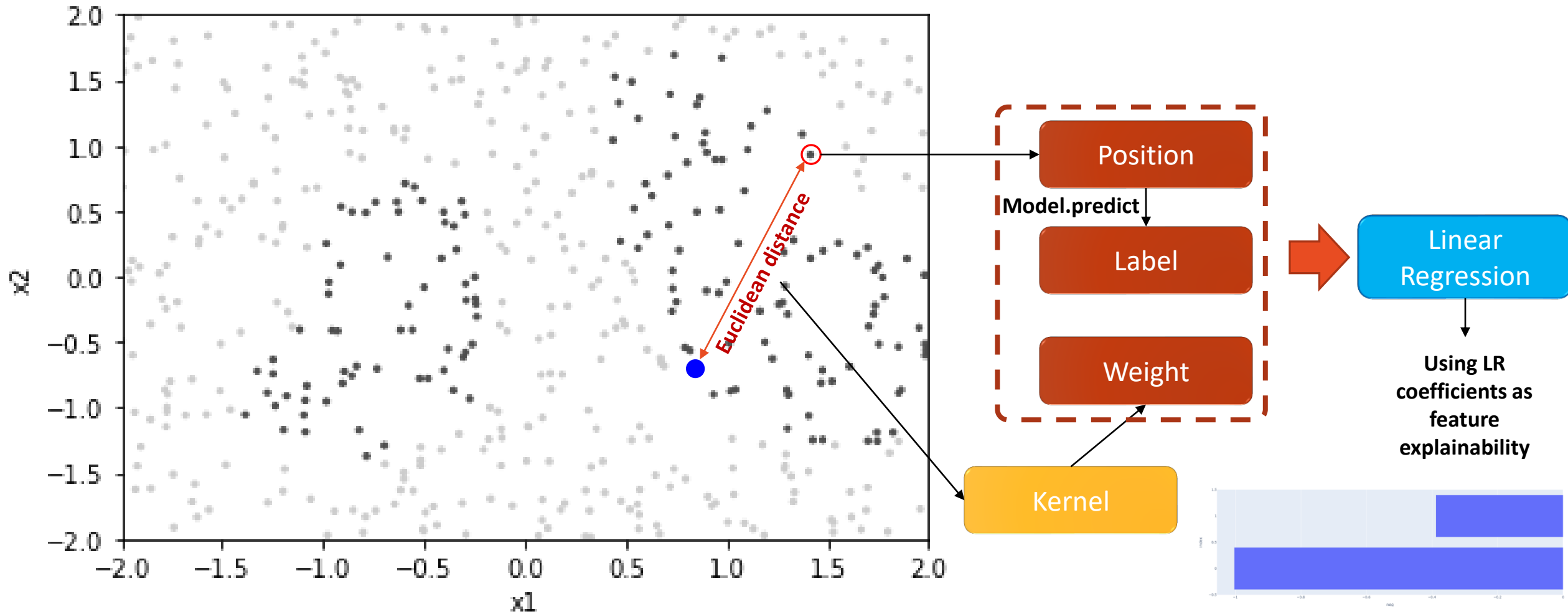


Generating uniform random numbers

Using model.predict(X_perturb)

**Source:** https://arteagac.github.io/blog

# LIME: Local Interpretable Model-agnostic Explanations

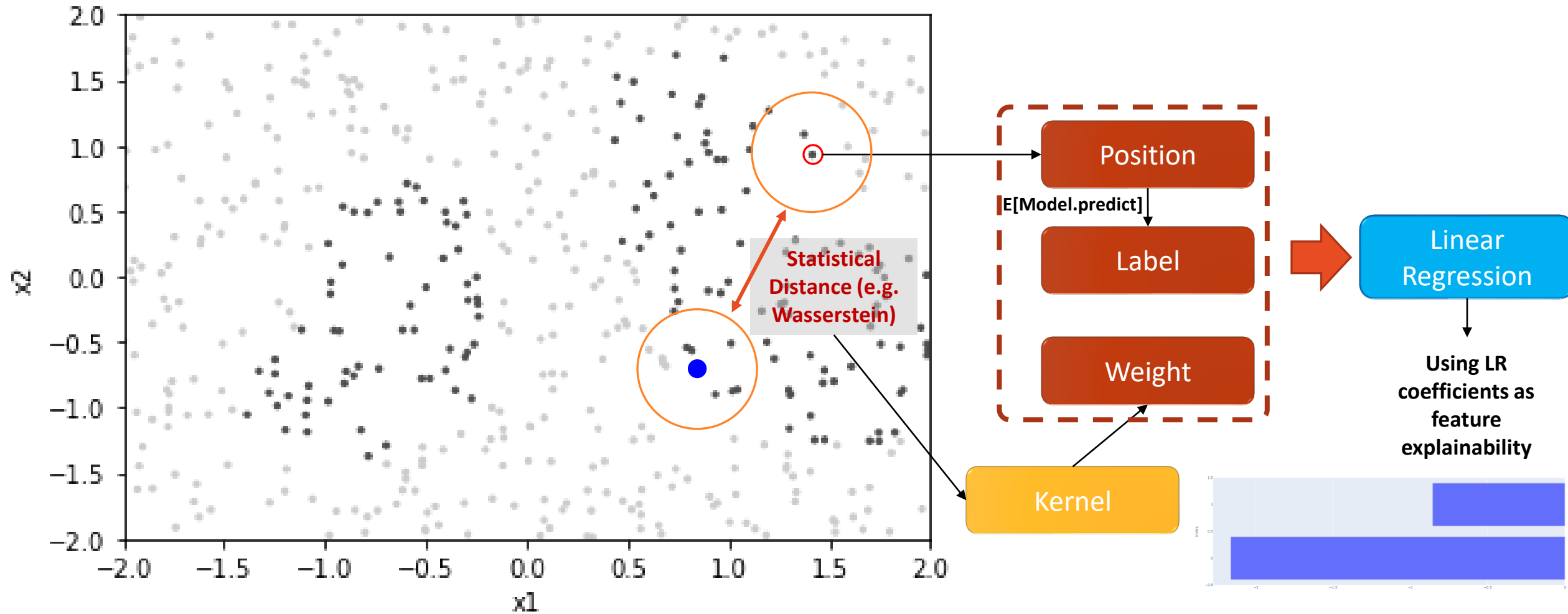Let's discuss a simple example of LIME on a tabular hypothetical dataset

# IDEA!

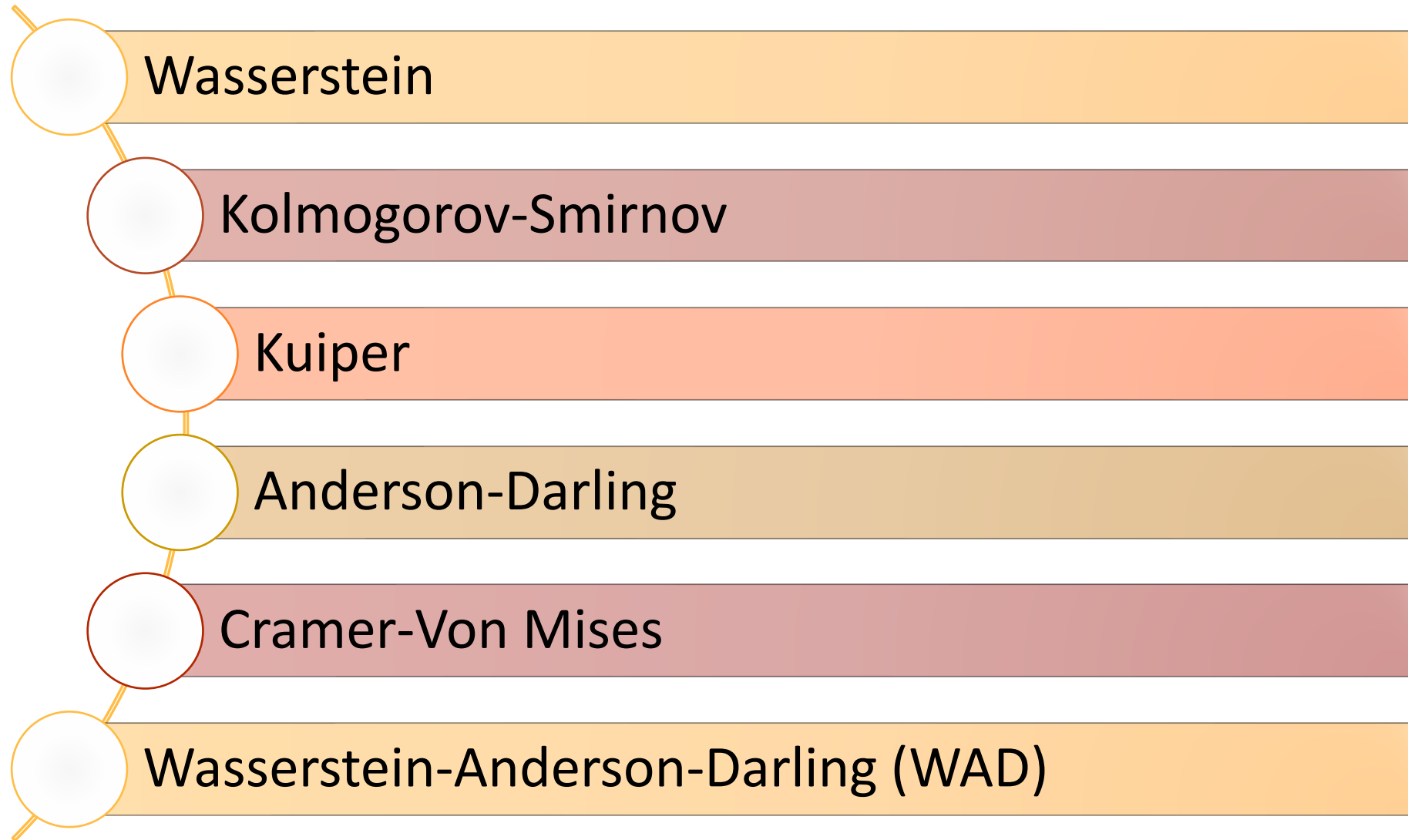SMILE: Statistical Model-agnostic Interpretability with Local Explanations

# SMILE: Statistical Model-agnostic Interpretability with Local Explanations

Let's discuss the same simple example with SMILE

# SMILE: Statistical Model-agnostic Interpretability with Local Explanations
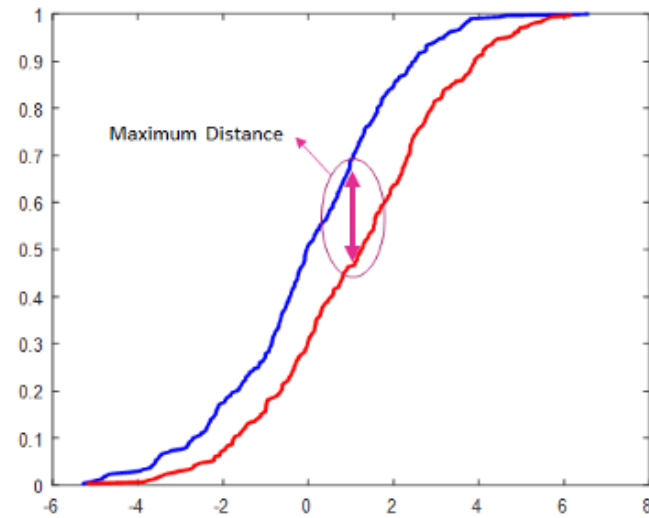
## Cumulative Distribution Function (CDF) Distance Measures

- Wasserstein
- Kolmogorov-Smirnov
- Kuiper
- Anderson-Darling
- Cramer-Von Mises
- Wasserstein-Anderson-Darling (WAD)

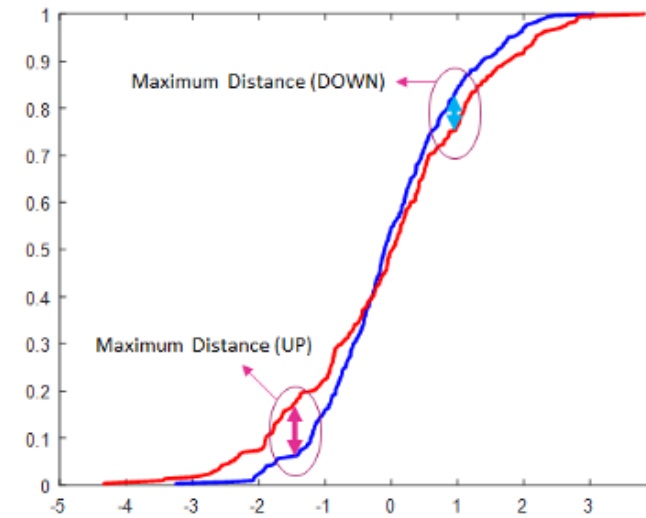# SMILE: Statistical Model-agnostic Interpretability with Local Explanations

Cumulative Distribution Function (CDF) Distance Measures

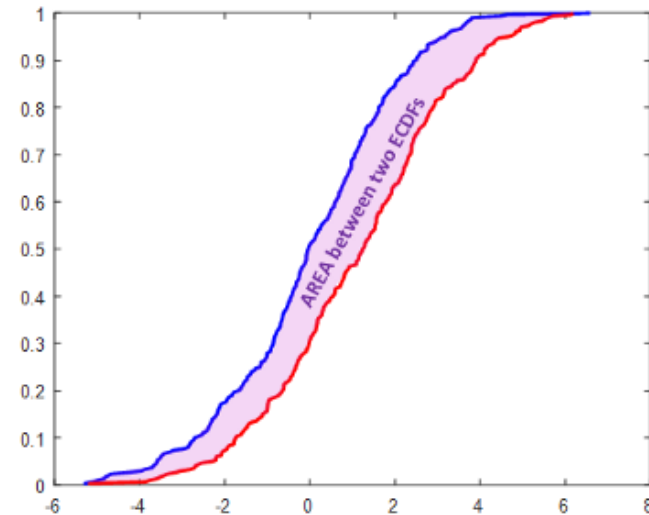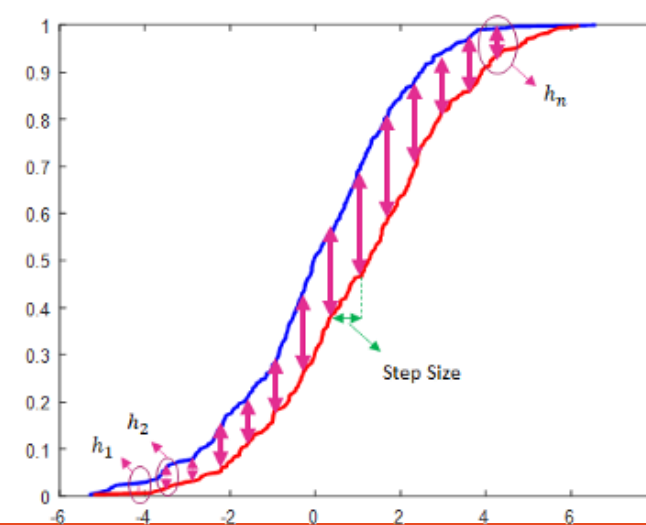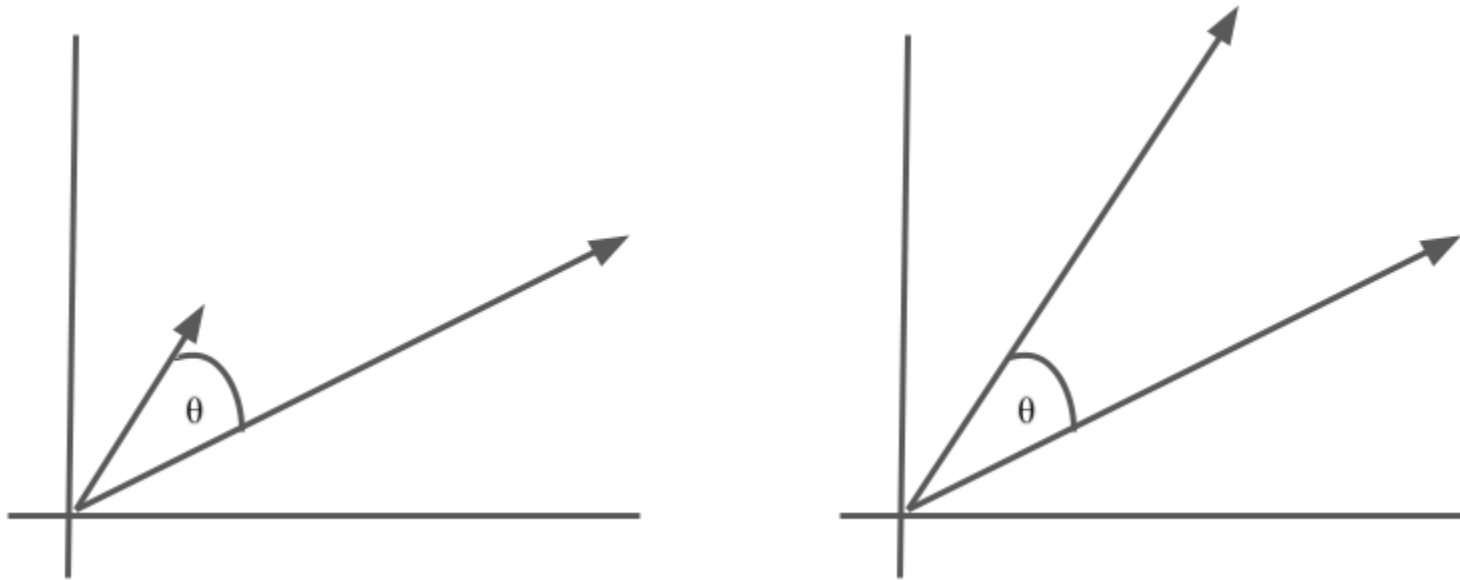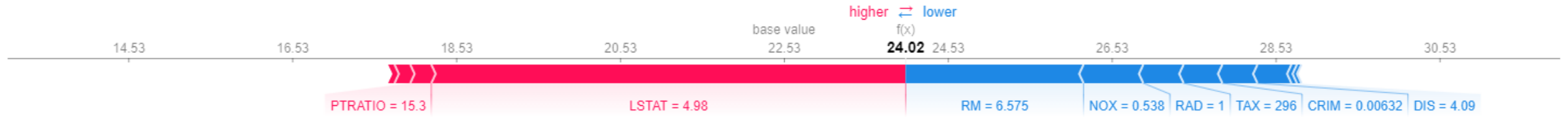# Wasserstein vs. Cosine

# Comparing SMILE with LIME and SHAP for Boston House Pricing Prediction



Comparing SHAP, LIME and SMILE

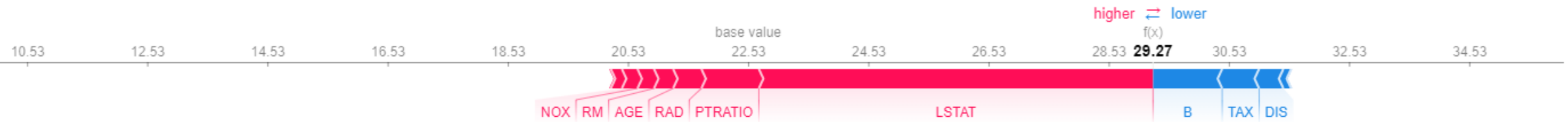# Comparing SMILE with LIME and SHAP with Force Plot



SHAP

LIME

SMILE

# Comparing with Human Feature Impact Estimates

Let's discuss a simple example of LIME on a tabular hypothetical dataset



Lundberg, S. M., & Lee, S. I. (2017, December). A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems* (pp. 4768-4777).

# Adversarial Attacks

Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods by Dylan Slack et al. (2020)

# Adversarial Attacks

## Fooling LIME and SHAP: Adversarial Attacks on Post-hoc Explanation Methods



Comparing SHAP, LIME and SMILE against Adversarial Attack - Unrealated Feature

\* COMPAS dataset for recidivism risk prediction

# SMILE for Images

Explaining the SMILE Procedure for local explainability in Image classifiers.

# LIME for Image Classifiers

How LIME provides explainability for image classifiers?



**Start**

**Input Image**

**Extracting Super-Pixels**

**Creating K Number of 0/1 Perturbations Sets**

**Switch off Super-Pixels based on Perturbations**

**Using Trained Classifier to Predict the Label**

**Computing Pairwise Cosine Distance Between Perturbation Vect. and Original Image Vect.**

**Using Kernel Function to Map Distances to [0 1] Weight Values (R, G,B)**

No

All Permutations Used?

Yes

**Fitting a weighed linear regression model using Perturbations, Predictions and Weights**

**Selecting M Super-Pixel that Gained Bigger Coef. In the Regression**

**END**

90% Class 2
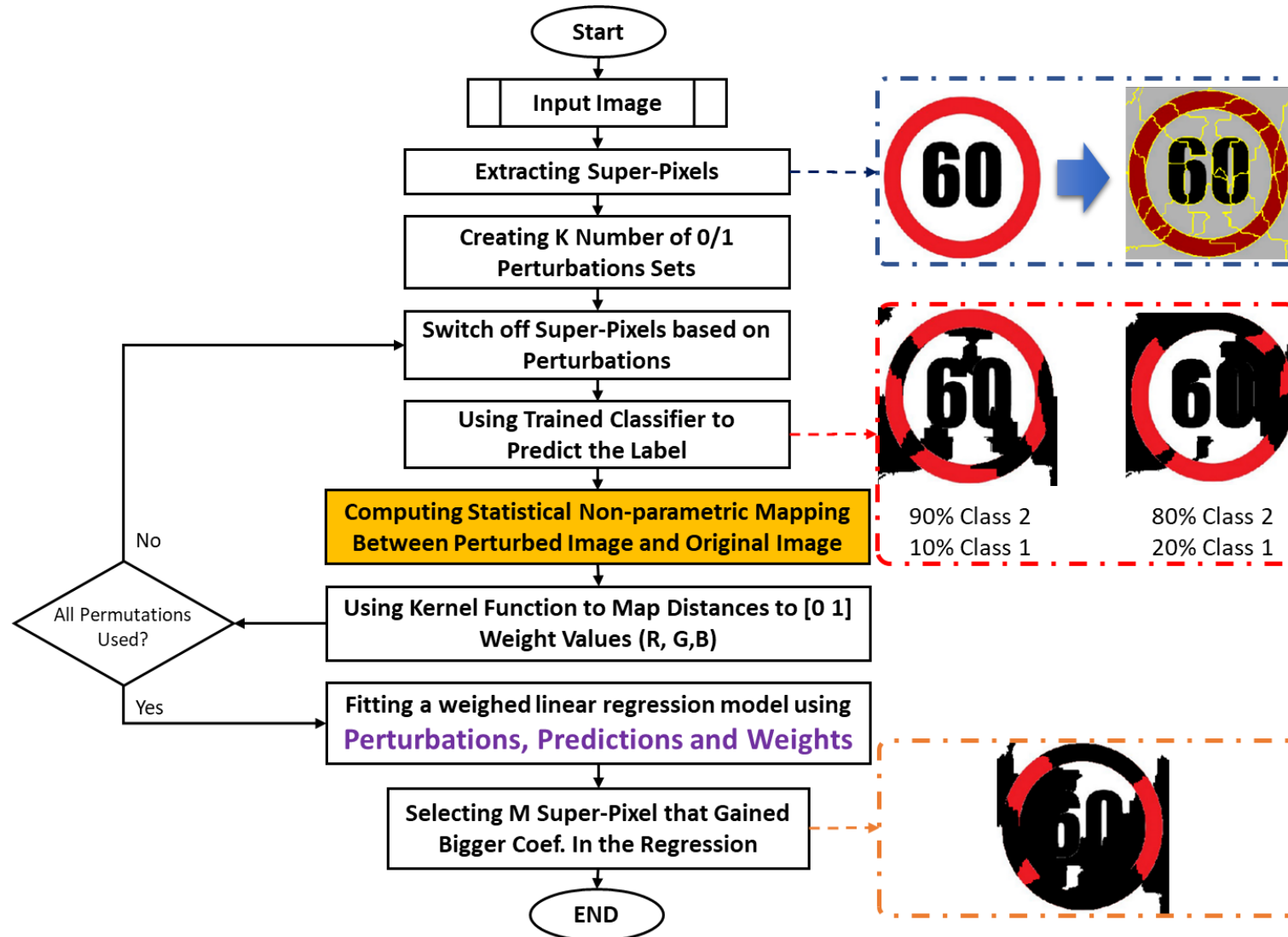10% Class 1

80% Class 2
20% Class 1

# SMILE for Image Classifiers

How SMILE provides explainability for image classifiers?

# SMILE for Image Classifiers
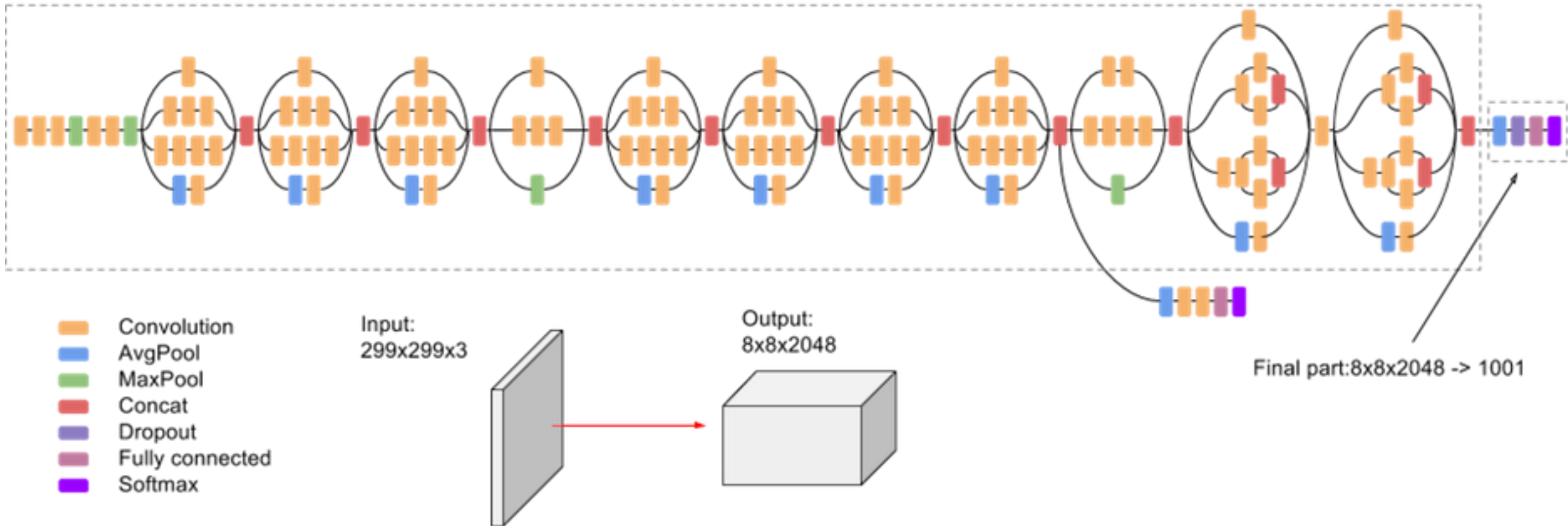
How SMILE provides explainability for image classifiers?

**InceptionV3**

Input: 299x299x3, Output:8x8x2048



Convolution
AvgPool
MaxPool
Concat
Dropout
Fully connected
Softmax

Input:
299x299x3

Output:
8x8x2048

Final part:8x8x2048 -> 1001

# Comparison between LIME and SMILE for Image Classification

**Original**

**LIME**

**SMILE**



[1] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). " Why should i trust you?" Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 1135-1144

25

# Comparison between LIME and SMILE for Image Classification

**Original**

**LIME**

**SMILE**



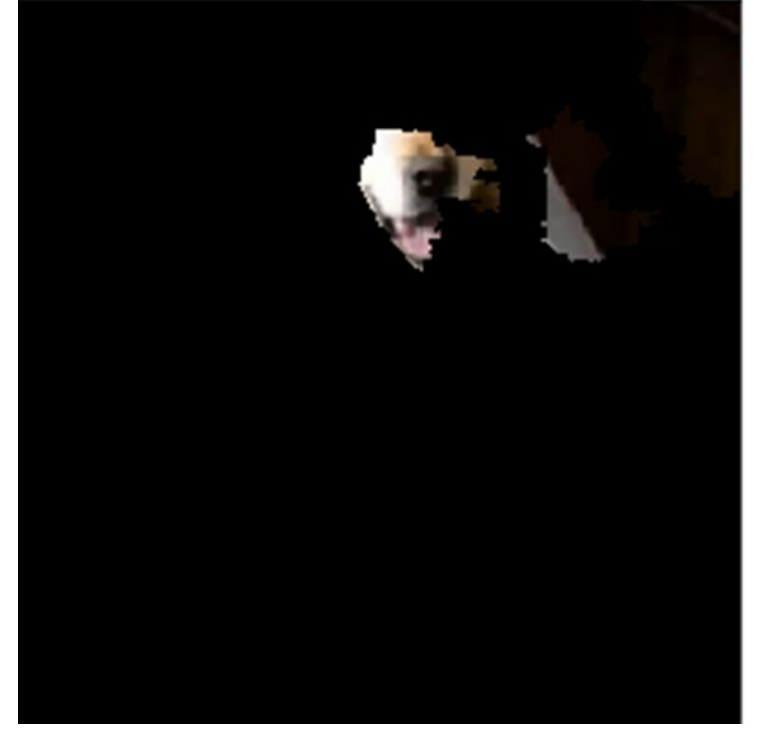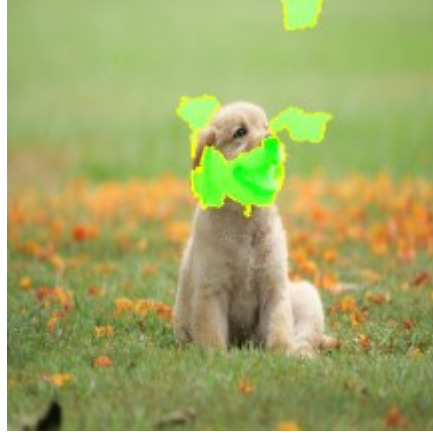In SMILE, as we process more information (images instead of perturbation vectors), we expect more accurate results.

# Comparison between BayLIME and SMILE for Image Classification

**Original**

**BayLIME – "non_Bay"**

**BayLIME – "Bay_non_info_prior"**

**SMILE**

**BayLIME – "Bay_info_prior"**

**BayLIME – "BayesianRidge_inf_prior_fit_alpha"**

# Heatmap Comparison between SMILE Distance Measure
# for Image Classification

Heatmap of LIME Coeffs - Cosine Distance

Heatmap of SMILE Coeffs - Kuiper Distance

Heatmap of SMILE Coeffs - KSD

Heatmap of SMILE Coeffs - WD

Heatmap of SMILE Coeffs - ADD

Heatmap of SMILE Coeffs - CVMD

# SMILE for Point Cloud

Explaining the SMILE Procedure for local explainability in Point Cloud classifiers.

**KernelSHAP***

**LIME****

**SMILE**

| C = 32 | C = 64 | C = 128 | C = 1024 |

References:
* Inspired by Lundberg, S.M. and S.-I. Lee, *A unified approach to interpreting model predictions*. Advances in neural information processing systems, 2017. **30**.
** Tan, H. and H. Kotthaus. *Surrogate model-based explainability methods for point cloud nns*. in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2022.

**LIME\***



**SMILE (all versions)**



| $\sigma = 0.1$ | $\sigma = 0.2$ | $\sigma = 0.3$ | $\sigma = 0.7$ |

References:
\* Tan, H. and H. Kotthaus. *Surrogate model-based explainability methods for point cloud nns*. in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2022.

# ModelNet40:

Misclassified:



| Predicted as 'Lamp' | Predicted as 'Plant' | Predicted as 'Bottle' | Predicted as 'Lamp' |

# SMILE for Texts

Explaining the SMILE Procedure for local explainability in Text classifiers.

# LIME vs SMILE for Text Classifiers

How LIME and SMILE provide explainability for text classifiers?



Left flowchart (LIME):

Start → Input Text Str → Text Pre-processing → Creating K Number of 0/1 Perturbations Sets → Remove words from the input text based on Perturbations → Using Trained Classifier to Predict the Label → Computing Pairwise Cosine Distance Between Perturbation Vect. and Original text Vect. → Using Kernel Function to Map Distances to Weight Values → All Permutations Used? (No → Remove words from the input text based on Perturbations; Yes → Fitting a weighed linear regression model using **Perturbations, Predictions and Weights**) → Selecting M Words that Gained Bigger Coef. In the Regression → END

Right flowchart (SMILE):

Start → Input Text Str → Text Pre-processing → Creating K Number of 0/1 Perturbations Sets → Remove words from the input text based on Perturbations → Using Trained Classifier to Predict the Label → **Computing Statistical Distance Between Perturbed Text and Original Text using Word2Vec Embeddings** → Using Kernel Function to Map Distances to Weight Values → All Permutations Used? (No → Remove words from the input text based on Perturbations; Yes → Fitting a weighed linear regression model using **Perturbations, Predictions and Weights**) → Selecting M Words that Gained Bigger Coef. In the Regression → END

# SMILE for Text Classifiers - Visualization

"For those who believe in God most of the big questions are answered But for those of us who can't readily accept the God formula the big answers don't remain stone-written. We adjust to new conditions and discoveries. We are pliable. Love need not be a command nor faith a dictum. I am my own god. We are here to unlearn the teachings of the church state, and our educational system. We are here to drink beer. We are here to kill war. We are here to laugh at the odds and live our lives so well that Death will tremble to take us - Charles Bukowski"

# SMILE vs LIME and SHAP for Text Classifiers:
## Quora Insincere Questions Classification

**"Is it just me or have you ever been in this phase wherein you became ignorant to the people you once loved, completely disregarding their feelings/lives so you get to have something go your way and feel temporarily at ease. How did things change?"** True Class: sincere Score=0.3597



Comparing SMILE and LIME (Negative: Sincere, Positive: Insincere)

# SMILE vs ELi5 Visualisation for Text Classifiers:
## Quora Insincere Questions Classification

**"Is it just me or have you ever been in this phase wherein you became ignorant to the people you once loved, completely disregarding their feelings/lives so you get to have something go your way and feel temporarily at ease. How did things change?"**



True Class: sincere Score=0.3597
y=**sincere** (probability **0.548**, score **-0.192**) top features

| Contribution[?] | Feature |
|---|---|
| +1.144 | <BIAS> |
| -0.951 | Highlighted in text (sum) |

is it just me or have you ever been in this phase wherein you became ignorant to the people you once loved, completely disregarding their feelings/lives so you get to have something go your way and feel temporarily at ease. how did things change?



Is it just me or have you ever been in this phase wherein you became ignorant to the people you once loved, completely disregarding their feelings/lives so you get to have something go your way and feel temporarily at ease. How did things change?

# SMILE for Generative AI

Explaining the SMILE Procedure for local explainability in Text classifiers.

# gSMILE for Generative Models (LLMs, VLMs, MLLMs)

Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657. (Submitted to IEEE Transactions on Artificial Intelligence).*

Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657. (Submitted to IEEE Transactions on Artificial Intelligence).*

# A Change in Fidelity Assessment for LLMs

Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657. (Submitted to IEEE Transactions on Artificial Intelligence).*

# The use of gSMILE for Bias and Fairness Evaluation

Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657. (Submitted to IEEE Transactions on Artificial Intelligence).*

# Model-agnostics vs. Model-specific

Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657. (Submitted to IEEE Transactions on Artificial Intelligence).*

# gSMILE for Multimodal Models

Dehghani, Z., Aslansefat, K., Khan, A., Rivera, A. R., George, F., & Khalid, M. (2024). Mapping the Mind of an Instruction-based Image Editing using SMILE. *arXiv preprint arXiv:2412.16277. (Submitted to Nature npj Artificial Intelligence)*.

Dehghani, Z., Aslansefat, K., Khan, A., Rivera, A. R., George, F., & Khalid, M. (2024). Mapping the Mind of an Instruction-based Image Editing using SMILE. *arXiv* *preprint arXiv:2412.16277. (Submitted to Nature npj Artificial Intelligence).*

2D t-SNE Scatter Plot of Image Embeddings

Box Plot of Keywords Related to Snow with Variations

Box Plot of Categorized Snow-Related Keywords

Dehghani, Z., Aslansefat, K., Khan, A., Rivera, A. R., George, F., & Khalid, M. (2024). Mapping the Mind of an Instruction-based Image Editing using SMILE. *arXiv preprint arXiv:2412.16277. (Submitted to Nature npj Artificial Intelligence)*.

# KG-SMILE



Moghaddam, Z. Z. S., Dehghani, Z., Rani, M., Aslansefat, K., Mishra, B. K., Kureshi, R. R., & Thakker, D. (2025). Explainable Knowledge Graph Retrieval-Augmented Generation (KG-RAG) with KG-SMILE. *arXiv preprint arXiv:2509.03626*. (Submitted to Information Systems Frontiers).

# Conclusion

- ✓ SMILE extends LIME by integrating statistical distance measures (e.g., Wasserstein, Cramér–von Mises, Kuiper) to produce more stable and reliable local explanations.

- ✓ The method is completely model-agnostic, making it suitable for explaining classical ML, deep learning, and Generative AI (gSMILE) models.

- ✓ SMILE inherits the visual interpretability of LIME and the quantitative depth of SHAP, combining their advantages while improving robustness.

- ✓ Across images, text, and 3D point cloud modalities, SMILE consistently delivers more coherent and less noisy feature attributions.

- ✓ Experimental results demonstrate that SMILE is more resilient to adversarial manipulation than LIME and SHAP, though not entirely immune.

- ✓ The framework now supports LLMs, VLMs, and multimodal models through gSMILE, enabling explainability, fidelity analysis, and bias/fairness evaluation for generative tasks.

- ✓ Future directions include expanding SMILE for various directions, including Image Captioning, QML, and Geo-SMILE.

# SMILE Reproducibility

**https://github.com/Dependable-Intelligent-Systems-Lab/xwhy**

**https://github.com/koo-ec/KG-SMILE**

**https://github.com/Sara068/Mapping-the-Mind-of-an-Instruction-based-Image-Editing-using-SMILE**

**https://github.com/Sara068/LLM-SMILE**

**https://github.com/koo-ec/Geo-SMILE (Not Public Yet)**

# Publications

[1] Aslansefat, K., Hashemian, M., Walker, M., Akram, M. N., Sorokos, I., & Papadopoulos, Y. (2023). Explaining black boxes with a SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *IEEE Software*, *41*(1), 87-97.

[2] Ahmadi, S. M., Aslansefat, K., Valcarce-Dineiro, R., & Barnfather, J. (2024). Explainability of Point Cloud Neural Networks Using SMILE: Statistical Model-Agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2410.15374 (Submitted to Engineering Applications of Artificial Intelligence)*.

[3] Dehghani, Z., Aslansefat, K., Khan, A., Rivera, A. R., George, F., & Khalid, M. (2024). Mapping the Mind of an Instruction-based Image Editing using SMILE. *arXiv preprint arXiv:2412.16277*. *(Submitted to Nature npj Artificial Intelligence)*.

[4] Dehghani, Z., Aslansefat, K., Khan, A., & Akram, M. N. (2025). Explainability of Large Language Models using SMILE: Statistical Model-agnostic Interpretability with Local Explanations. *arXiv preprint arXiv:2505.21657*. *(Submitted to IEEE Transactions on Artificial Intelligence)*.

[5] Moghaddam, Z. Z. S., Dehghani, Z., Rani, M., Aslansefat, K., Mishra, B. K., Kureshi, R. R., & Thakker, D. (2025). Explainable Knowledge Graph Retrieval-Augmented Generation (KG-RAG) with KG-SMILE. *arXiv preprint arXiv:2509.03626*. (Submitted to Information Systems Frontiers).

# References

Here you can find the references that has been used in this presentation

[1] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). " Why should i trust you?" Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 1135-1144).

[2] Lundberg, S. M., & Lee, S. I. (2017, December). A unified approach to interpreting model predictions. In Proceedings of the 31st international conference on neural information processing systems (pp. 4768-4777).

[3] Slack, D., Hilgard, S., Jia, E., Singh, S., & Lakkaraju, H. (2020, February). Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (pp. 180-186).

[4] Ghorbani, A., Abid, A., & Zou, J. (2019, July). Interpretation of neural networks is fragile. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 33, No. 01, pp. 3681-3688).

[5] Belle, V., & Papantonis, I. (2021). Principles and practice of explainable machine learning. Frontiers in big Data, 39.

[6] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. Information Fusion, 58, 82-115.

# Thank You

If you have any question, please feel free to ask



Koorosh Aslansefat
Assistant Professor | Data Scientist
| Co-Lead of Dependable Intelligent Syste...

k.aslansefat@hull.ac.uk