# A Framework for Understanding Unintended Consequences of Machine Learning

Robert Różański

unintended and potentially harmful consequences:
- underperforming consumer products/services
- unequal access to resources
- legal failure/injustice

unintended and potentially harmful consequences:
- underperforming consumer products/services
- unequal access to resources
- legal failure/injustice

no universal solution:

e.g. in the context of protected attributes, sometimes their use is prohibited (sex & job application), while in others they must be used to ensure good treatment (sex & medical diagnosis)
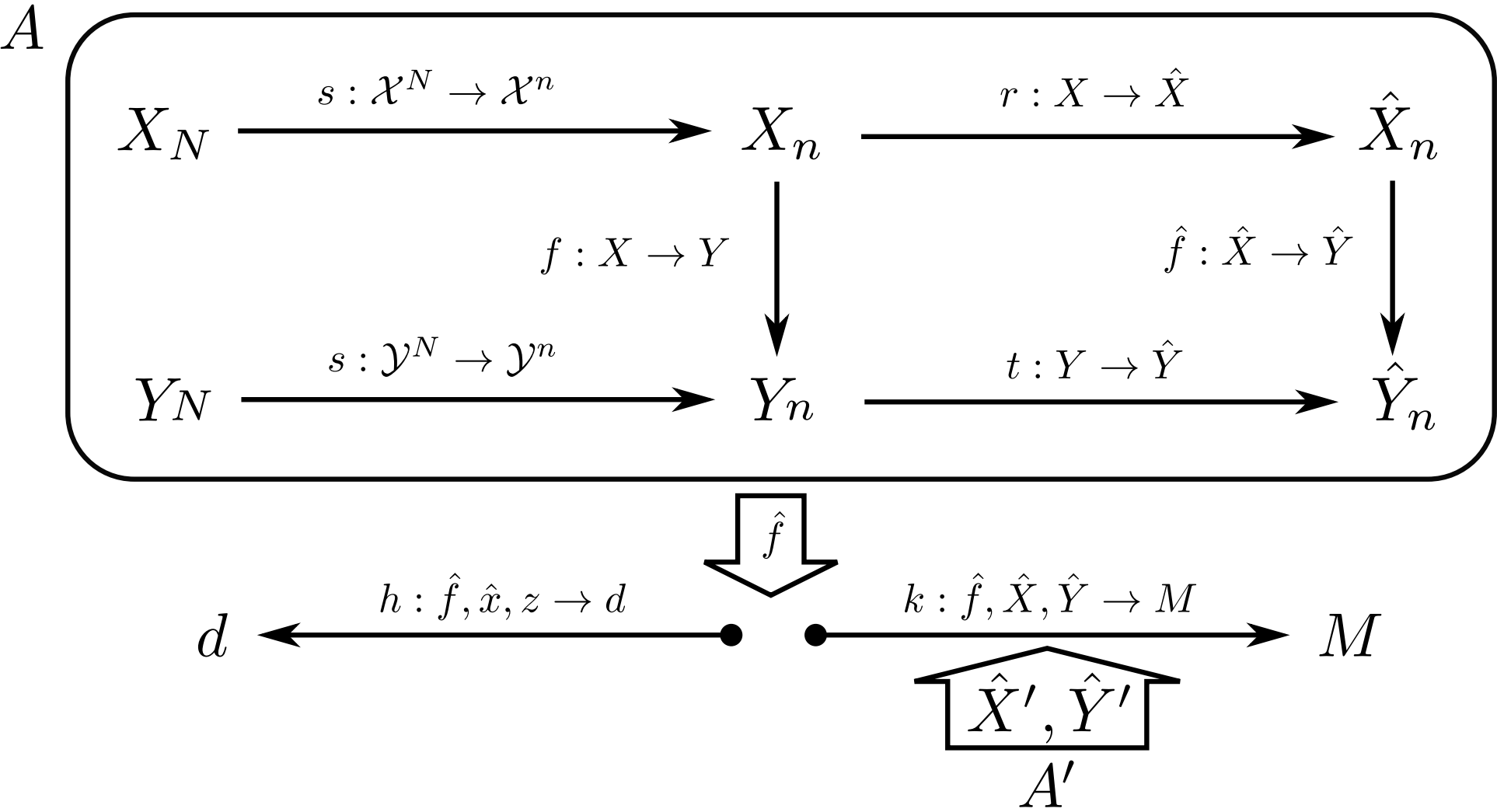
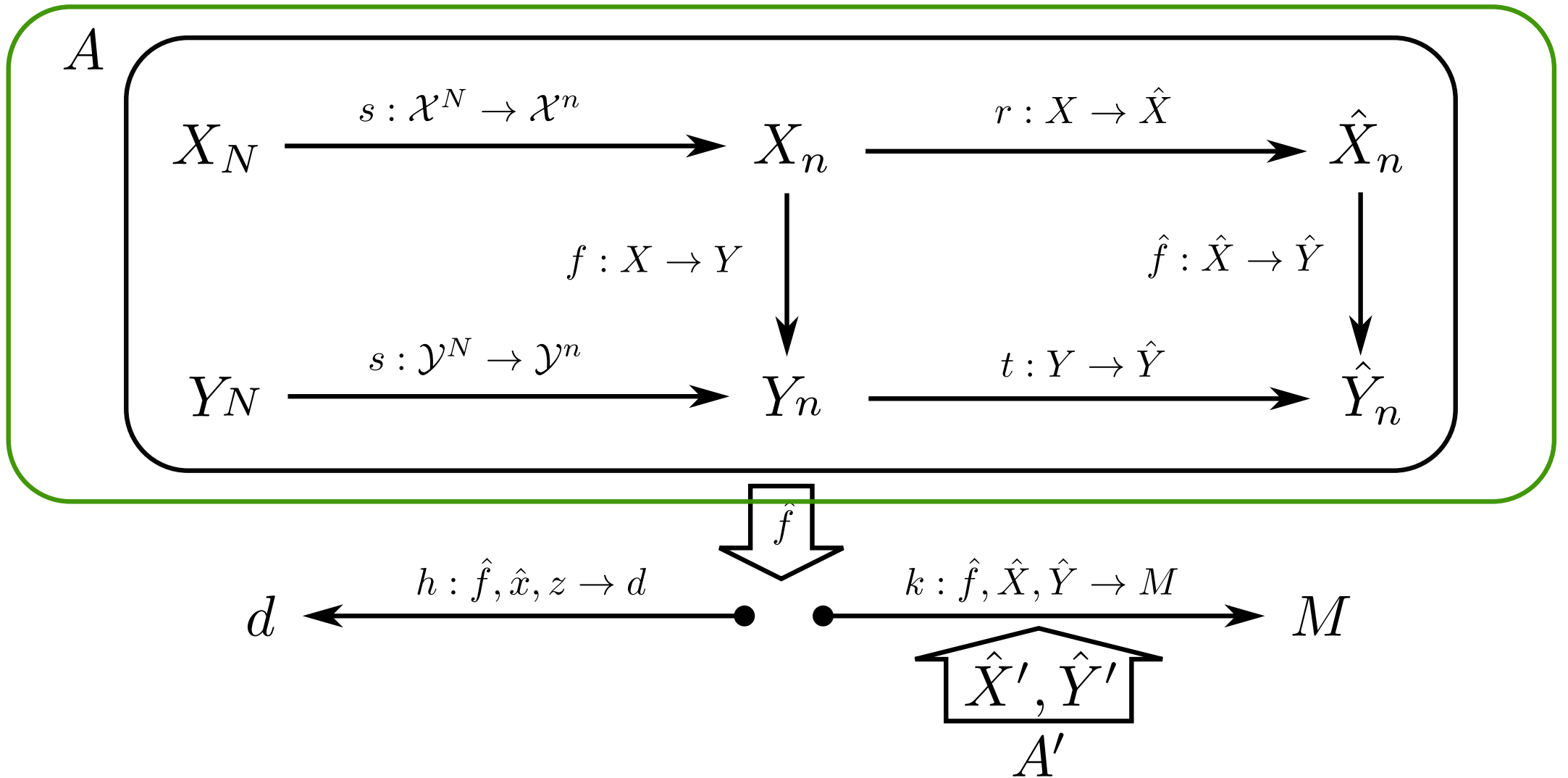unintended and potentially harmful consequences:
- underperforming consumer products/services
- unequal access to resources
- legal failure/injustice
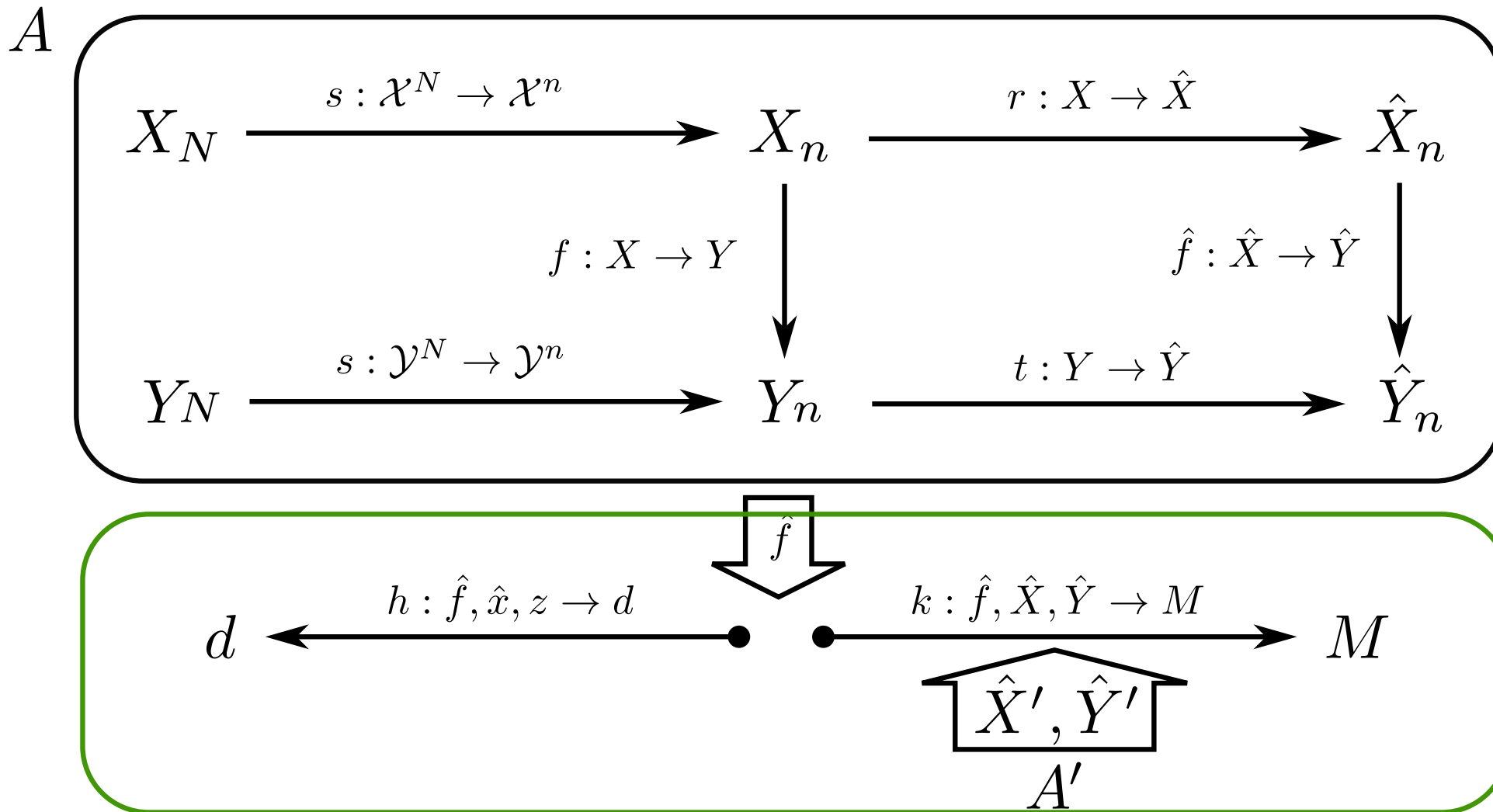
no universal solution:

e.g. in the context of protected attributes, sometimes their use is prohibited (sex & job application), while in others they must be used to ensure good treatment (sex & medical diagnosis)

there is a lack of common framework that would allow to compare problems and solutions

$A$

$X_N \xrightarrow{s : \mathcal{X}^N \to \mathcal{X}^n} X_n \xrightarrow{r : X \to \hat{X}} \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{s : \mathcal{Y}^N \to \mathcal{Y}^n} Y_n \xrightarrow{t : Y \to \hat{Y}} \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{h : \hat{f}, \hat{x}, z \to d} \bullet \quad \bullet \xrightarrow{k : \hat{f}, \hat{X}, \hat{Y} \to M} M$

$\hat{X}', \hat{Y}'$

$A'$

$$A$$

$$X_N \xrightarrow{s:\mathcal{X}^N \to \mathcal{X}^n} X_n \xrightarrow{r:X \to \hat{X}} \hat{X}_n$$

$$f:X \to Y \qquad \hat{f}:\hat{X} \to \hat{Y}$$

$$Y_N \xrightarrow{s:\mathcal{Y}^N \to \mathcal{Y}^n} Y_n \xrightarrow{t:Y \to \hat{Y}} \hat{Y}_n$$

$$\hat{f}$$

$$d \xleftarrow{h:\hat{f},\hat{x},z \to d} \bullet \quad \bullet \xrightarrow{k:\hat{f},\hat{X},\hat{Y} \to M} M$$

$$\hat{X}', \hat{Y}'$$

$$A'$$

model building
process

$A$

$$X_N \xrightarrow{\ s : \mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r : X \to \hat{X}\ } \hat{X}_n$$

$$f : X \to Y \qquad\qquad \hat{f} : \hat{X} \to \hat{Y}$$

$$Y_N \xrightarrow{\ s : \mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t : Y \to \hat{Y}\ } \hat{Y}_n$$

$\hat{f}$

$$d \xleftarrow{\ h : \hat{f}, \hat{x}, z \to d\ } \bullet \quad \bullet \xrightarrow{\ k : \hat{f}, \hat{X}, \hat{Y} \to M\ } M$$

$$\hat{X}', \hat{Y}'$$

$A'$

evaluation &
deployment

$A$

$X_N \xrightarrow{\;s : \mathcal{X}^N \to \mathcal{X}^n\;} X_n \xrightarrow{\;r : X \to \hat{X}\;} \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{\;s : \mathcal{Y}^N \to \mathcal{Y}^n\;} Y_n \xrightarrow{\;t : Y \to \hat{Y}\;} \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{\;h : \hat{f}, \hat{x}, z \to d\;} \bullet \quad \bullet \xrightarrow{\;k : \hat{f}, \hat{X}, \hat{Y} \to M\;} M$

$\hat{X}', \hat{Y}'$

$A'$

Ideal, underlying features on the whole population

$A$

$X_N \xrightarrow{\quad s : \mathcal{X}^N \to \mathcal{X}^n \quad} X_n \xrightarrow{\quad r : X \to \hat{X} \quad} \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{\quad s : \mathcal{Y}^N \to \mathcal{Y}^n \quad} Y_n \xrightarrow{\quad t : Y \to \hat{Y} \quad} \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{\quad h : \hat{f}, \hat{x}, z \to d \quad} \bullet \quad \bullet \xrightarrow{\quad k : \hat{f}, \hat{X}, \hat{Y} \to M \quad} M$

$\hat{X}', \hat{Y}'$

$A'$

sampling functions

$A$

$X_N \xrightarrow{\ s : \mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r : X \to \hat{X}\ } \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{\ s : \mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t : Y \to \hat{Y}\ } \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{\ h : \hat{f}, \hat{x}, z \to d\ } \bullet \quad \bullet \xrightarrow{\ k : \hat{f}, \hat{X}, \hat{Y} \to M\ } M$

$\hat{X}', \hat{Y}'$

$A'$

Ideal, underlying features on the sample

$A$

$X_N \xrightarrow{\ s : \mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r : X \to \hat{X}\ } \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{\ s : \mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t : Y \to \hat{Y}\ } \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{\ h : \hat{f}, \hat{x}, z \to d\ } \bullet \quad \bullet \xrightarrow{\ k : \hat{f}, \hat{X}, \hat{Y} \to M\ } M$

$\hat{X}', \hat{Y}'$

$A'$

model based on the
ideal, underlying
features

$A$

$X_N \xrightarrow{\ s : \mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r : X \to \hat{X}\ } \hat{X}_n$

$f : X \to Y$

$\hat{f} : \hat{X} \to \hat{Y}$

$Y_N \xrightarrow{\ s : \mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t : Y \to \hat{Y}\ } \hat{Y}_n$

$\hat{f}$

$d \xleftarrow{\ h : \hat{f}, \hat{x}, z \to d\ } \bullet \quad \bullet \xrightarrow{\ k : \hat{f}, \hat{X}, \hat{Y} \to M\ } M$

$\hat{X}', \hat{Y}'$

$A'$

projections from the
ideal features to the
measured ones

Actual, measured features on the sample

$A$

$X_N \xrightarrow{\ s:\mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r:X \to \hat{X}\ } \hat{X}_n$

$f:X \to Y$

$\boxed{\hat{f}:\hat{X} \to \hat{Y}}$

$Y_N \xrightarrow{\ s:\mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t:Y \to \hat{Y}\ } \hat{Y}_n$

$\hat{f}$

$d \xleftrightarrow{\ h:\hat{f},\hat{x},z \to d\ } \bullet \quad \bullet \xleftrightarrow{\ k:\hat{f},\hat{X},\hat{Y} \to M\ } M$

$\hat{X}',\hat{Y}'$

$A'$

actual model, based
on the measured
features

scoring function
test data
measure of success

real world decision
process
decision

historical bias: if the world as it is or was leads a model to produce outcomes that are not wanted.

Example: representation bias in image search results:

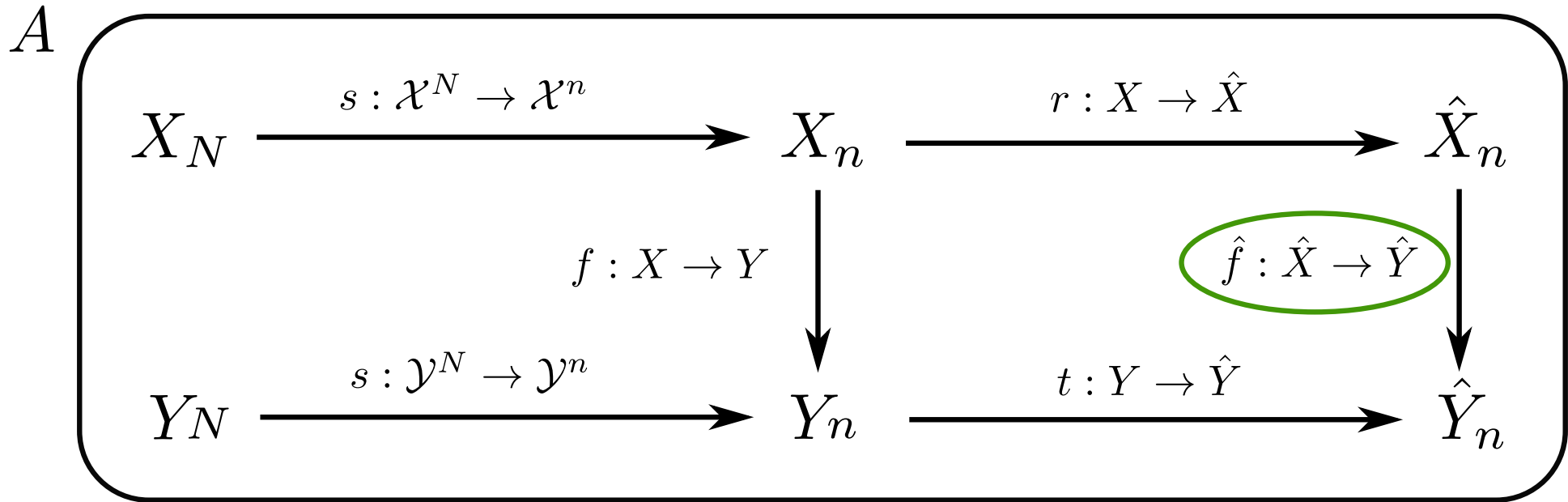https://www.theguardian.com/technology/2016/dec/05/google-alters-search-autocomplete-remove-are-jews-evil-suggestion

representation bias: certain parts of the input space are underrepresented. This could be caused not only by deficient sampling methods, but, for example, also because the population of interest has changed, etc.

Example: geographic diversity in the ImageNet (45% from USA)

measurement bias: available/measureable features and labels are noisy proxies for the features and labels of interest. Could be a result of differences in measurement process or data quality accross groups, etc.

Example: in predictive policing and recidivism prediction proxy variable "arrest" is used to measure "crime"
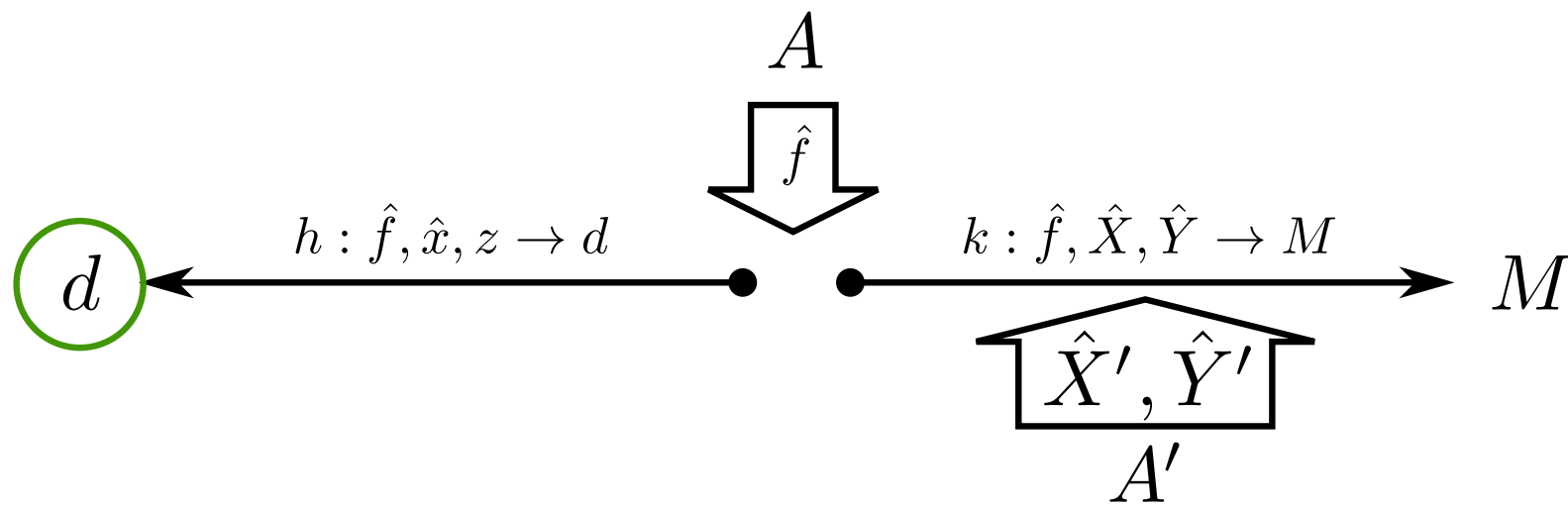
$A$

$$X_N \xrightarrow{\ s:\mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r:X \to \hat{X}\ } \hat{X}_n$$

$f:X \to Y$

$\hat{f}:\hat{X} \to \hat{Y}$

$$Y_N \xrightarrow{\ s:\mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t:Y \to \hat{Y}\ } \hat{Y}_n$$

aggregation bias: when a one-size-fit-all model is used
for groups with different conditional distributions. This can lead to
underperforming models even if the sample is balanced.

Example: in diabetes diagnosis/monitoring, clinical meaning of
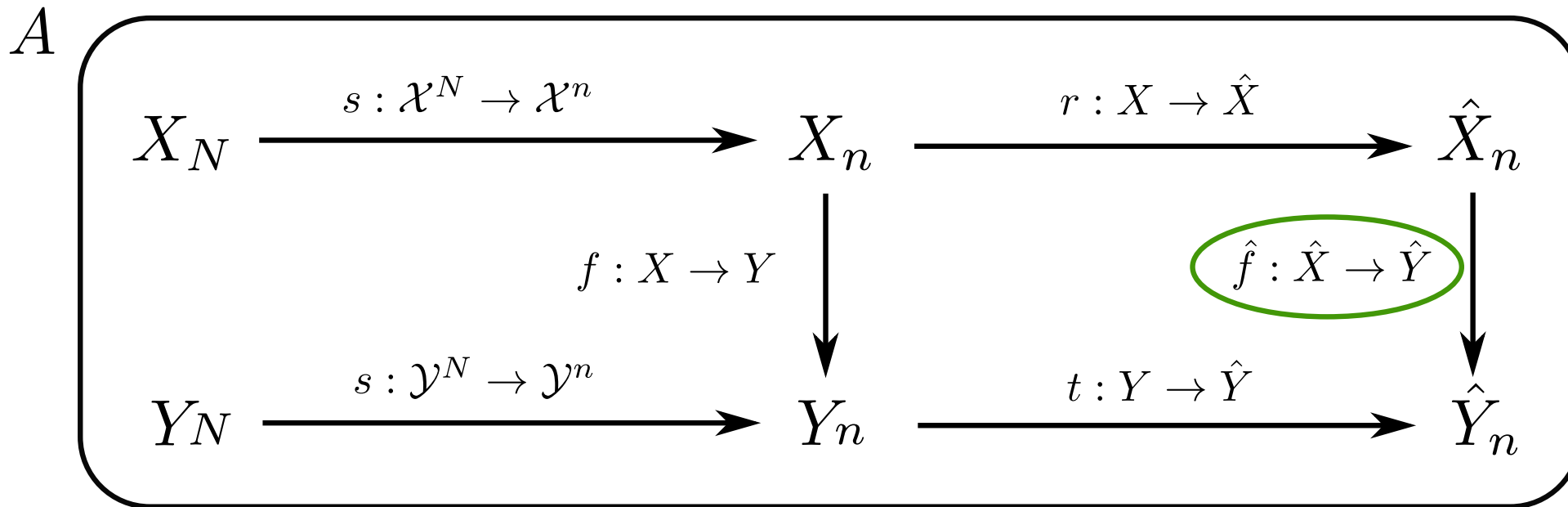HbA1c levels differs accross ethnicities and genders

evaluation bias: when the evaluation and/or benchmark data for an algorithm doesn't represent the target population.

Example: underperformance of facial analysis models on dark-skinned females (7.4% and 4.4% of the images in benchmark datasets such as Adience and IJB-A are of dark-skinned female faces.)
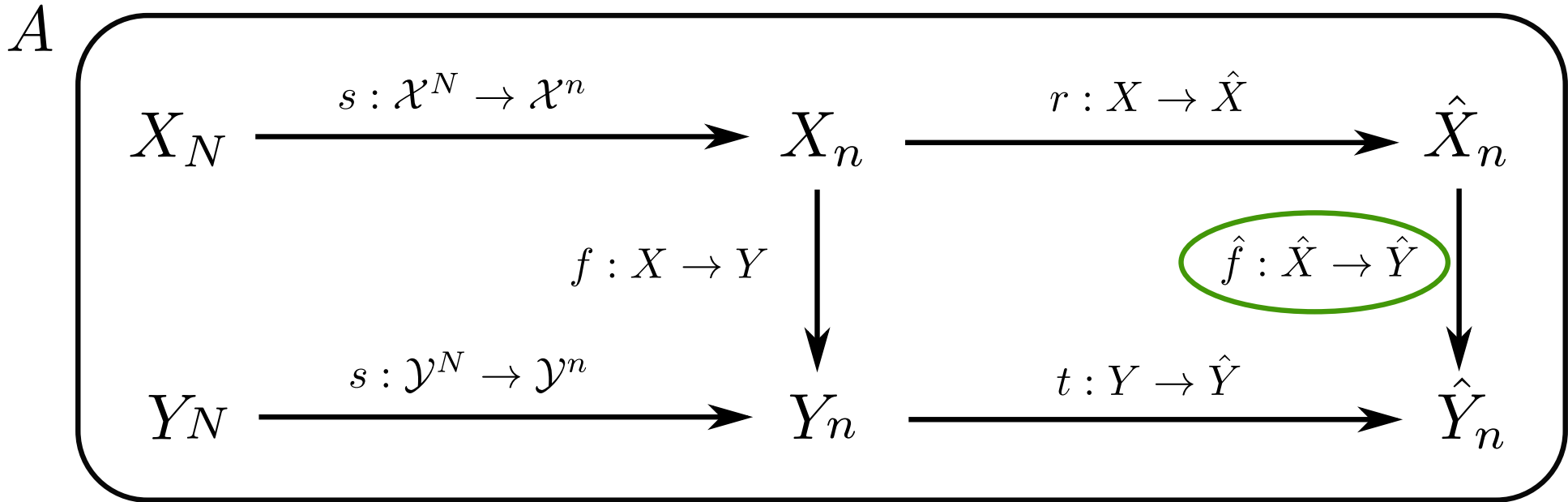
deployment bias: there is a mismatch between the problem a model is intended to solve and the way in which it is actually used.

Example: models predicting person's likelihood of commiting a future crime used in "off label" ways, e.g. to help determine the length of a sentence.
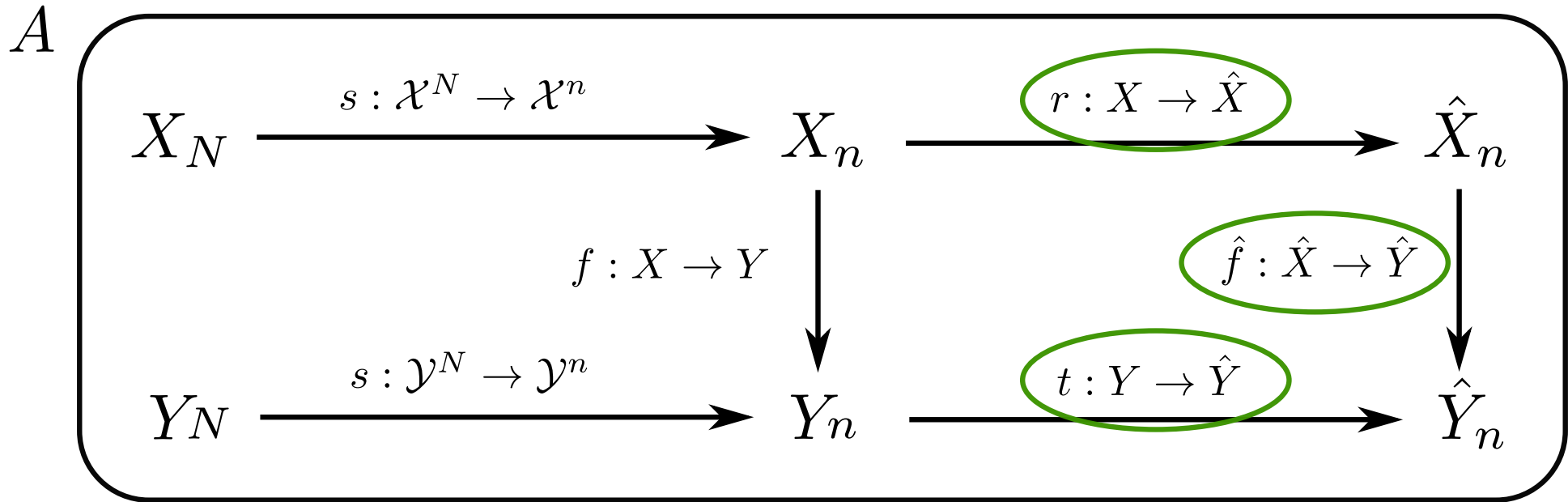
case study: aggregation bias

case study: aggregation bias

adjusting the model:
   - using multitask learning
   - using more complex model, which can capture the differences in conditional distribution between the groups (if enough data available)

A

$$X_N \xrightarrow{\ s : \mathcal{X}^N \to \mathcal{X}^n\ } X_n \xrightarrow{\ r : X \to \hat{X}\ } \hat{X}_n$$

$$f : X \to Y \qquad \hat{f} : \hat{X} \to \hat{Y}$$

$$Y_N \xrightarrow{\ s : \mathcal{Y}^N \to \mathcal{Y}^n\ } Y_n \xrightarrow{\ t : Y \to \hat{Y}\ } \hat{Y}_n$$

case study: aggregation bias

adjusting the model:
    - using multitask learning
    - using more complex model, which can capture the differences in conditional distribution between the groups (if enough data available)

adjusting the data:
    - changing the projection functions $r$ & $t$ (representation learning)