

Trustworthy AI

Fadwa Idlahcen

Plan

- Motivation

- What's AI trustworthiness?

- Robustness, generalization, explainability, transparency, reproducibility, fairness and privacy preservation : definitions, examples and widely used AI models

- AI lifecycle and where trustworthiness should be checked for

- Challenges and conclusion

Motivation

- Accuracy is the only metric for AI models
- Biased, lack in user privacy protection, vulnerable to attacks...



Insight - Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

October 11, 2018 2:50 AM GMT+2 · Updated 5 years ago

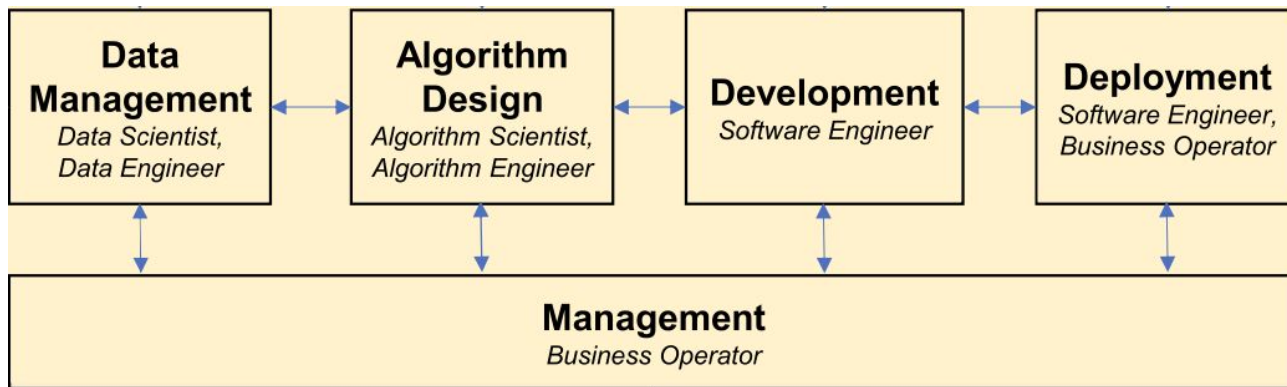


Trustworthy AI

Practitioners need to address AI trustworthiness, including :

- Robustness, generalization, explainability, transparency, reproducibility, fairness and privacy preservation

Optimizing trustworthiness throughout AI lifecycle, rather than at each step.

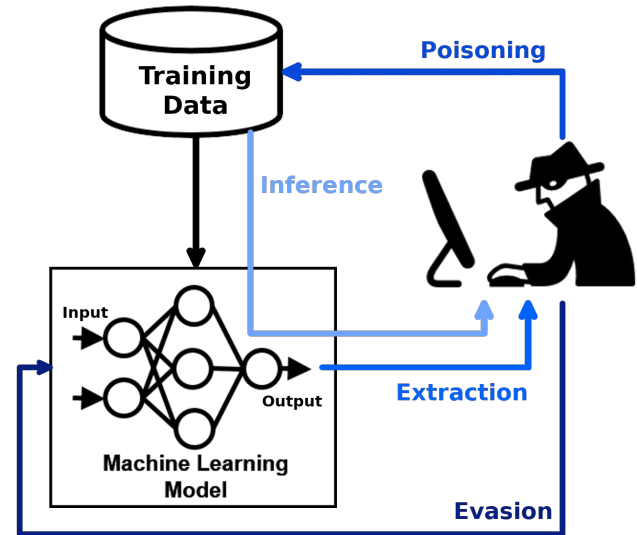


Robustness

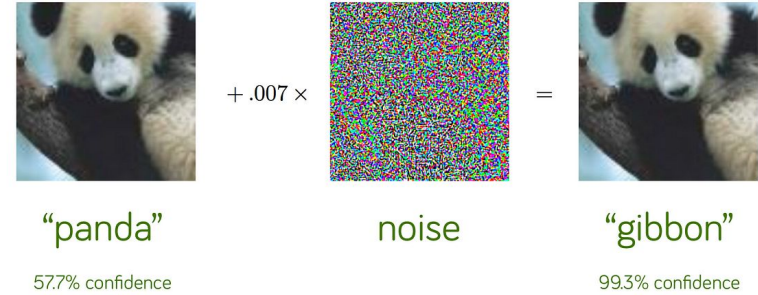
It's not just about reliability against errors, erroneous inputs and unseen data but also safety.

3 levels :

- Data: when an AI model is trained without considering the diverse distributions of data in different scenarios.
- Algorithms : adversarial attacks
- Systems : the use of illegal inputs



Adversarial attacks



-A carefully computed example to be misclassified

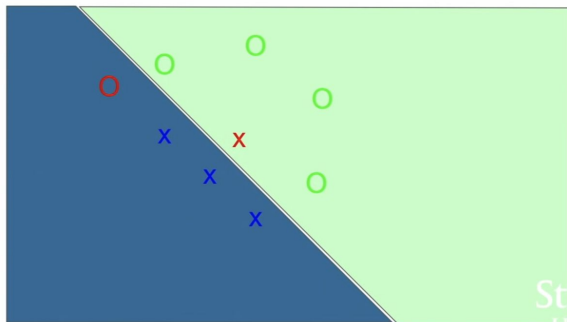
-Some use cases :

- Image classification: misclassification. E.g: autonomous vehicles, security systems, and medical imaging.
- Text generation: generating deceptive or misleading text that can manipulate sentiment analysis algorithms, spam filters, or automated content moderation systems.
- Malware evasion: designing malware that evades detection by antivirus or intrusion detection systems.

Adversarial attacks

Piecewise linearity of NN \implies Fast gradient sign method

Adversarial Examples from
Excessive Linearity



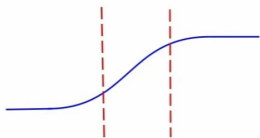
Rectified linear unit



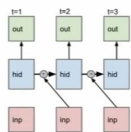
Maxout



Carefully tuned sigmoid

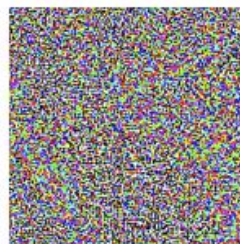


LSTM



\mathbf{x}
“panda”
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$
“nematode”
8.2% confidence

=



$\mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$
“gibbon”
99.3 % confidence

Robustness

-Tests:

- Test datasets with various distributions are used
- Deriving a lower bound of the minimum distortion to an attack on an AI model

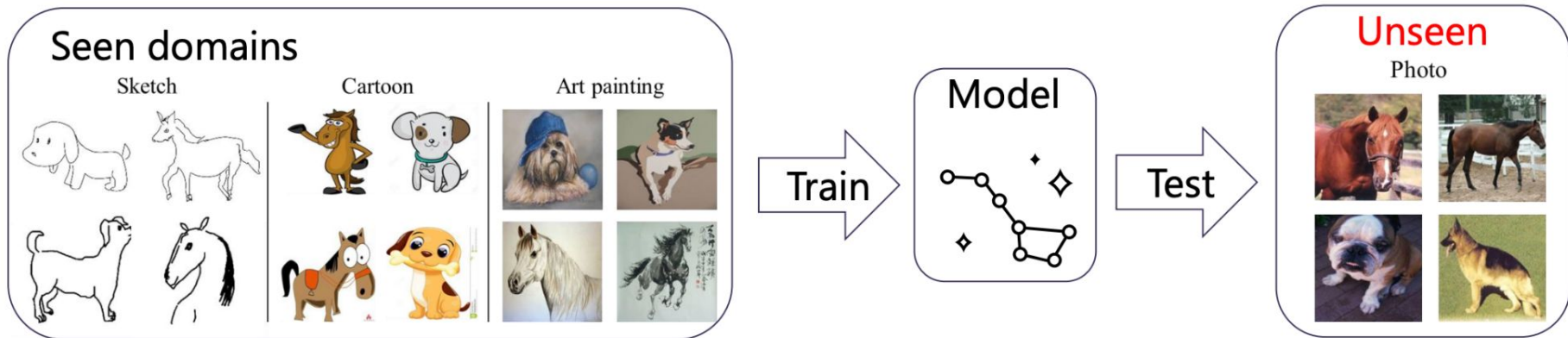
-Defense:

- defensive distillation: removing the gradient from the model to be protected;
- adversarial training (prone to overfitting)

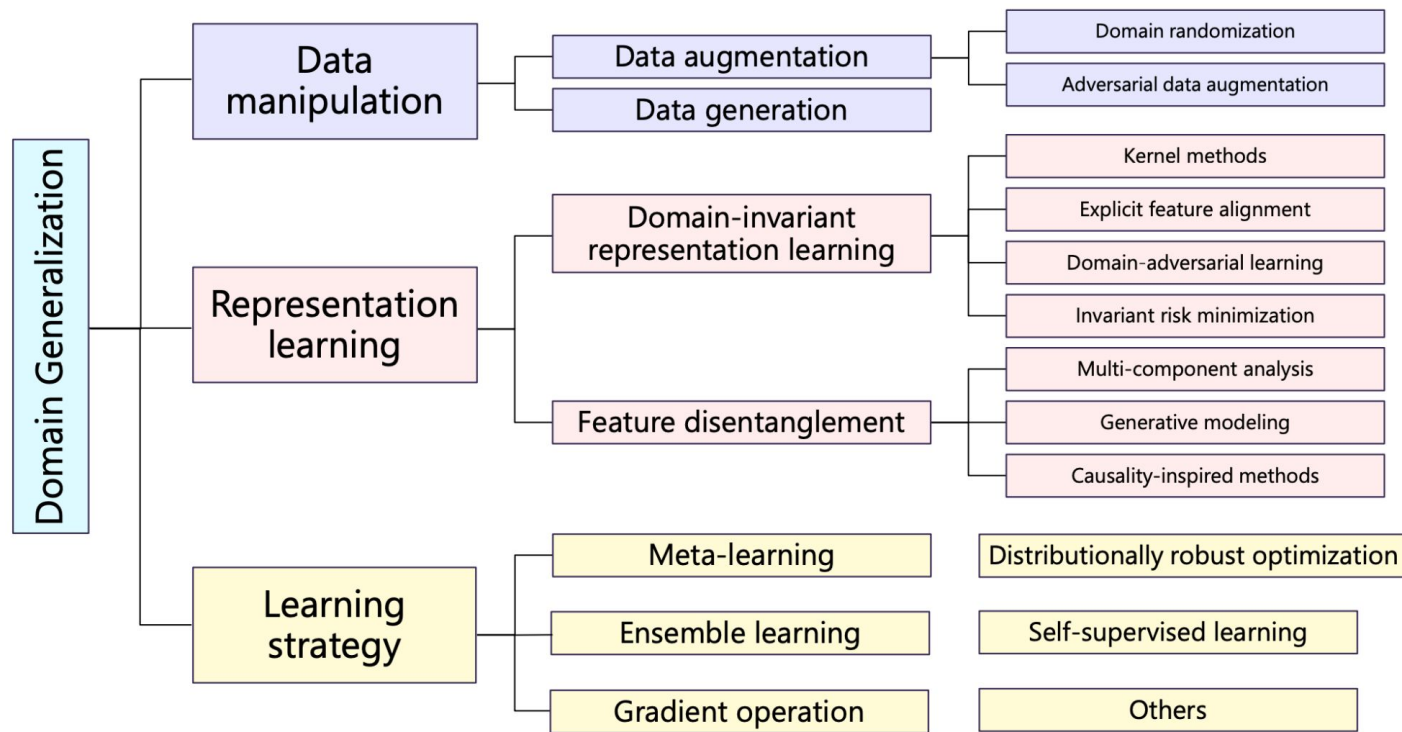
⇒ Not effective!!!

Generalization

- From training limited data to accurate predictions on unseen data
- Robustness against distributional shifts = generalization problem
 - Double effect : an algorithm that is robust against small perturbations has better generalization \neq adversarial training may reduce the testing accuracy

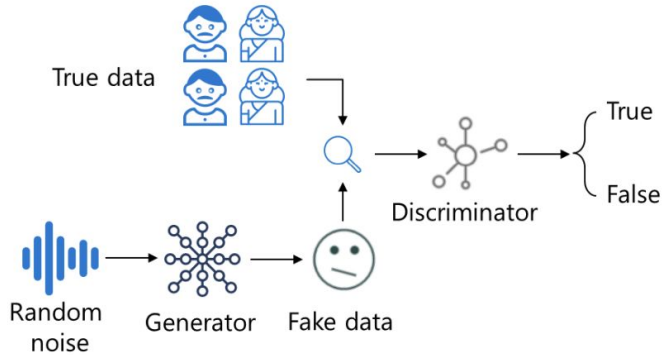


Domain generalization



Data manipulation

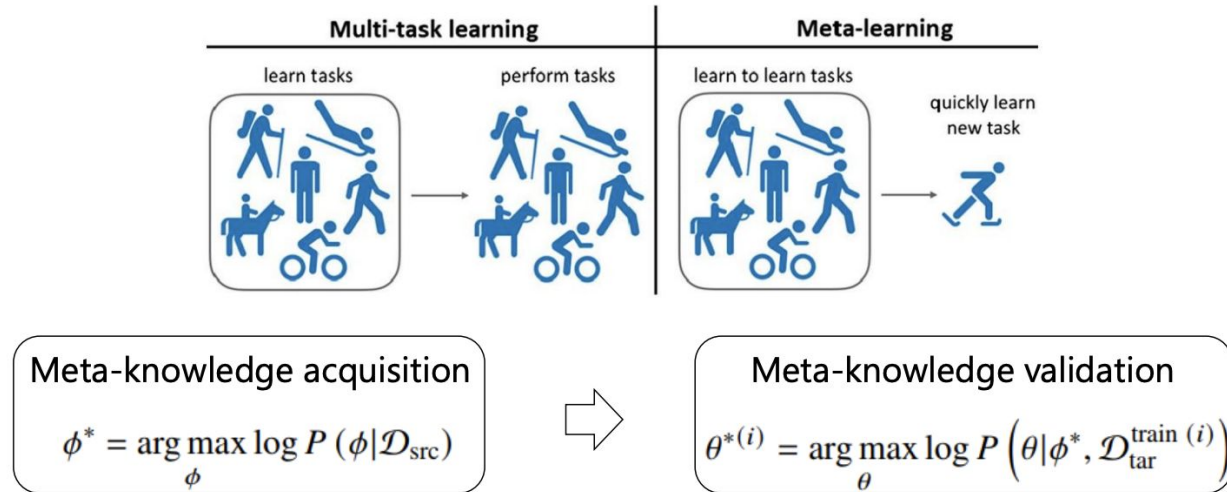
- Insufficient data in the target domain
- Adversarial data augmentation : via gradient training
- Data generation : GANs



Learning strategy

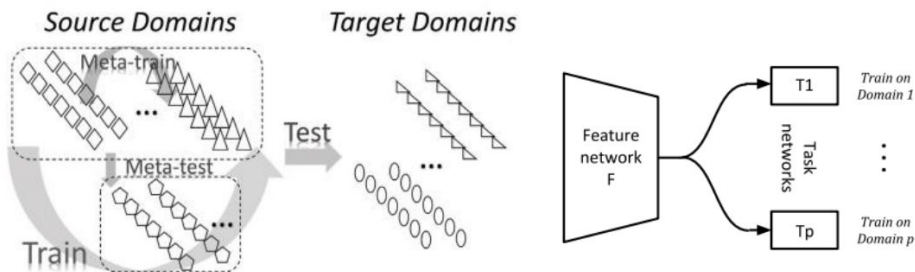
-Meta learning : Divide domains into several tasks, then use meta-learning to learn general knowledge

- Learning to learn, or meta-learn the general knowledge
 - Instead of the original tasks, meta-learning wants to acquire knowledge about **new tasks**



Learning strategy

- How to adopt meta-learning for DG?
 - Key: Old tasks to new tasks in meta-learning → Old domains to new domains
- MLDG: Meta-learning for DG
- MetaReg: meta-learning for regularization



Algorithm 1 Meta-Learning Domain Generalization

```

1: procedure MLDG
2:   Input: Domains  $\mathcal{S}$ 
3:   Init: Model parameters  $\Theta$ . Hyperparameters  $\alpha, \beta, \gamma$ .
4:   for ite in iterations do
5:     Split:  $\bar{\mathcal{S}}$  and  $\check{\mathcal{S}} \leftarrow \mathcal{S}$ 
6:     Meta-train: Gradients  $\nabla_{\Theta} = \mathcal{F}'_{\Theta}(\bar{\mathcal{S}}; \Theta)$ 
7:     Updated parameters  $\Theta' = \Theta - \alpha \nabla_{\Theta}$ 
8:     Meta-test: Loss is  $\mathcal{G}(\check{\mathcal{S}}; \Theta')$ .
9:     Meta-optimization: Update  $\Theta$ 

$$\Theta = \Theta - \gamma \frac{\partial(\mathcal{F}(\bar{\mathcal{S}}; \Theta) + \beta \mathcal{G}(\check{\mathcal{S}}; \Theta - \alpha \nabla_{\Theta}))}{\partial \Theta}$$

10:   end for
11: end procedure

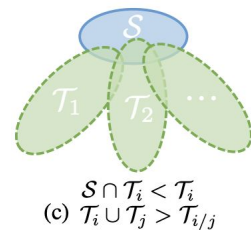
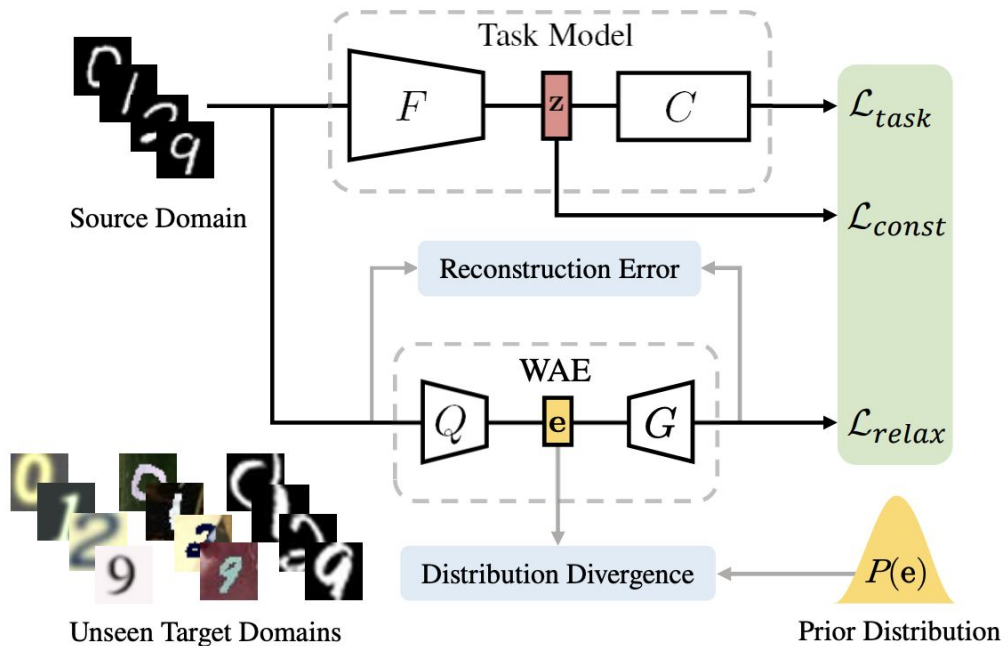
```

- Li D, Yang Y, Song Y Z, et al. Learning to generalize: Meta-learning for domain generalization. AAAI 2018.
- Balaji Y, Sankaranarayanan S, Chellappa R. Metareg: Towards domain generalization using meta-regularization. NeurIPS 2018.

Data manipulation - Meta-learning

-VAE for data generation

F : feature extractor; C : classifier ; z : latent representation of x ; \mathcal{L}_{task} : classification loss; \mathcal{L}_{const} : worst-case loss; \mathcal{L}_{relax} : relaxation loss



\mathcal{S} : Source domain(s)

\mathcal{T} : Target domain(s)

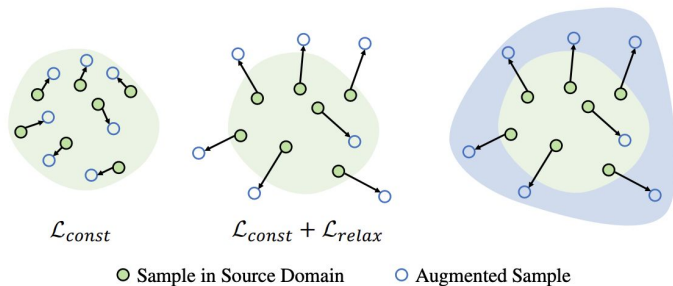
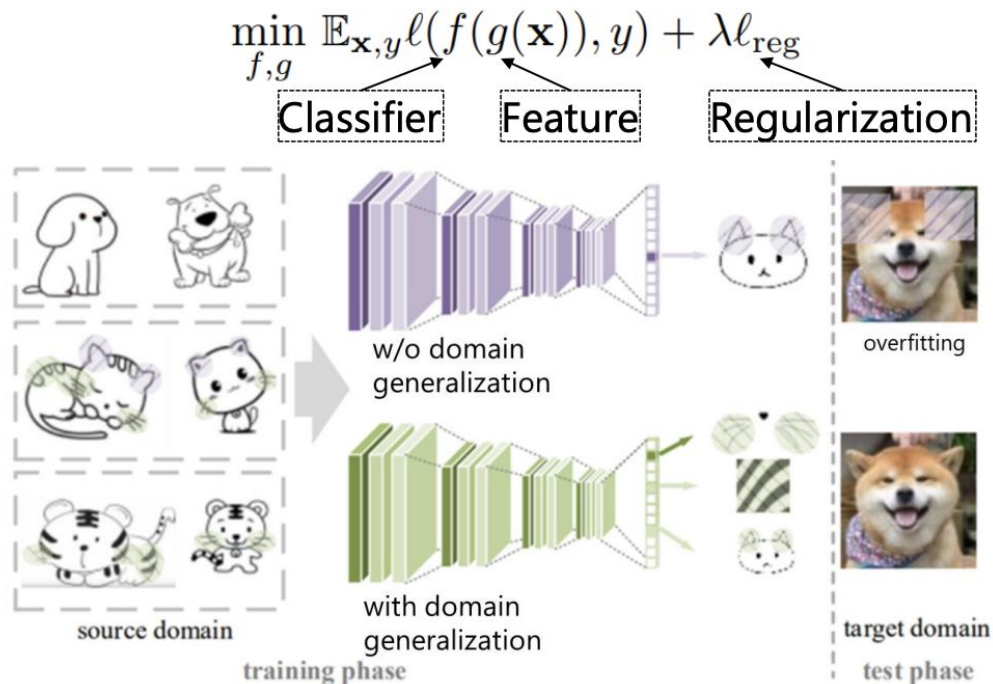


Figure 3. Motivation of \mathcal{L}_{relax} . **Left:** The augmented samples may be close to the source domain if applying \mathcal{L}_{const} . **Middle:** We expect to create out-of-domain augmentations by incorporating \mathcal{L}_{relax} . **Right:** This would yield an enlarged training domain.

Representation learning

-Learning invariant features



Representation learning

How?

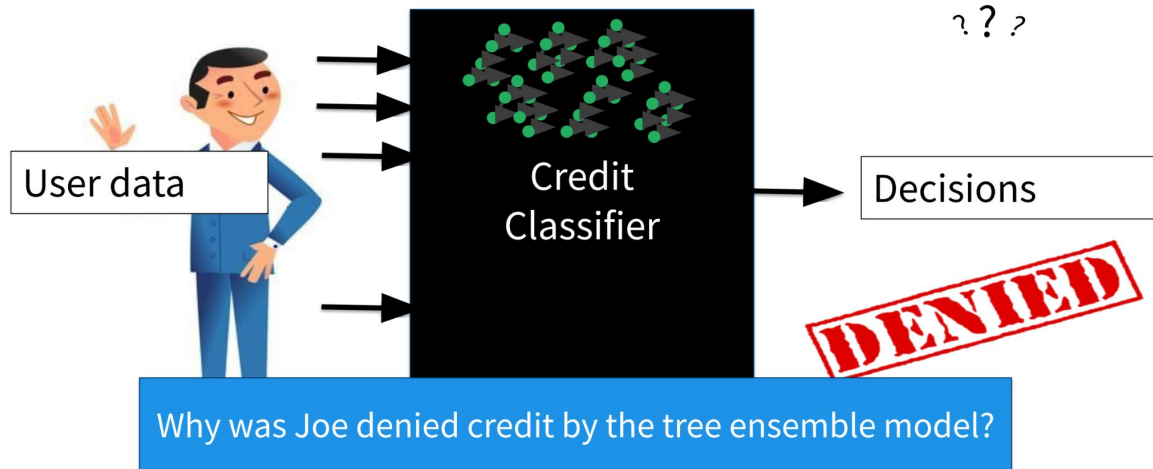
- kernel based models:
 - transfer component analysis
- Domain adversarial learning
- ...

Explainability

-Explainability: understanding how an AI model makes decisions

- Model explainability by design
- Post hoc model explainability

-Explainability vs Interpretability



Explainable ML

Design

-self-explainable models :

KNN, linear/logistic regression, decision trees/rules, and probabilistic graphical models

⇒ Complex structure = unexplainable

-Hybrid combinations of self-explainable models and black-box models

Post hoc

-Explainer approximation

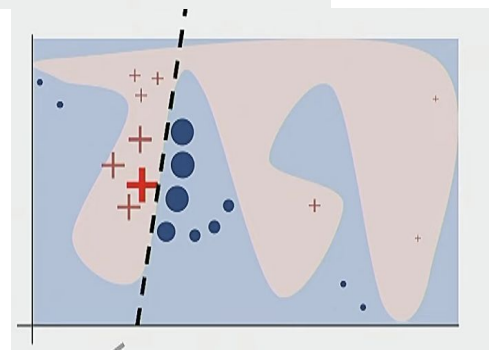
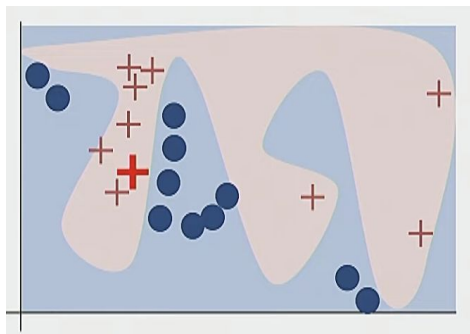
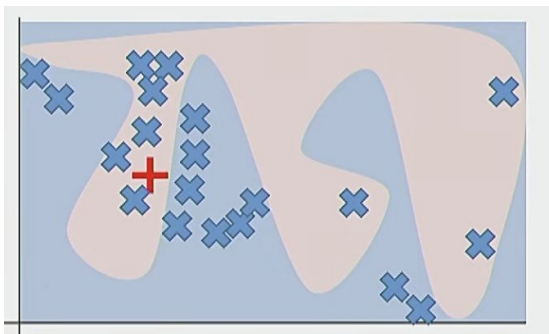
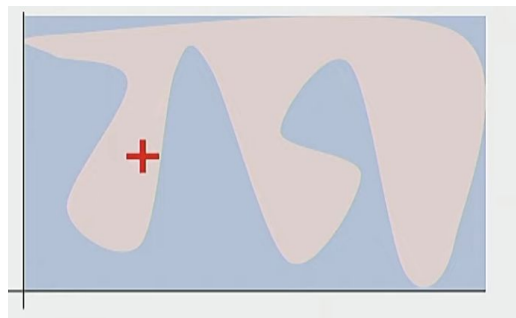
-Feature based explanation

-Two types : local and global

Local explainability

-Feature importance: LIME : explaining individual predictions at a time

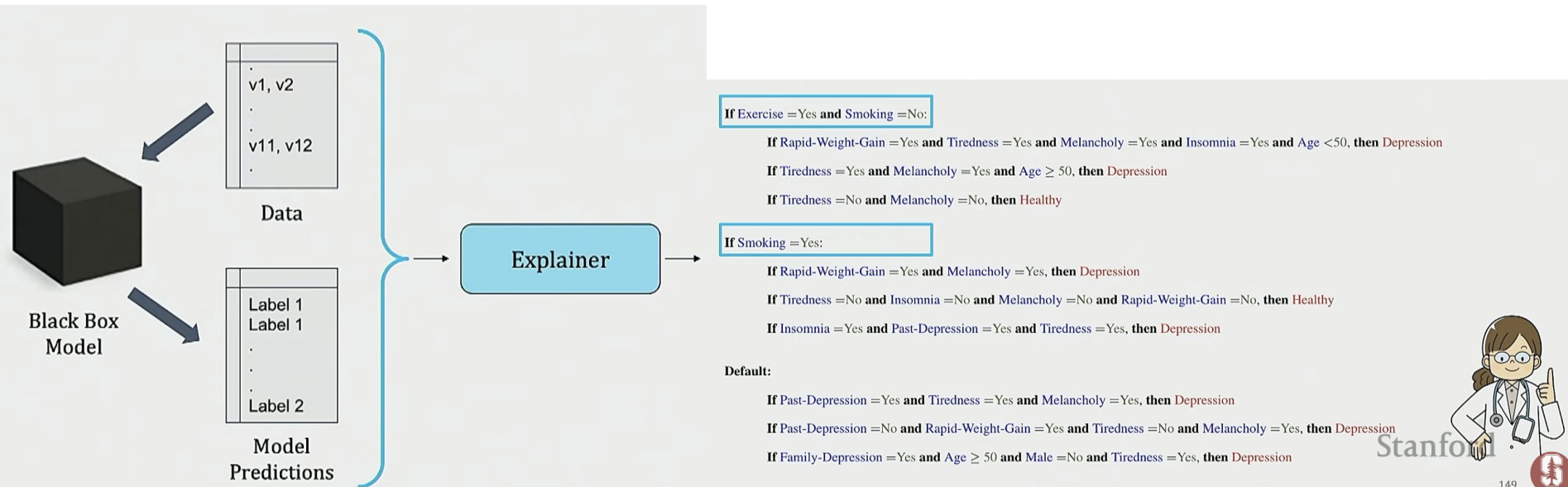
- sample points around x_i
- use model to predict labels for each sample
- weigh samples according to distance to x_i
- learn simple linear model on weighted samples
- use simple linear model to explain



Global explainability

-Combination of local explanations

-Model Distillation



Transparency


-Transparency considers AI as a software system, and seeks to disclose information regarding its entire lifecycle



Documentation



Auditing



Information
Sharing

Reproducibility

- Essential step to verify AI research
- Should be considered over the entire lifecycle (data, methods, and experiments)



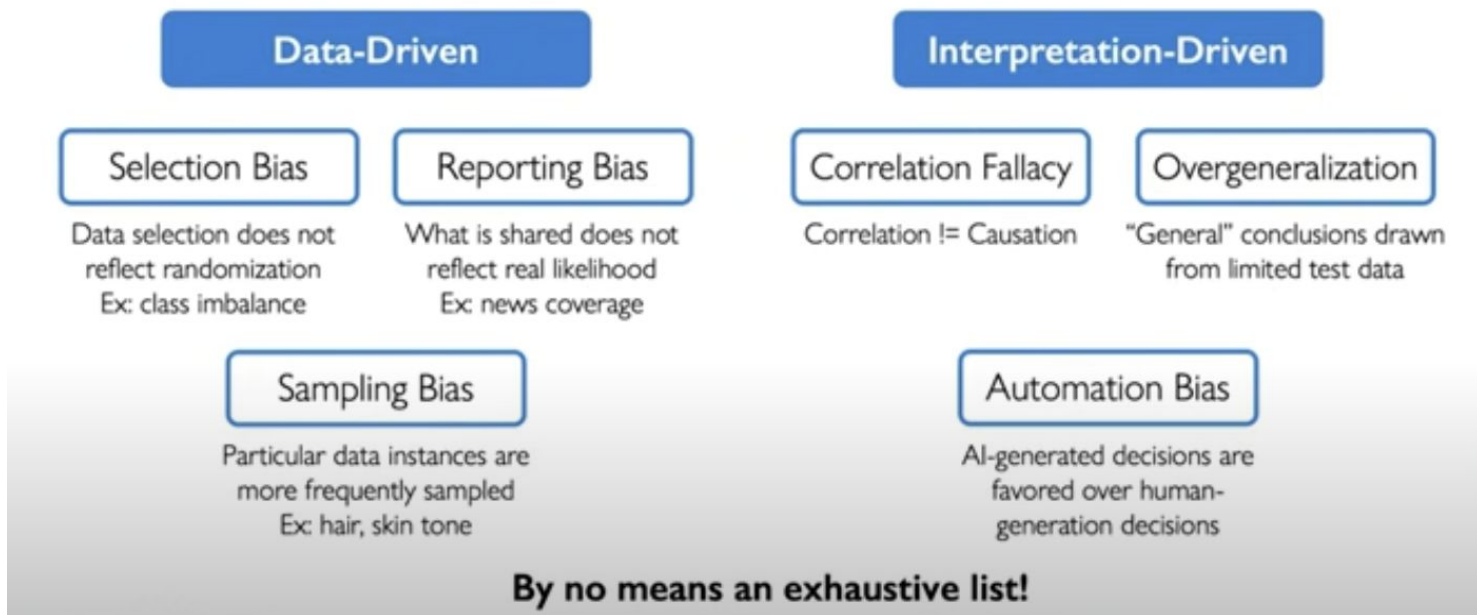
-ChatGPT-3 and reproducibility issues

Fairness

- Hiring
- Financial risk assessment
- Face identification

The image shows a screenshot of a news article from The New York Times. The article title is "Who Is Making Sure the A.I. Machines Aren't Racist?" and the sub-headline is "When Google forced out two well-known artificial intelligence experts, a long-simmering research controversy burst into the open." To the right of the text is a portrait of a woman with dark curly hair and a light-colored scarf. Below the article is a Google Translate interface. The source text is in Turkish: "O bir doktor" and "O bir hemşire". The target text is in English: "He is a doctor" and "She is a nurse". The interface includes language selection dropdowns for English, Spanish, French, Turkish, Arabic, and a "Translate" button. At the bottom of the page, there are links for "About Google Translate", "Community", "Mobile", "G+", "About Google", "Privacy & Terms", "Help", and "Send feedback".

Common Biases



Improving Fairness

-Bias mitigation : removing problematic signal

-Inclusion: adding signal for desired features (counterfactual data augmentation)

Example : Language modelling

Adversarial Multi-Task Learning to Mitigate Bias

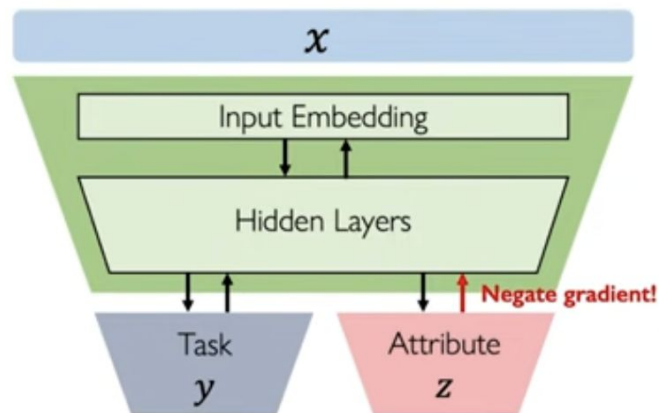
Setup: specify attribute z for which we seek to mitigate bias. Jointly predict output y and z .

Two discriminator output heads:

1. Target / class label y
2. Sensitive attribute z

Train adversarially:

1. Predict sensitive attribute z
2. Negate gradient for z head
3. "Remove" effect of z on task decision

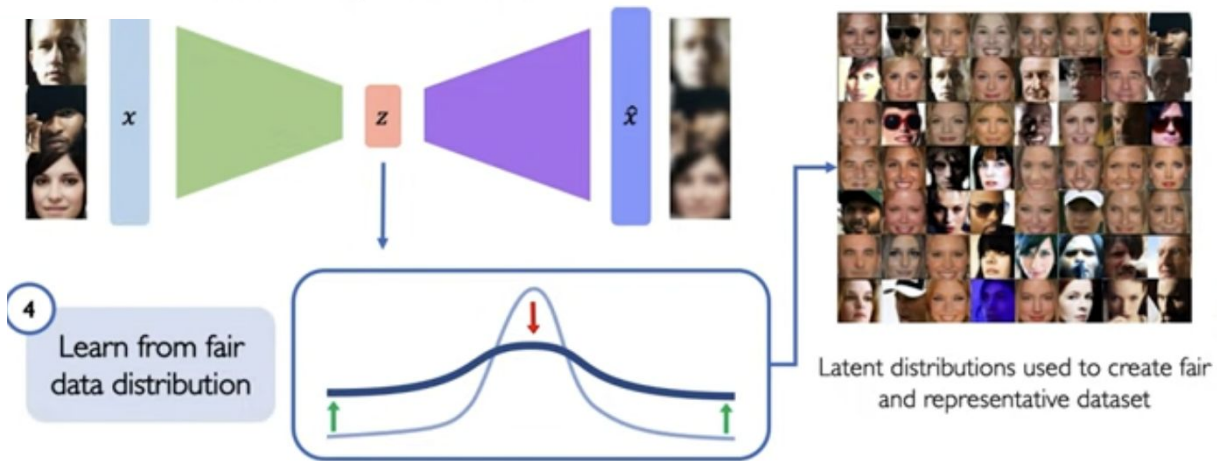
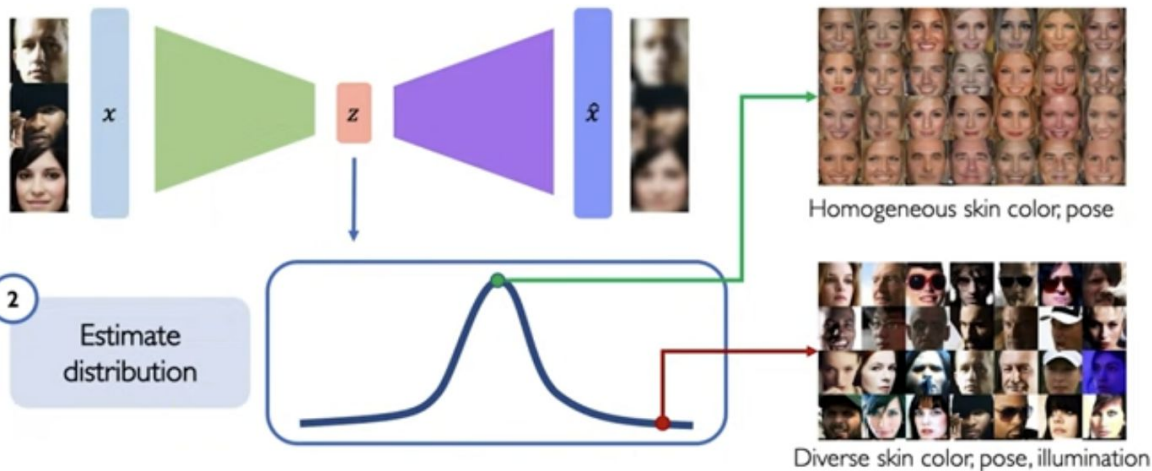


Jointly predict output label y and sensitive attribute z to remove from decision



Hidden features;
counterfactual effect
of adding a bias

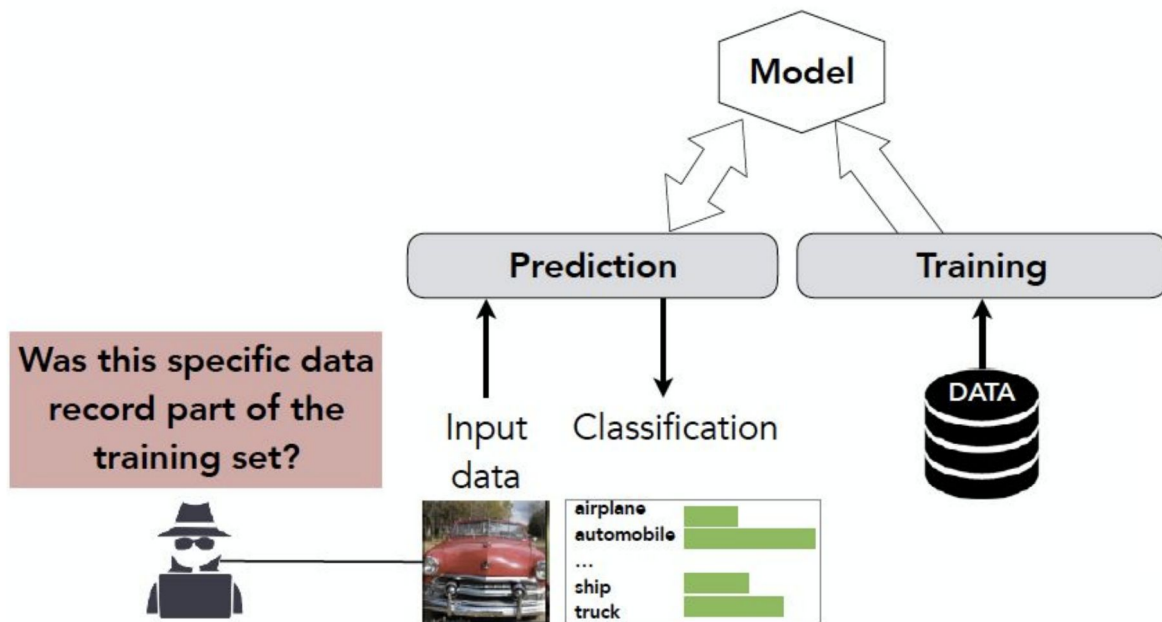
Mitigating Bias through Learned Latent Structure



Privacy Protection

-Protection against information leakage

Membership inference attack



AI Lifecycle

Lifecycle	Approaches
Data Preparation	Data Collection
	Data Preprocessing
Algorithm Design	Adversarial Robustness
	Explainability ML
	Model Generalization
	Algorithmic Fairness
	Privacy Computing

Development	Functional Testing
	Performance Benchmarking
	Simulation
	Formal Verification
Deployment	Anomaly Monitoring
	Human-AI Interaction
	Fail-Safe Mechanism
	Hardware Security
Management	Documentation
	Auditing
	Cooperation
Workflow	MLOps

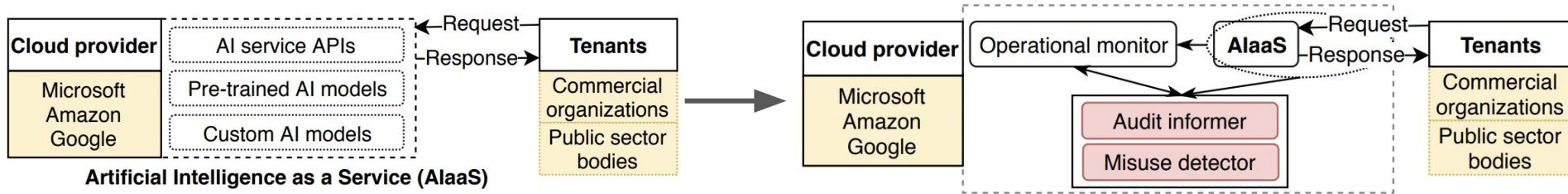
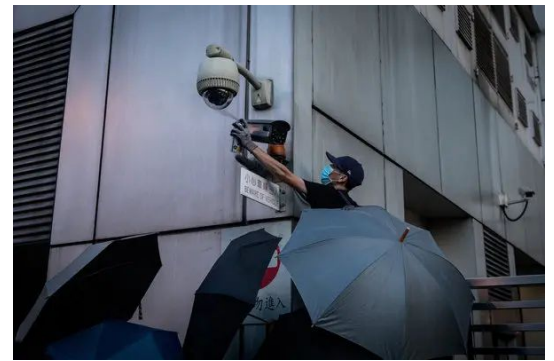
Deployment

-Anomaly monitoring :

police use of facial recognition in Hong Kong

'Deepfakes'

-Microsoft limits the request rate to their Face API, Amazon prevents more than 100 faces from being detected in single image



<i>Technical indicator for vision and face services</i>	<i>Potential implications</i>
High request rate for face detection	Population surveillance
Large number of faces in an image/video	Population surveillance
Large number of different faces are analysed	Population surveillance
Large number of identification attempts for particular individual(s)	Privacy threats to an individual
Detection of 'black-listed' objects	Controversial application

Deployment

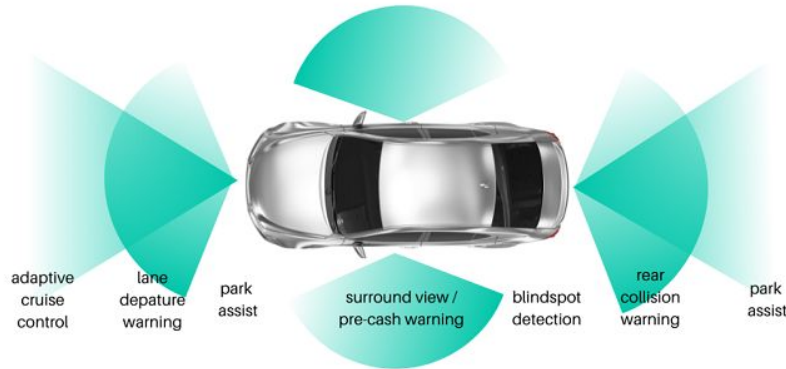
-Human-AI interaction :

User interface: visualization of ML models and interactive parameter-tuning

Human intervention :

- participating in decisions : advanced driver-assistance system (ADAS)
- monitoring failure

-Fail-Safe Mechanisms : if an AI system is causing harm



MLOps

-similar to DevOps

-A start point to build the workflow for trustworthy AI

-Properties:

- Aligned principles of trustworthiness
- Extensive management of artifacts
- Continuous feedback loops
- Close collaboration between interdisciplinary roles

Challenges and Conclusion

- Shift of focus from performance-driven AI to trust-driven AI
- The good and the bad : Large-scale Pre-trained Models
- Limitations in current evaluations of trustworthiness
- End-User Awareness of the Importance of AI Trustworthiness
- Inter-disciplinary and International Cooperation

Thank you.