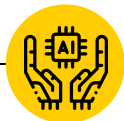


# Quantifying the Trustworthiness Level of Artificial Intelligence Models and Decisions

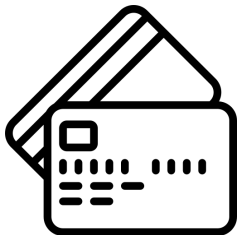


## Agenda

- Motivation
- Taxonomy
- Trusted AI Algorithm
- Validation Scenarios
- Demo



## Motivation



Credit  
Scoring



Medical  
Diagnostics



Applicant  
Screening



Autonomous  
Driving



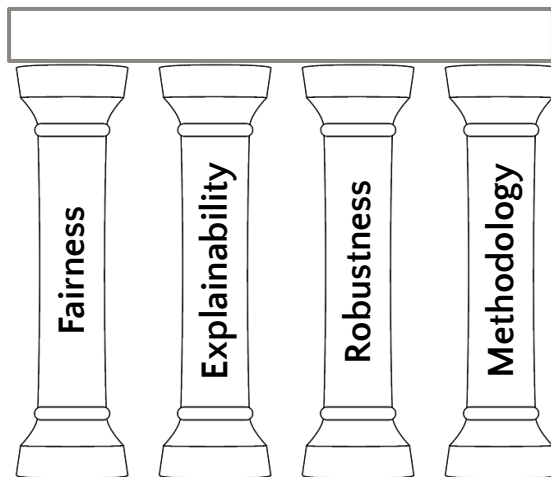
## Thesis Goals

- 1. Identify metrics and key dimensions of trust
- 2. Compile a general taxonomy
- 3. Develop an algorithm to automatically compute the trust score of machine learning models
- 4. Prototypical implementation
- 5. Validation of the proposed solution



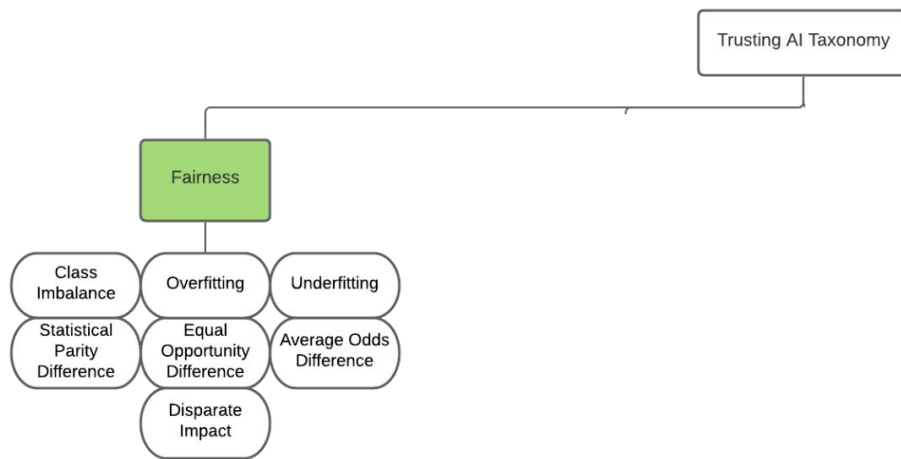
## Pillars of Trust

# Trust in AI



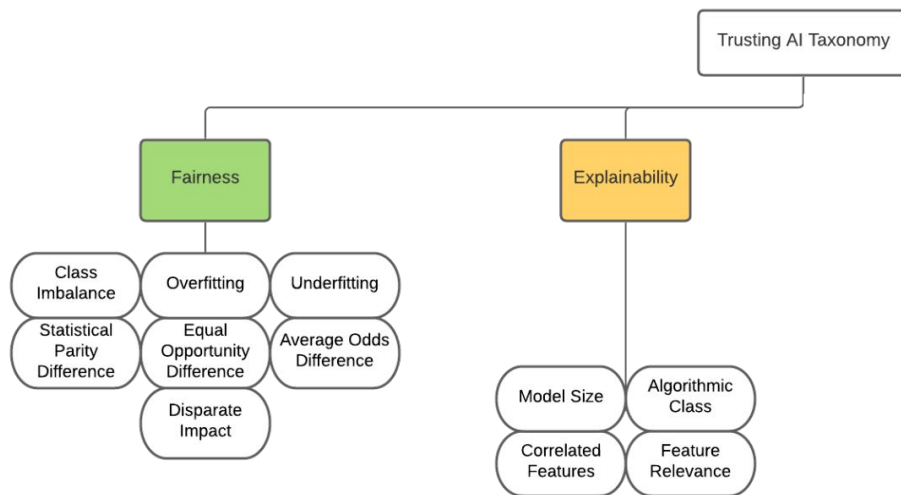


# Taxonomy - Trusted AI Metrics



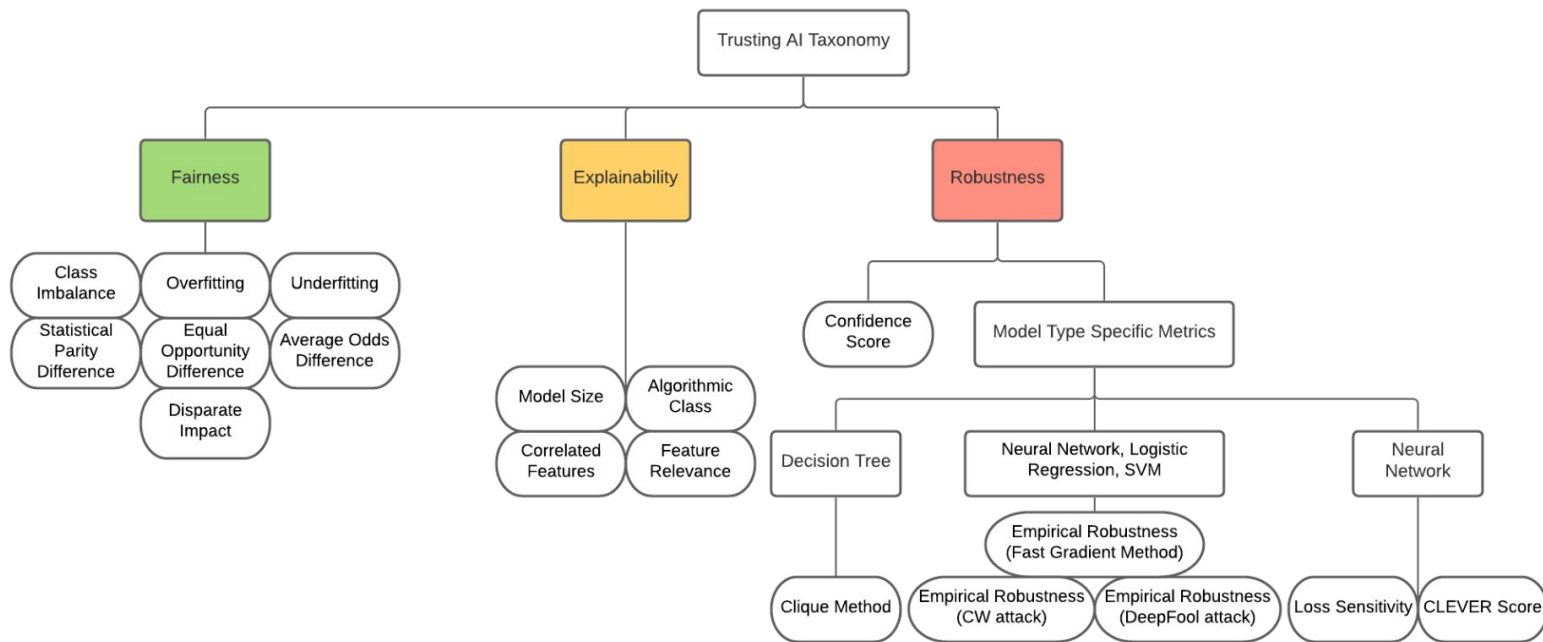


# Taxonomy - Trusted AI Metrics



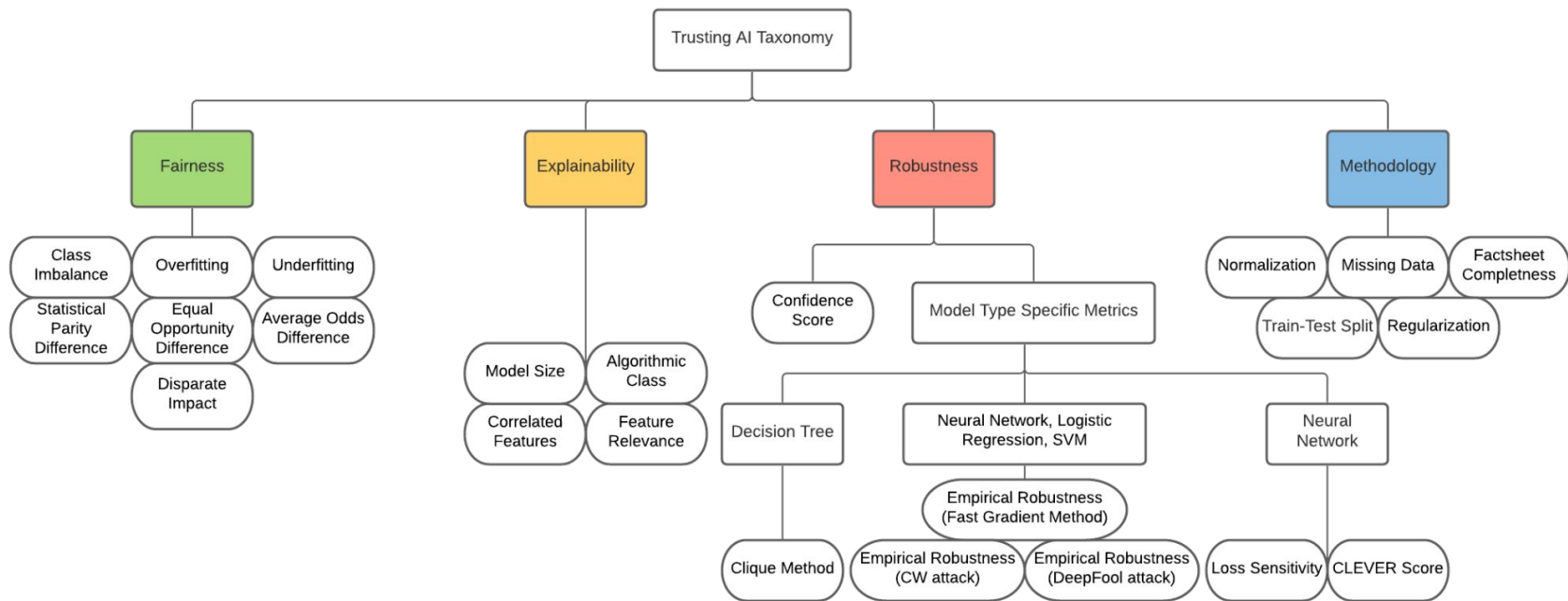


# Taxonomy - Trusted AI Metrics





# Taxonomy - Trusted AI Metrics







## **Metric Dependencies**

---

- ML Model
- Training Data
- Test Data
- Factsheet



# Trusted AI Algorithm Design

Pillar Specific Metrics

$\text{metric}_{1,E}$

$\text{metric}_{2,E}$

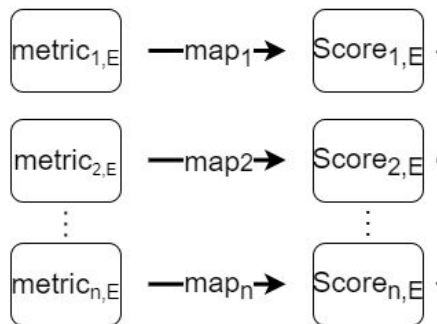
$\vdots$

$\text{metric}_{n,E}$



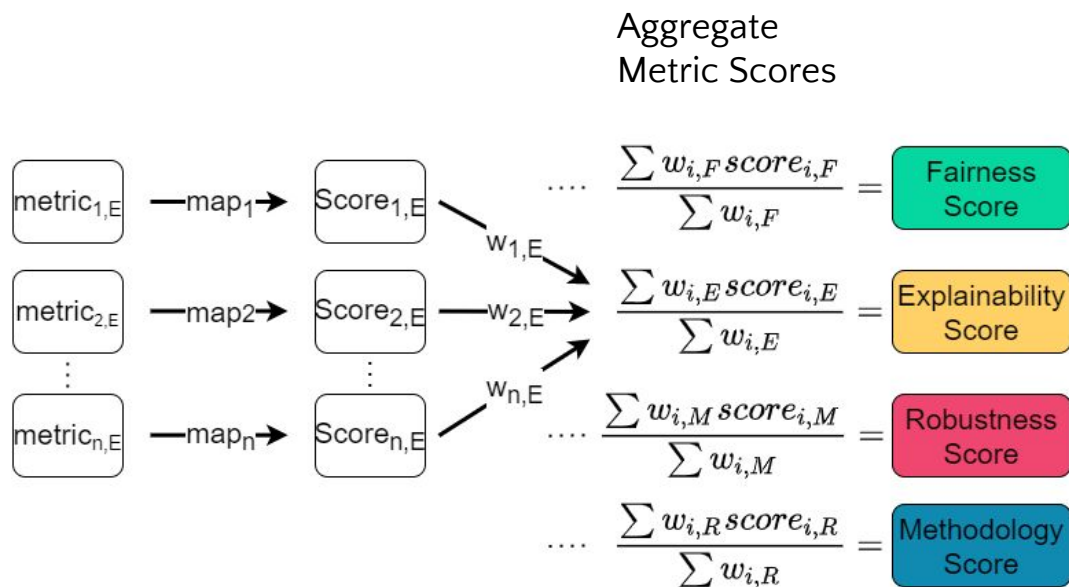
# Trusted AI Algorithm Design

Metric Scores  
between 1-5



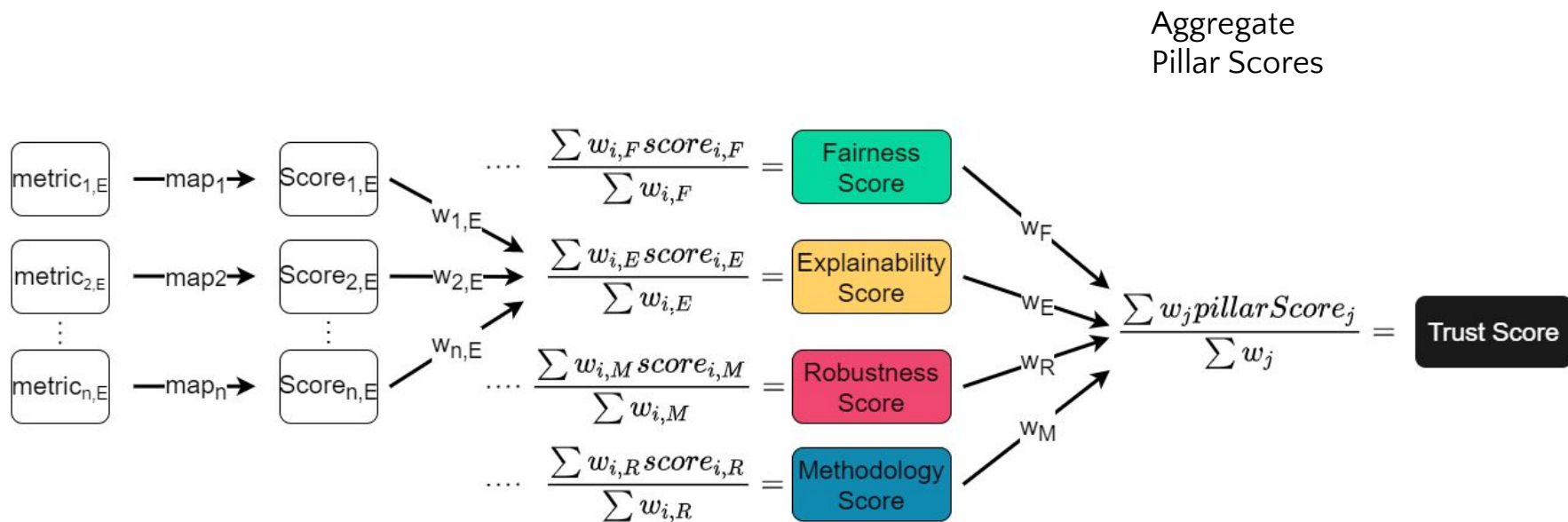


# Trusted AI Algorithm Design



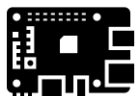


# Trusted AI Algorithm Design



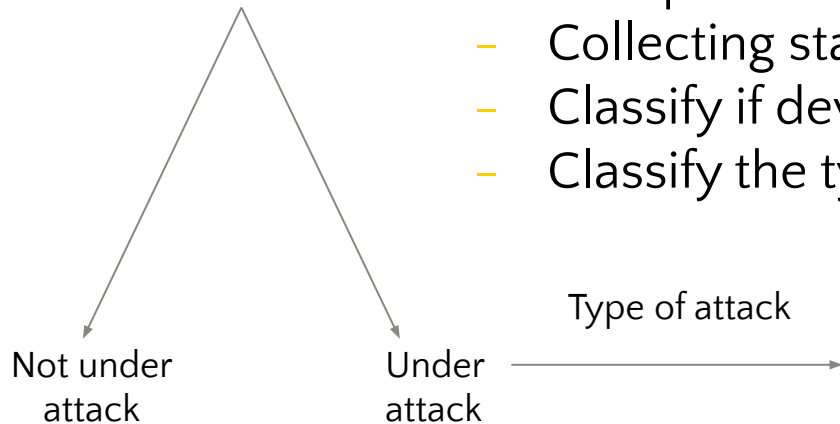


## Validation Scenario 1



### IoT Attack Classification:

- Multiple client devices
- Collecting status information
- Classify if device is subject to attack
- Classify the type of attack



1. *Fakepsd attack*
2. *Sendout attack*
3. *Write attack*
4. *Random attack*
5. *Exchange attack*
6. *Hide attack*



## Validation Scenario 2



### Credit Card Approval:

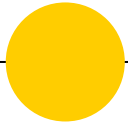
- Personal user information
- Credit scoring
- Classify if a credit card can be given

Credit card  
approved

Credit card  
rejected

# Demo Prototype

<https://www.csg.uzh.ch/trusted-ai/>







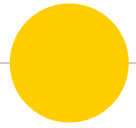
## Conclusion and Future Work

---

- Compile a general taxonomy
- Develop an algorithm to automatically compute the trust score of machine learning models
- Validation of the proposed solution
- Add a new pillar (e.g. Privacy)
- Extend the set of metrics
- Add support for regression model
- Add support for other ML libraries

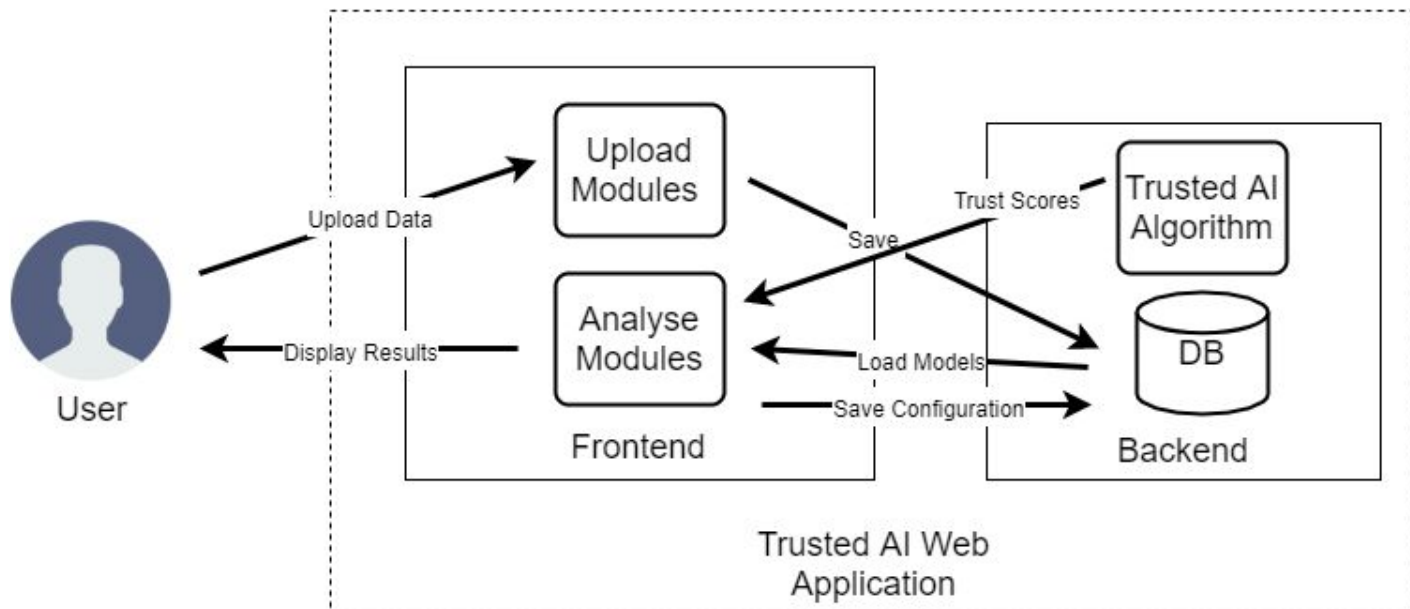


**Questions?**



# **Backup Slides**

# App Architecture





# Trusted AI Algorithm

---

<b>Algorithm</b>	Trusted AI Algorithm
------------------	----------------------

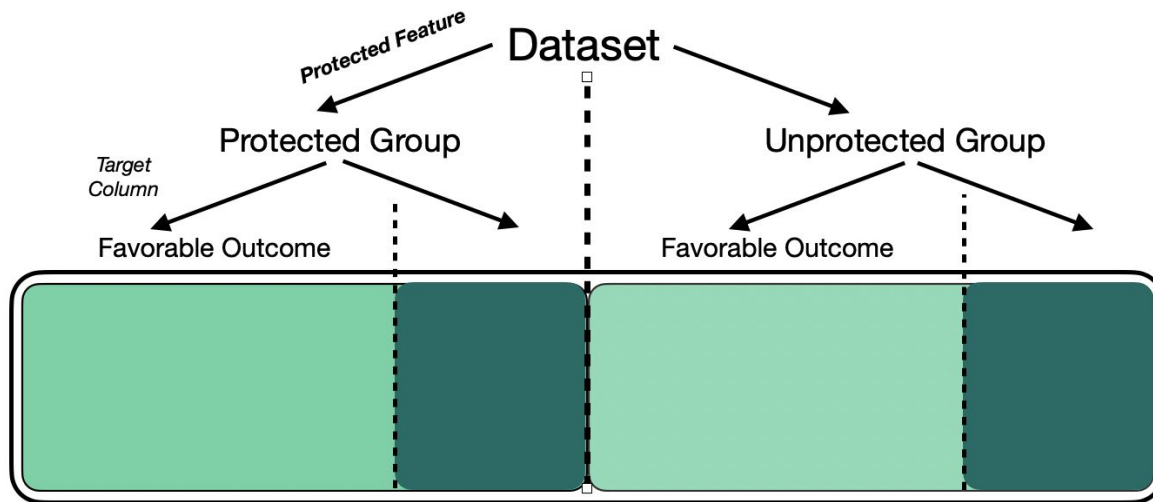
---

```
1: function TRUSTED_AI( $model, data_{train}, data_{test}, factsheet, config_{map}, config_{weights}$ )
2:    $trust\_score \leftarrow 0$ 
3:    $scores\_pillars \leftarrow \text{empty dictionary}$ 
4:   for  $p$  in  $pillars$  do
5:      $score_p \leftarrow 0$ 
6:      $scores\_metrics \leftarrow \text{empty dictionary}$ 
7:      $metrics \leftarrow get\_metrics(p)$ 
8:     for  $m$  in  $metrics$  do
9:        $args \leftarrow config_{map}[m]$ 
10:       $scores\_metrics[m] \leftarrow get\_score_m(args)$ 
11:      for ( $m \in keys, v \in values$ ) in  $scores\_metrics$  do
12:         $w_m \leftarrow config_{weights}[m]$ 
13:         $score_p \leftarrow score_p + w_m * v$ 
14:       $scores\_pillars[p] \leftarrow score_p$ 
15:    for ( $p \in keys, v \in values$ ) in  $scores\_pillars$  do
16:       $w_p \leftarrow config_{weights}[p]$ 
17:       $trust\_score \leftarrow trust\_score + w_p * v$ 
18:    return  $trust\_score$ 
```

---



# Fairness Computation





# Statistical Parity Difference

---

**Algorithm 10** Statistical Parity Difference Metric

---

```
1: function STATISTICAL_PARITY_DIFFERENCE(model, data, factsheet)
2:   protected_feature  $\leftarrow$  factsheet["protected_feature"]
3:   protected_values  $\leftarrow$  factsheet["protected_values"]
4:
5:   protected_group  $\leftarrow$  filter(data[protected_feature].is_in(protected_values))
6:   unprotected_group  $\leftarrow$  filter(data[protected_feature].not_in(protected_values))
7:
8:   protected_group_size  $\leftarrow$  size(protected_group)
9:   unprotected_group_size  $\leftarrow$  size(unprotected_group)
10:
11:   target_column  $\leftarrow$  factsheet["target_column"]
12:   favorable_outcomes  $\leftarrow$  factsheet["favorable_outcomes"]
13:
14:   favored_protected_group  $\leftarrow$  filter(protected_group[target_column].is_in(favorable_outcomes))
15:   favored_unprotected_group  $\leftarrow$  filter(unprotected_group[target_column].is_in(favorable_outcomes))
16:
17:   favored_protected_group_size  $\leftarrow$  size(favored_protected_group)
18:   favored_unprotected_group_size  $\leftarrow$  size(favored_unprotected_group)
19:
20:   favored_protected_ratio  $\leftarrow$  favored_protected_group_size/protected_group_size
21:   favored_unprotected_ratio  $\leftarrow$  favored_unprotected_group_size/unprotected_group_size
22:
23:   statistical_parity_difference  $\leftarrow$  |favored_unprotected_ratio - favored_protected_ratio|
24:   return statistical_parity_difference
```

---



# Equal Opportunity Difference

---

**Algorithm 11** Equal Opportunity Difference Metric

---

```
1: function EQUAL_OPPORTUNITY_DIFFERENCE(model, data, factsheet)
2:   protected_feature  $\leftarrow$  factsheet["protected_feature"]
3:   protected_values  $\leftarrow$  factsheet["protected_values"]
4:
5:   protected_group  $\leftarrow$  filter(data[protected_feature].is_in(protected_values))
6:   unprotected_group  $\leftarrow$  filter(data[protected_feature].not_in(protected_values))
7:
8:   target_column  $\leftarrow$  factsheet["target_column"]
9:   favorable_outcomes  $\leftarrow$  factsheet["favorable_outcomes"]
10:
11:   favorable_protected_group  $\leftarrow$  filter(protected_group[target_column].is_in(favorable_outcomes))
12:   favorable_unprotected_group  $\leftarrow$  filter(unprotected_group[target_column].is_in(favorable_outcomes))
13:
14:   favorable_protected_group_size  $\leftarrow$  size(favorable_protected_group)
15:   favorable_unprotected_group_size  $\leftarrow$  size(favorable_unprotected_group)
16:
17:   predicted_favorable_protected_group  $\leftarrow$  filter(favorable_protected_group[y_pred].is_in(favorable_outcomes))
18:   predicted_favorable_unprotected_group  $\leftarrow$  filter(favorable_unprotected_group[y_pred].is_in(favorable_outcomes))
19:
20:   predicted_favorable_protected_group_size  $\leftarrow$  size(predicted_favorable_protected_group)
21:   predicted_favorable_unprotected_group_size  $\leftarrow$  size(predicted_favorable_unprotected_group)
22:
23:   predicted_favorable_protected_ratio  $\leftarrow$  predicted_favorable_protected_group_size / favorable_protected_group_size
24:   predicted_favorable_unprotected_ratio  $\leftarrow$  predicted_favorable_unprotected_group_size / favorable_unprotected_group_size
25:
26:   equal_opportunity_difference  $\leftarrow$  |predicted_favorable_protected_ratio - predicted_favorable_unprotected_ratio|
27:   return equal_opportunity_difference
```

---





## Equal Opportunity Computation Example

Unprotected Group			Protected Group		
P=0	$\hat{Y} = 0$	$\hat{Y} = 1$	P=1	$\hat{Y} = 0$	$\hat{Y} = 1$
Y = 0	8000	3000	Y = 0	7000	9000
Y = 1	4000	6000	Y = 1	6000	9000

$$TPR_{unprotected} = Pr\{\hat{Y} = 1|P = 0, Y = 1\} = 6000/10000 = 60\%$$

$$TPR_{protected} = Pr\{\hat{Y} = 1|P = 1, Y = 1\} = 9000/15000 = 60\%$$

$$FPR_{unprotected} = Pr\{\hat{Y} = 1|P = 0, Y = 0\} = 3000/11000 \approx 27.27\%$$

$$FPR_{protected} = Pr\{\hat{Y} = 1|P = 1, Y = 0\} = 9000/16000 \approx 56.25\%$$

$$EOD(\hat{Y}, P) = \frac{9000}{16000} - \frac{3000}{11000} \approx 56.25\% - 27.27\% \approx 28.98\%$$