

IBM Research

Generative Perturbation Analysis for Probabilistic Black-Box Anomaly Attribution

Tsuyoshi (Ide-san) Ide, Naoki Abe
{tide, nabe}@us.ibm.com
IBM T. J. Watson Research Center

Agenda

- What is the task, “Anomaly Attribution”?
- What’s wrong with the existing attribution methods?
- What is the new idea?
- Illustrative examples

“Anomaly attribution” is an important topic in XAI (explainable AI) research.

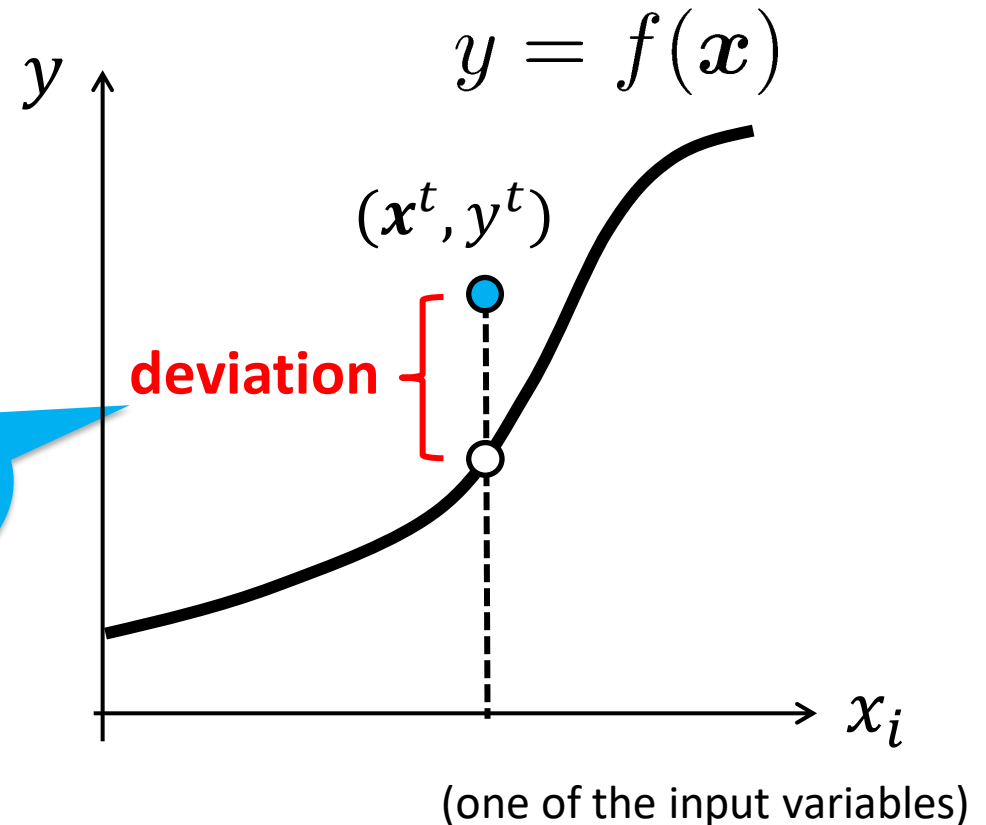
Given:

- Black-box regression model $y = f(x)$ and a (set of) test sample (x^t, y^t)
 - No access to the model beyond API
 - No access to the training data

Explain:

- The deviation $f(x^t) - y^t$
- by computing the attribution score (responsibility score) for *each* of the input variables x .

Why did I get this?



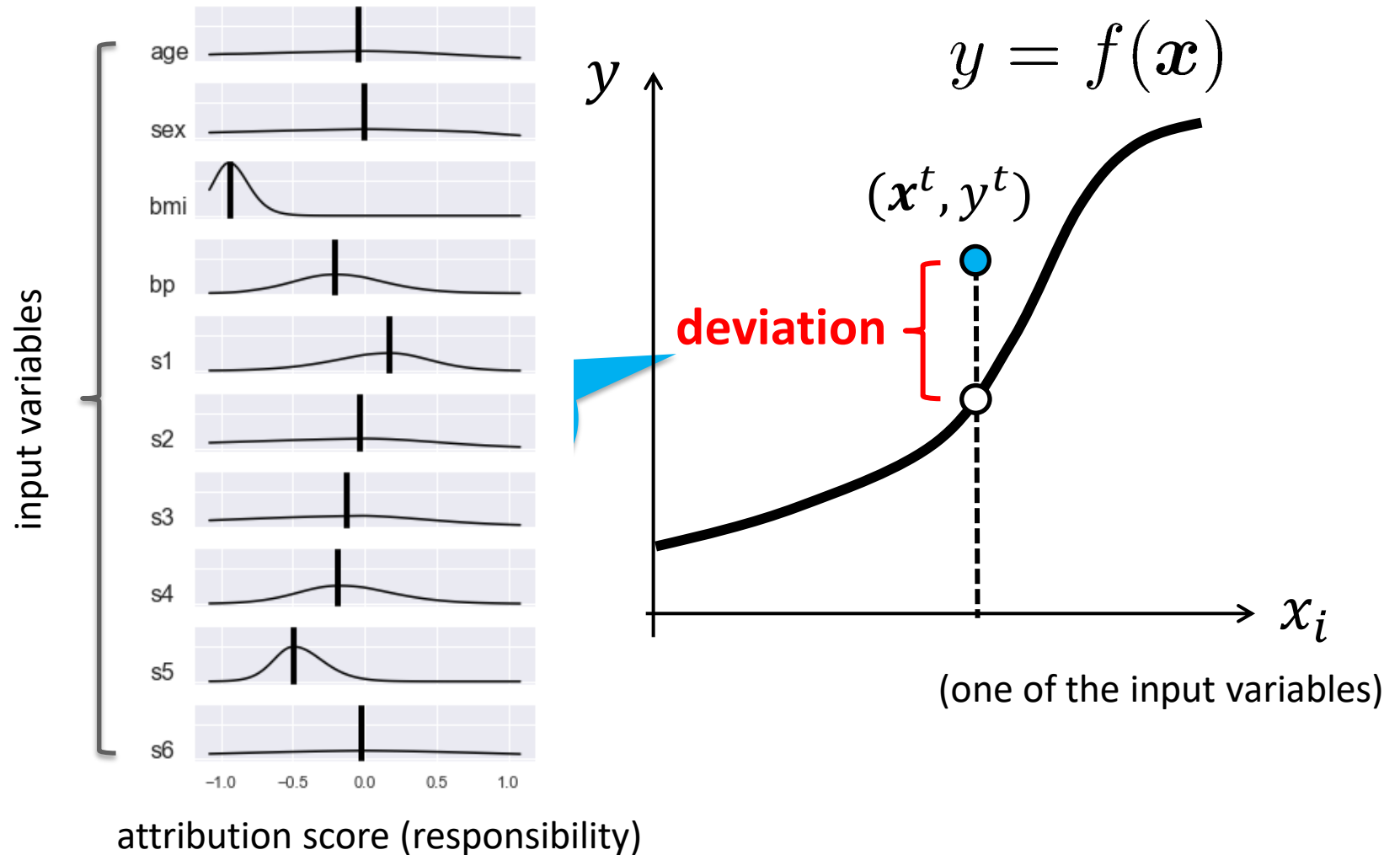
“Anomaly attribution” is an important topic in XAI (explainable AI) research.

New capabilities that GPA has enabled

GPA = generative perturbation analysis (proposed method)

deviation-sensitive

probabilistic



Agenda

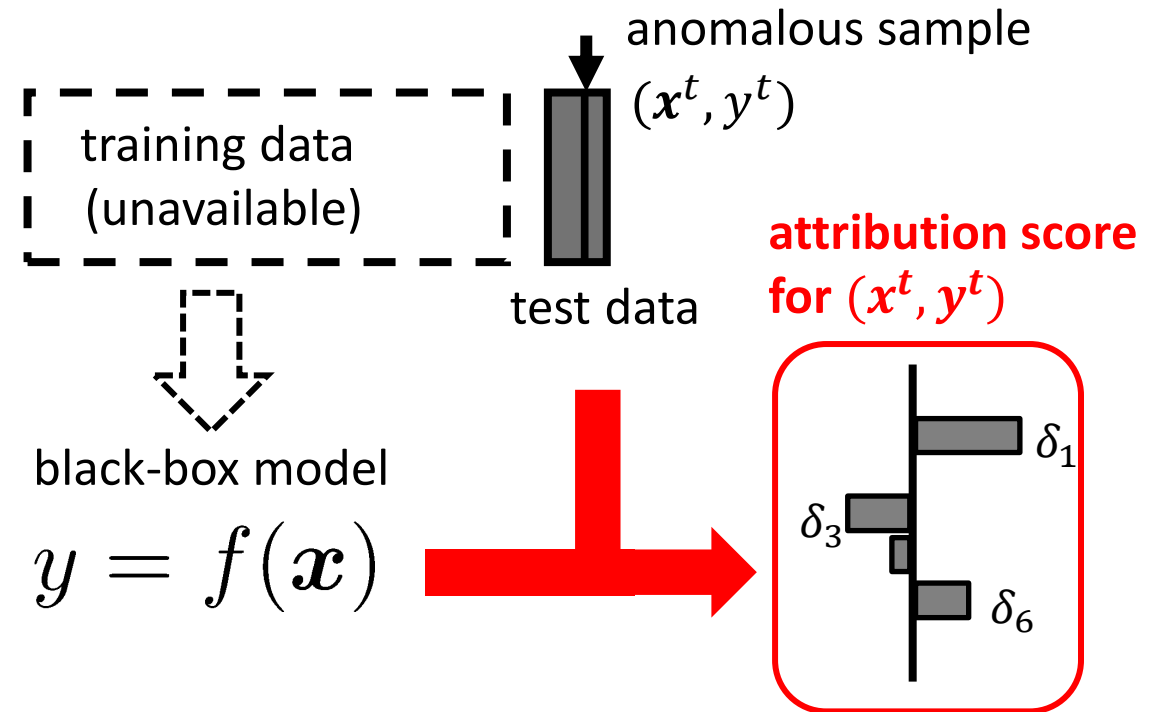
- What is the task, “Anomaly Attribution”?
- What’s wrong with the existing attribution methods?
- What is the new idea?
- Illustrative examples

LIME, Shapley values (SV), and integrated gradient (IG) are three major existing black-box attribution methods.

- LIME, SV, IG have the same in/output
 - In: black-box $y = f(\mathbf{x})$ and test sample.
 - Out: attribution score for each variable
- Why bother to develop a new method?

They are deviation-agnostic.

They can't compute score's uncertainty



LIME, SV, and IG are to explain a black-box function itself locally.

- LIME = local gradient at \mathbf{x}^t
 - Gradient is numerically estimated via sampling.

- IG = increment from a reference point \mathbf{x}^0

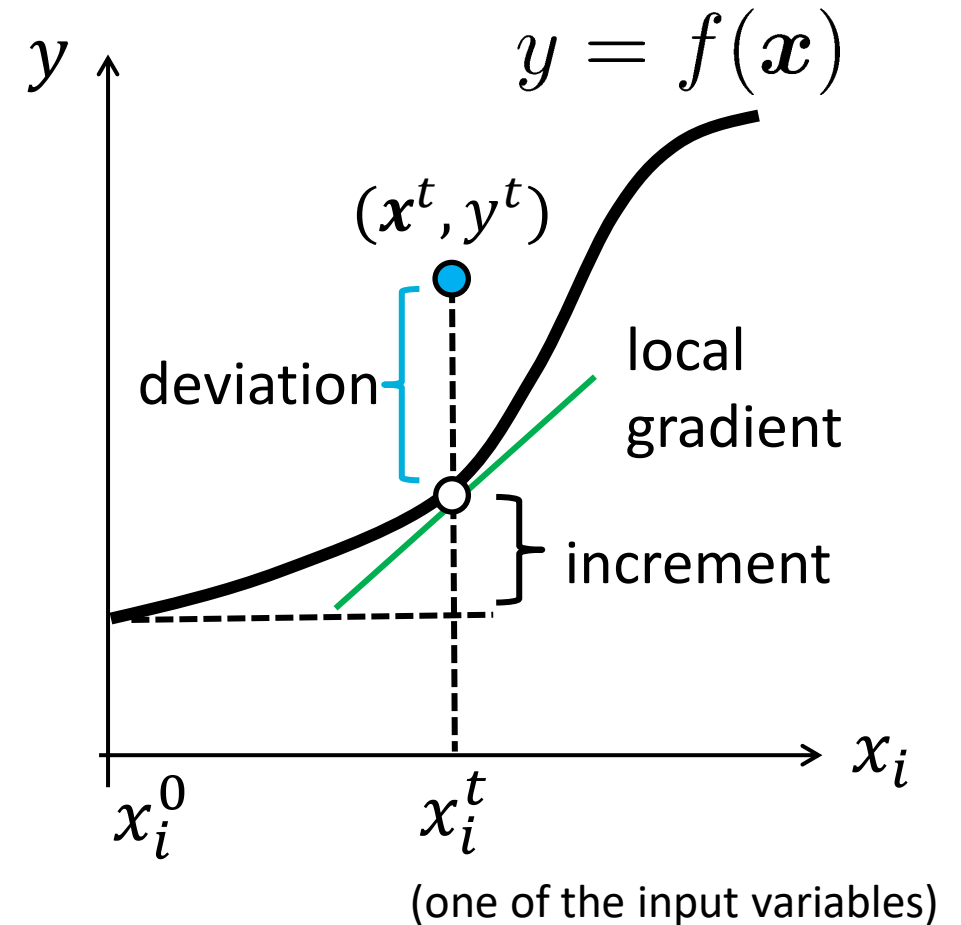
$$\text{IG}_i(\mathbf{x}^t | \mathbf{x}^0) \triangleq (x_i^t - x_i^0) \int_0^1 d\alpha \left. \frac{\partial f}{\partial x_i} \right|_{\mathbf{x}^0 + (\mathbf{x}^t - \mathbf{x}^0)\alpha}$$

- EIG = expected IG

- Computed by marginalizing \mathbf{x}^0

- SV = (something mysterious)

$$\text{SV}_i(\mathbf{x}^t) = \frac{1}{M} \sum_{k=0}^{M-1} \binom{M-1}{k}^{-1} \sum_{\mathcal{S}_i: |\mathcal{S}_i|=k} [\langle f | \mathbf{x}_i^t, \mathbf{x}_{\mathcal{S}_i}^t \rangle - \langle f | \mathbf{x}_{\mathcal{S}_i}^t \rangle]$$



Can they be used to explain deviations? – No.

Summary of theoretical results.

- Result 1: LIME, SV, IG, and EIG are deviation-agnostic
 - This is obvious from the original definition.
 - ✓ They explain $f(\mathbf{x})$ locally at $\mathbf{x} = \mathbf{x}^t$, independently y .
 - The conclusion still holds even when the target function is $f(\mathbf{x}) - y$ rather than $f(\mathbf{x})$.
- Result 2: SV is equivalent to EIG up to the second order of power expansion.

$$SV_i(\mathbf{x}^t, y^t) \approx \text{EIG}_i(\mathbf{x}^t, y^t)$$

- Result 3: LIME is equivalent to the derivative of IG and EIG

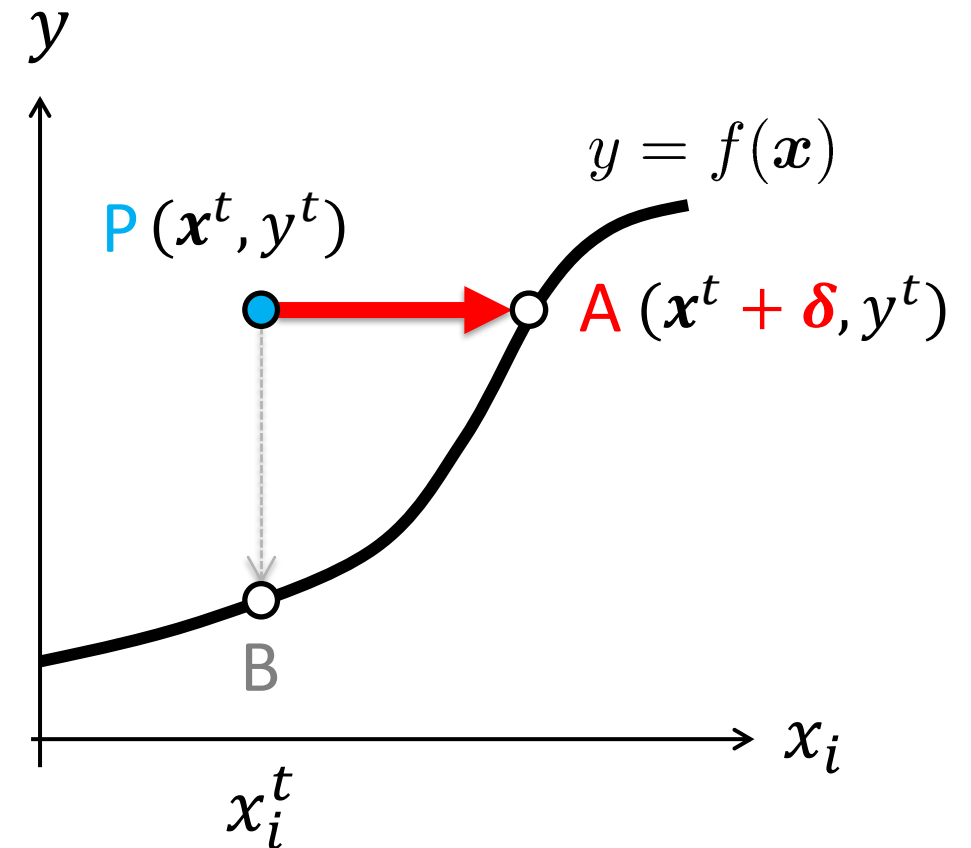
$$\text{LIME}_i(\mathbf{x}^t, y^t) = \frac{\partial \text{EIG}_i(\mathbf{x}^t, y^t)}{\partial x_i}$$

Agenda

- What is the task, “Anomaly Attribution”?
- What’s wrong with the existing attribution methods?
- What is the new idea?
- Illustrative examples

Given a test point (x^t, y^t) being anomalous, we ask: How much “work” would we need to bring it to the normalcy?

- The “work” required for each variable should be a natural attribution score.
- The outlier **P** wouldn't have been anomalous if it were at **A**.
- Hence, the amount of shift, δ , can be viewed as the “work,” indicating the responsibility of each variable.
- How about B? We need a help of $p(y | \mathbf{x})$.

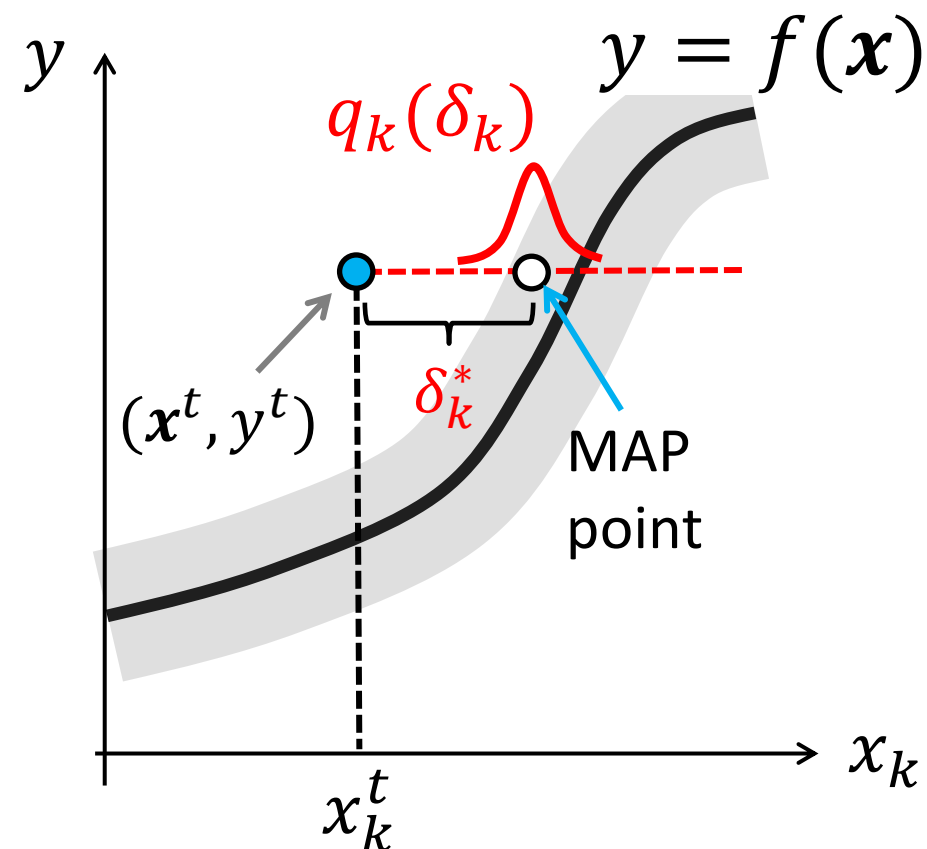


Perturbation as explanation:

Our goal is to find the posterior distribution of δ .

- We need a generative model to handle the ambiguity in prediction.
 - The on-the-curve points may not represent normalcy.
- Generative process with δ as model parameter.
 - $p(y | \mathbf{x}, \delta, \lambda) = \mathcal{N}(y | f(\mathbf{x} + \delta), \lambda^{-1})$
 - priors (η, a_0, b_0 are hyperparameters):
 - ✓ $p(\delta) = \mathcal{N}(\delta | \mathbf{0}, \eta \mathbf{I})$
 - ✓ $p(\lambda) = \text{Gam}(\lambda | a_0, b_0)$
- Formal solution of the posterior (typically $N_{\text{test}} = 1$)

$$Q(\delta) \propto p(\delta) \prod_{t=1}^{N_{\text{test}}} \int_0^{\infty} d\lambda p(y^t | \mathbf{x}^t, \delta, \lambda) p(\lambda)$$



Using variational Bayesian approach combined with a mean-field-like approximation to get variable-wise posteriors.

- Formal solution of the posterior (typically $N_{\text{test}} = 1$)

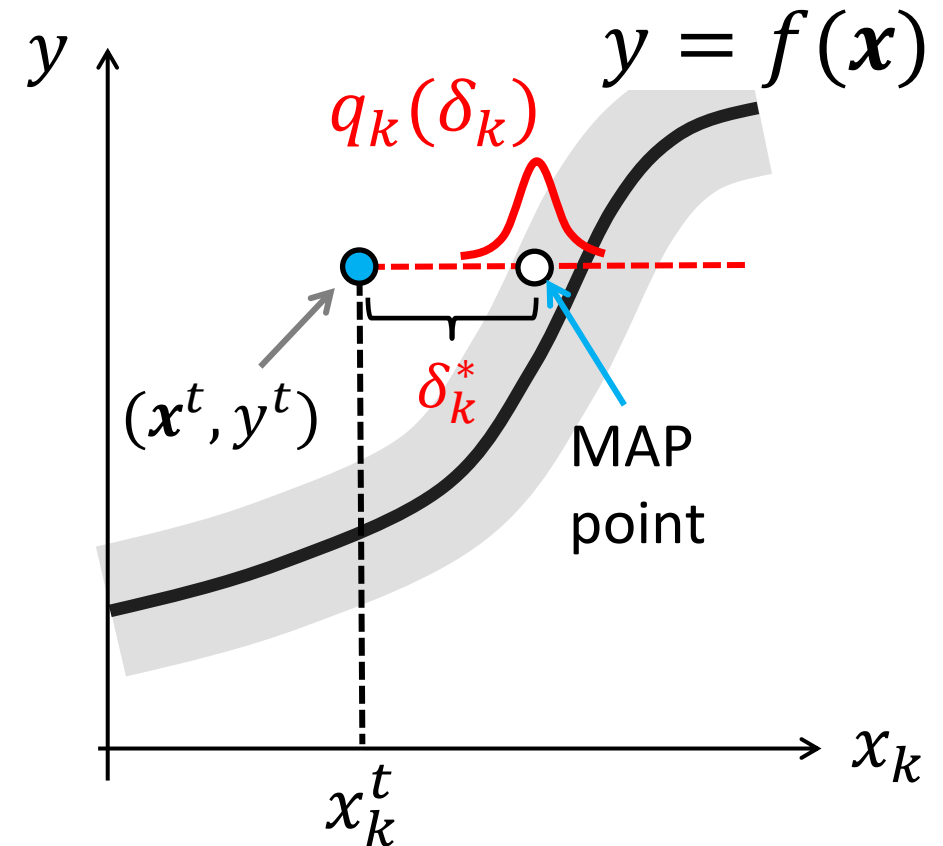
$$Q(\boldsymbol{\delta}) \propto p(\boldsymbol{\delta}) \prod_{t=1}^{N_{\text{test}}} \int_0^\infty d\lambda p(y^t | \mathbf{x}^t, \boldsymbol{\delta}, \lambda) p(\lambda),$$
$$\propto p(\boldsymbol{\delta}) \prod_{t=1}^{N_{\text{test}}} \frac{1}{\sqrt{b_0}} \left\{ 1 + \frac{[y^t - f(\mathbf{x}^t + \boldsymbol{\delta})]^2}{2b_0} \right\}^{-(a_0 + \frac{1}{2})},$$

- How do we get a variable-wise distribution?

- We find an approximated solution by minimizing the KL divergence between $Q(\boldsymbol{\delta})$ and a factorized form:

$$Q(\boldsymbol{\delta}) = Q(\delta_1, \dots, \delta_M) \approx \prod_{k=1}^M q_k(\delta_k),$$

- We also use a mean-field-like approximation to get an explicit form of $\{q_k(\delta_k)\}$. → paper



(For ref.) How the GPA algorithm works

- GPA algorithm has two parts:
 - MAP (maximum a posteriori) estimation
 - Distribution estimation
- MAP estimation solves:

$$\min_{\delta} \left\{ \frac{\eta}{2} \|\delta\|_2^2 + \ln \left\{ 1 + \frac{[y^t - f(\mathbf{x}^t + \delta)]^2}{2b(\mathbf{x}^t)} \right\}^{\frac{2a_0+1}{2}} \right\}$$
 - Use proximal gradient (with ℓ_1 regularizer)
 - The gradient is estimated via local sampling (like LIME)
- Distribution estimation uses a mean-field approximation
 - “Think of the others fixed to the MAP value and focus on yourself.”

Algorithm 2 Generative Perturbation Analysis

Require: $f(x)$, $\mathcal{D}_{\text{test}}$, parameters $\eta, \nu, \kappa, a_0, \{b(\mathbf{x}^t)\}$.

1: randomly initialize $\delta \approx \mathbf{0}$.

```

2: repeat MAP
3:   set  $\mathbf{g} = \mathbf{0}$ 
4:   for all  $(y^t, \mathbf{x}^t) \in \mathcal{D}_{\text{test}}$  do
5:     Compute the local gradient  $\frac{\partial f(\mathbf{x}^t + \delta)}{\partial \delta}$ 
6:     Update  $\mathbf{g} \leftarrow \mathbf{g} + \frac{\partial f(\mathbf{x}^t + \delta)}{\partial \delta} \frac{y^t - f(\mathbf{x}^t + \delta)}{2b(\mathbf{x}^t) + [y^t - f(\mathbf{x}^t + \delta)]^2}$ 
7:   end for
8:    $\mathbf{g} \leftarrow (1 - \kappa\eta)\delta + \kappa(2a_0 + 1)\mathbf{g}$ 
9:    $\delta = \text{sign}(\mathbf{g}) \max\{0, |\mathbf{g}| - \eta\nu\}$ 
10: until convergence
  
```

11: set $\delta^* = \delta$

```

12: for all  $k$  do distribution
13:    $q_k(\delta) = Q(\delta_1^*, \dots, \delta_{k-1}^*, \delta, \delta_{k+1}^*, \dots, \delta_M^*)$ 
14:    $q_k(\cdot) \leftarrow q_k(\cdot) / \int d\delta' q_k(\delta')$  with Eq. (18)
15: end for
  
```

16: **return** $\{q_k(\cdot) \mid k = 1, \dots, M\}$ and δ^*

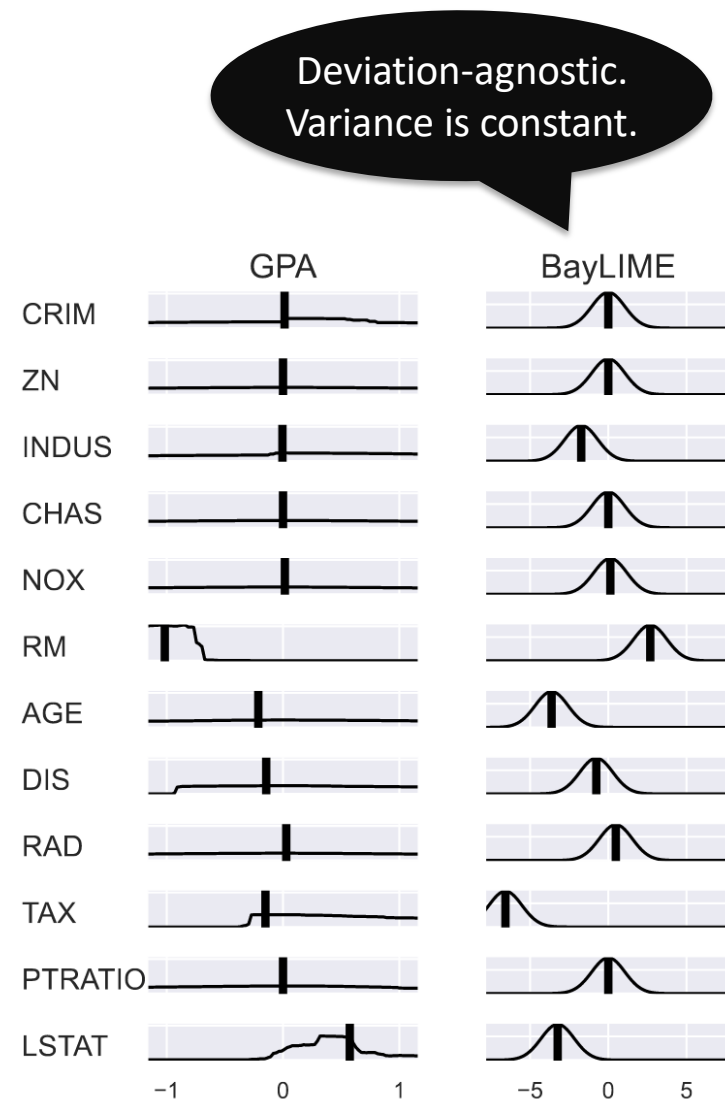
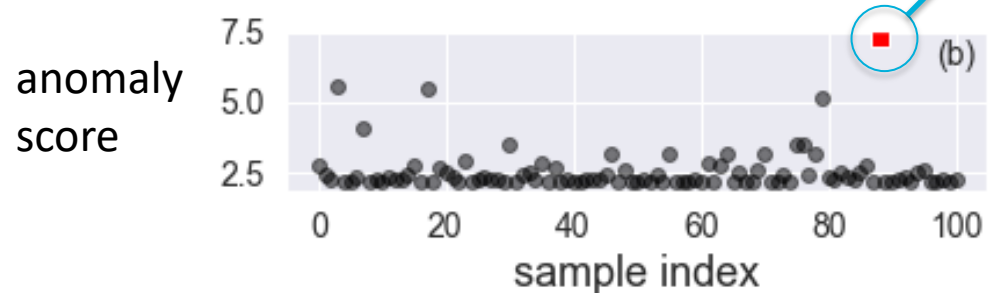
Agenda

- What is the task, “Anomaly Attribution”?
- What’s wrong with the existing attribution methods?
- What is the new idea?
- Illustrative examples

“Why does this house look so unusual?”

House hunting use-case

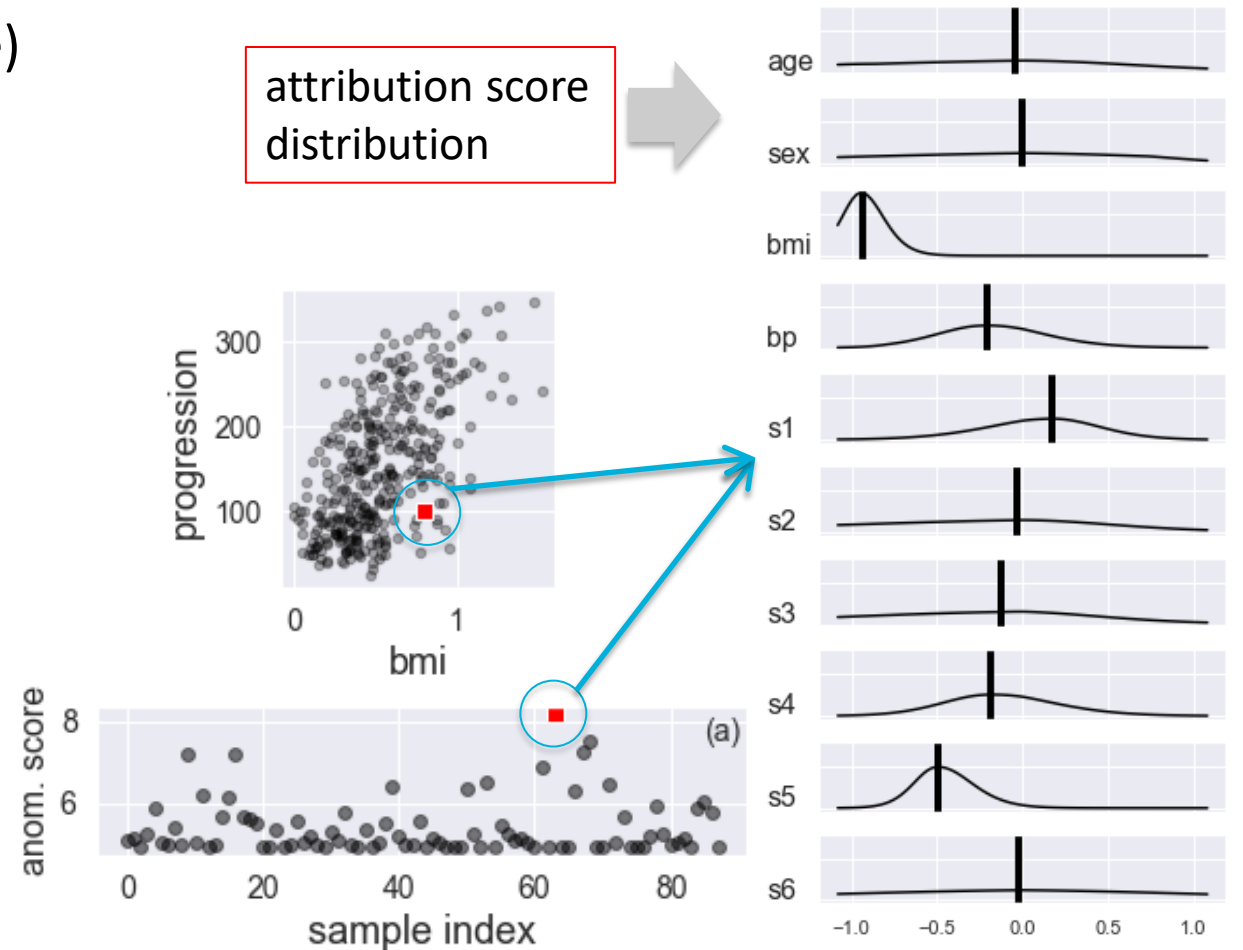
- Boston Housing data
 - y : house price
 - x : house age, # rooms, neighborhood crime rate, etc.
- Computed attribution scores for the top outlier.
 - GPA is able to provide variable-specific distributions in contrast to BayLIME
- Is it a bargain? Probably yes.
 - It's got unusually larger #rooms (RM) and lower poor neighbors (LSTAT) than the peers in the same price range.



“Why does this patient look so unusual?”

Healthcare use-case

- Diabetes data
 - y : diabetes' progression (numerical score)
 - x : biomarkers (BMI, blood pressure, etc.)
- Computed attribution score for the top outlier (patient # 63).
 - Found a large negative score in BMI
 - ✓ The high and narrow pdf translates to high confidence
 - For his progression level, he would look like a regular patient if BMI were much smaller:
 - ✓ “He is overweight but healthy (low progression)” or “He is healthy despite overweight”



Summary

- GPA is the first black-box attribution framework allowing probabilistic attribution.
- We have showed a strong impossibility result: LIME, SV, and IG are deviation-agnostic, and hence, not suitable for anomaly attribution.
- We have also uncovered a relationship between LIME, SV, and IG for the first time.

	model-agnostic	training-data-free	baseline-input-free	y -sensitive	built-in UQ
LIME [34]	yes	yes	yes	no	yes/no
SV [43,44]	yes	no	yes	no	no
IG [39,46]	yes	yes	no	no	no
EIG [6]	yes	no	yes	no	no
Z-score [5]	yes	no	yes	no	no
LC [20]	yes	yes	yes	yes	no
GPA	yes	yes	yes	yes	yes

Thank you!