

slides available at:



Strategic ML: How to Learn With Data That **'Behaves'**

Nir Rosenfeld

Technion CS

tutorial @ ICML 2024

outside → in:



input



output



this is machine learning:



this is machine learning on **images**:



this is machine learning on *text*:



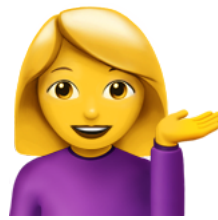
this is machine learning on... **humans?**



this is machine learning on... **humans!**



this is machine learning on... **humans!**

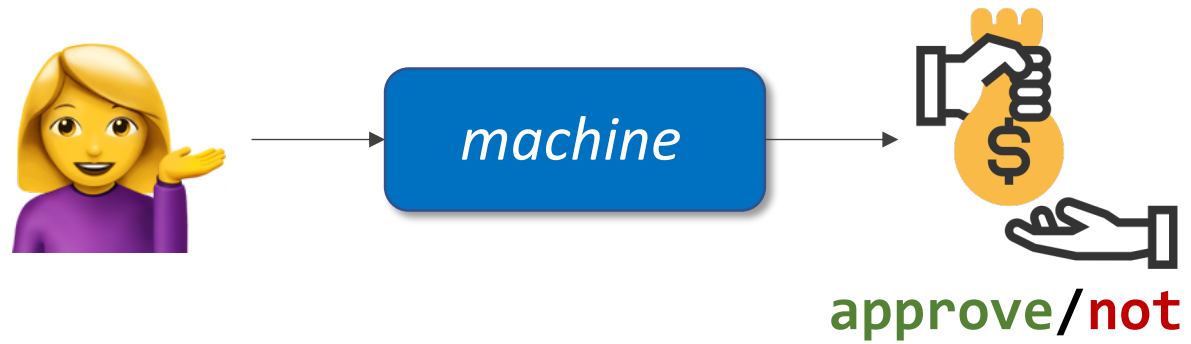


machine

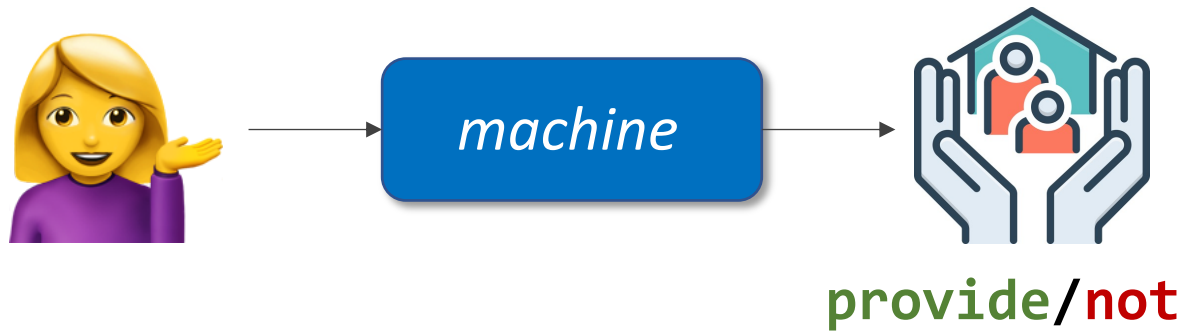


admit/not

this is machine learning on... **humans!**



this is machine learning on... **humans!**



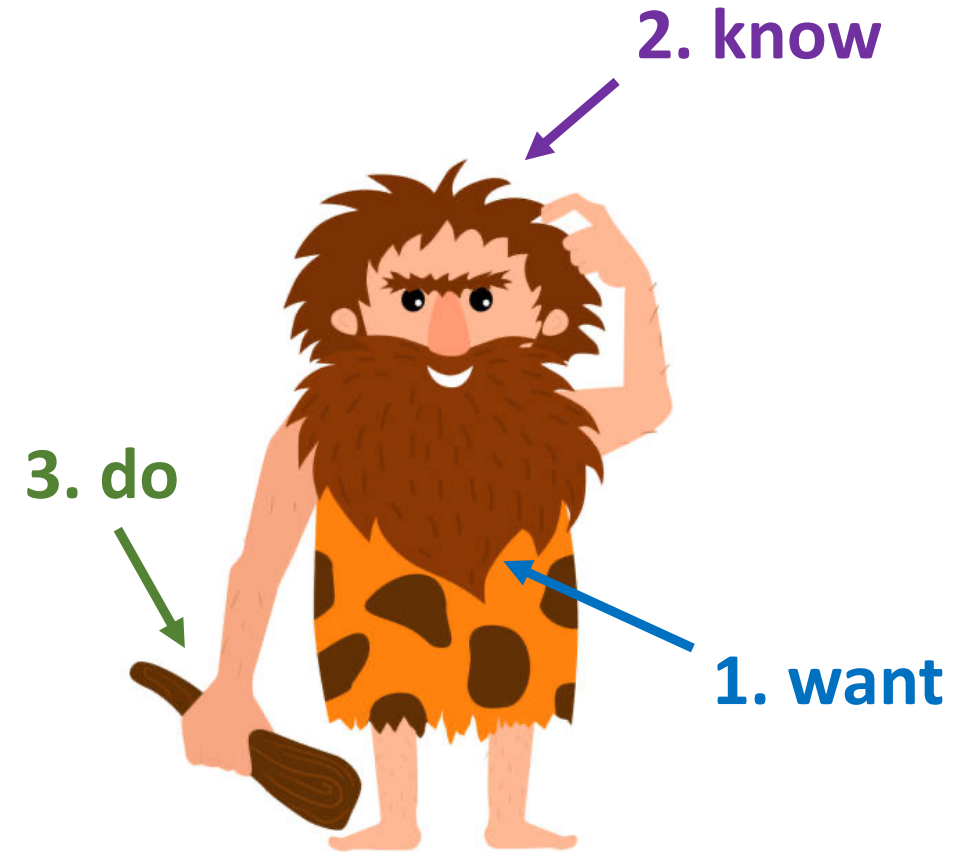
what could possibly go wrong?

(**or:** how does human behavior change learning and its outcomes?)

Strategic classification

- Builds on **conventional binary classification**
- **Augments** to account for **human behavior**
- Models humans as **inputs with agency**

- Allows (and requires!) to encode what humans:
- Basic elements of economic modeling
- Together, combine to determine **how humans behave**

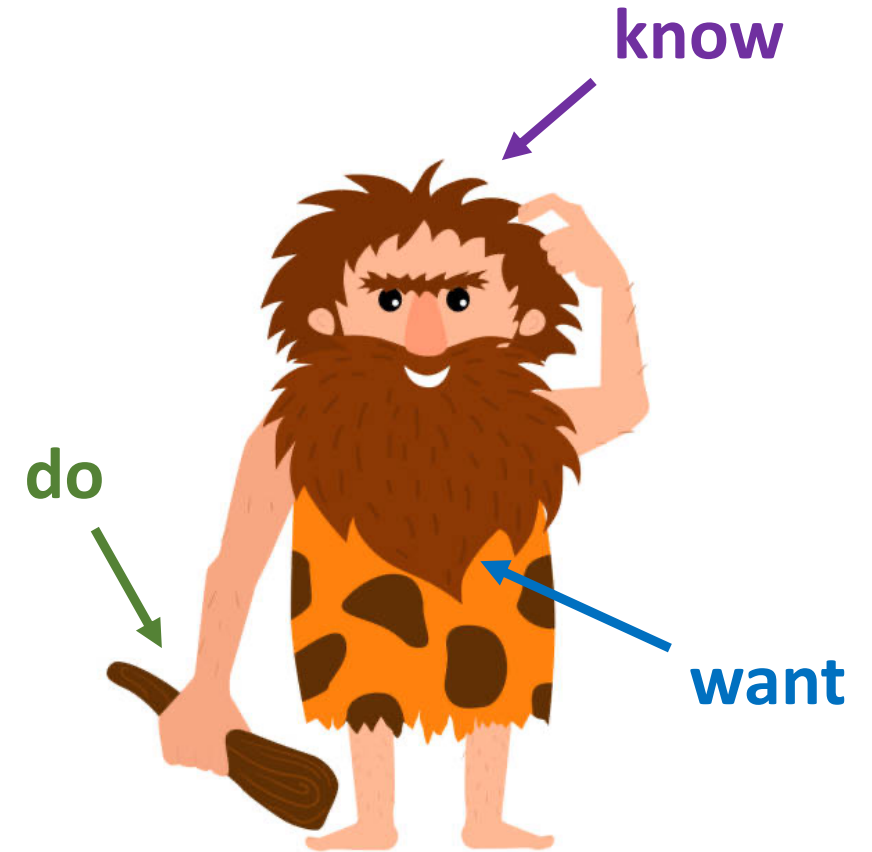


homo-sapiens

modeling challenge: weave these into learning setup

Strategic classification

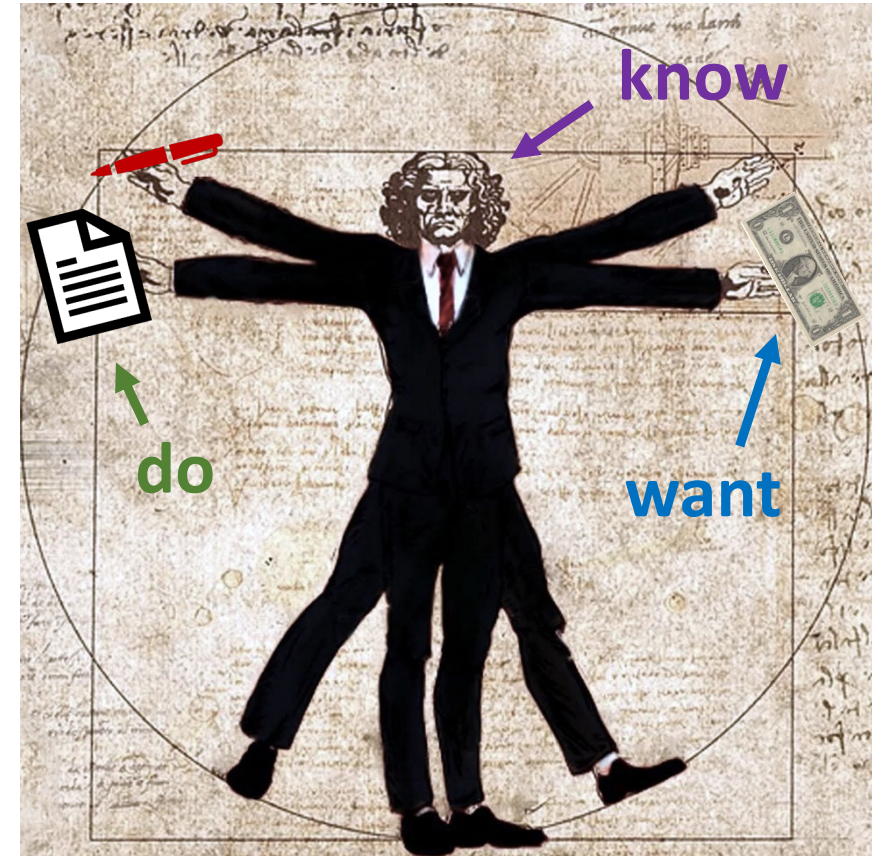
- **SC is great** because it is:
 - **simple enough** to permit tractable analysis
 - **powerful enough** to introduce novel challenges
 - **meaningful enough** to have social implications
 - **flexible enough** to permit extensions, variations, and generalizations
- Start with **rigid assumptions** – e.g., rationality:



homo-sapiens

Strategic classification

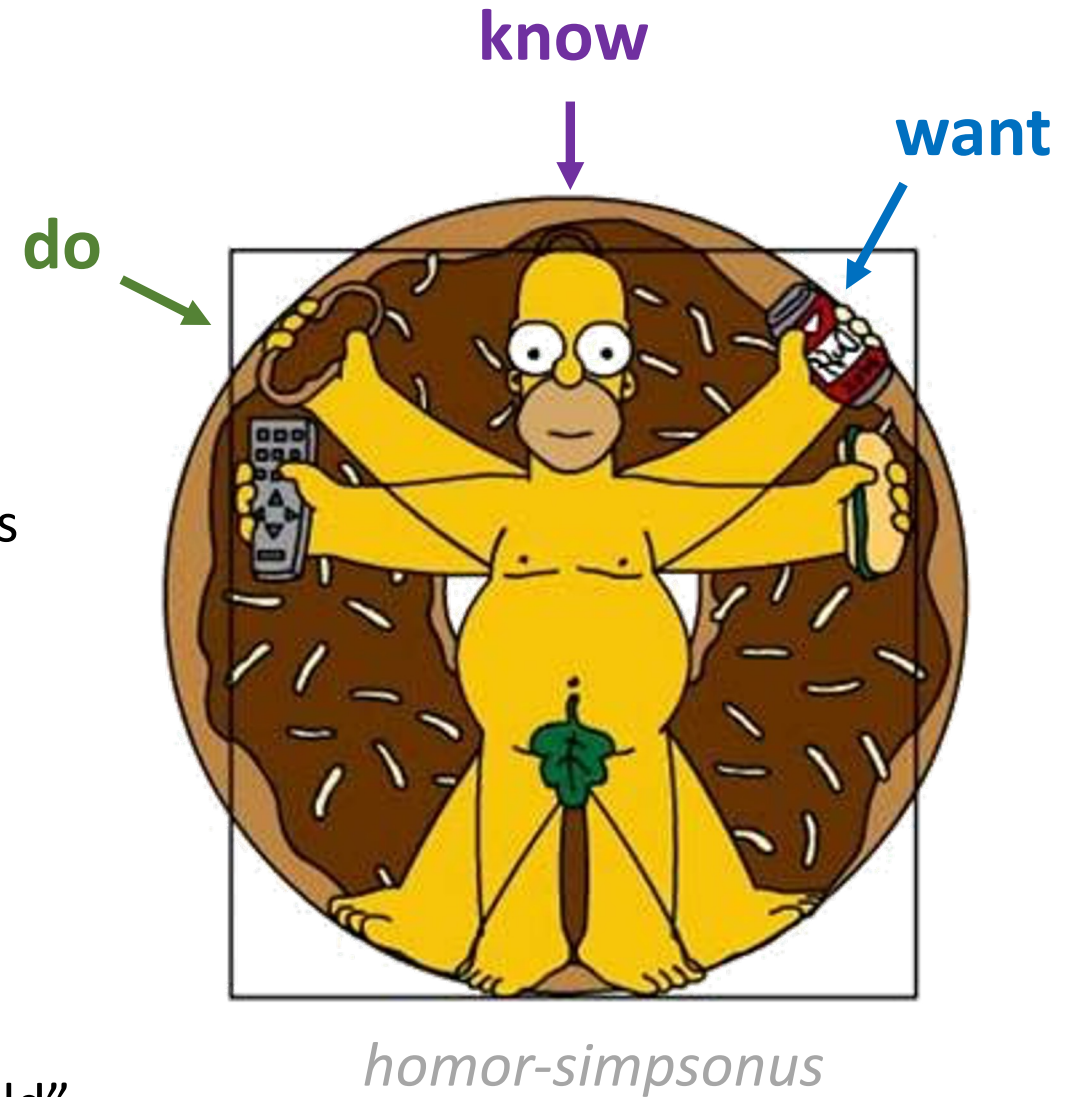
- **SC is great** because it is:
 - **simple enough** to permit tractable analysis
 - **powerful enough** to introduce novel challenges
 - **meaningful enough** to have social implications
 - **flexible enough** to permit extensions, variations, and generalizations
- Start with **rigid assumptions** – e.g., rationality:



homo-economicus

Strategic classification

- **SC is great** because it is:
 - **simple enough** to permit tractable analysis
 - **powerful enough** to introduce novel challenges
 - **meaningful enough** to have social implications
 - **flexible enough** to permit extensions, variations, and generalizations
- Start with **rigid assumptions** – e.g., rationality
- **Ultimate goal**: capture realistic behavior “in the wild”



Outline

- **Introduction**

- Three main sections:

- I) **ML aspects:** (~40 min)

- strategic learning – setup
 - as learning vs. as a game
 - optimization
 - generalization (stats)
 - modeling

- II) **Econ/GT aspects:** (~60 min)

- incentives (=want)
 - information (=know)
 - actions (=do)
 - limited resources
 - social welfare

- III) **Beyond:** (~20 min)

- causality
 - dependencies
 - over time

- **Challenges and opportunities**

- **Summary**

Tutorial theme and goals

- Introduction to **emerging new field**
- Many **open research questions**
- Much **potential for application**
- **Main theme:** transitioning from **theory** \mapsto **practice**
- Focus on **supervised batch setting** (covers “half” of literature; other part being online)

- More **breadth** (less depth) \rightarrow *see references*
- More **modeling** (less results)
- More **questions** (less answers)
- More **content** (less time) \rightarrow *fast paced!*

Strategic classification

an introduction

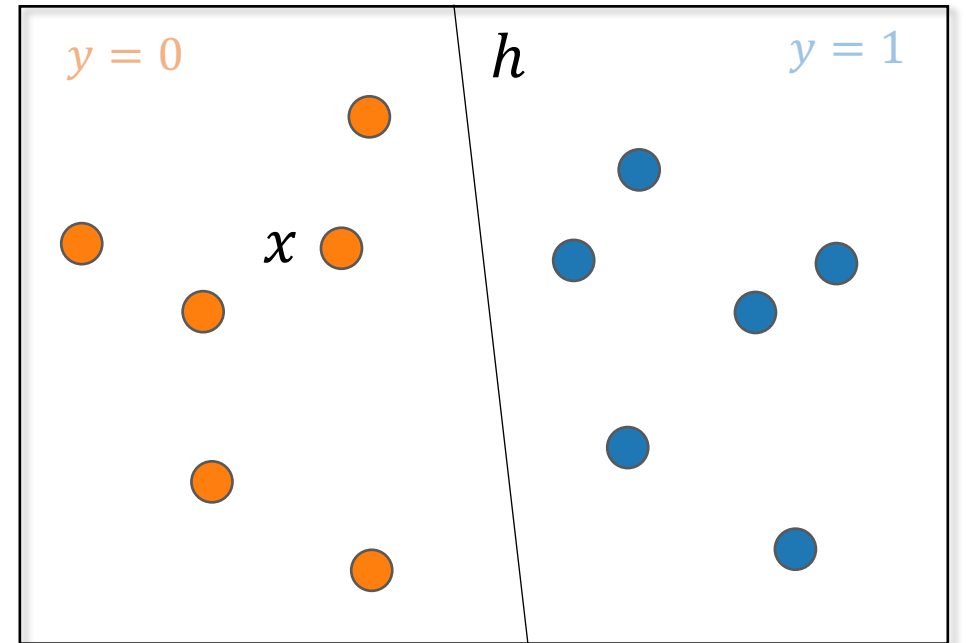
**standard
classification:**

train:

$$\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$$

learned model

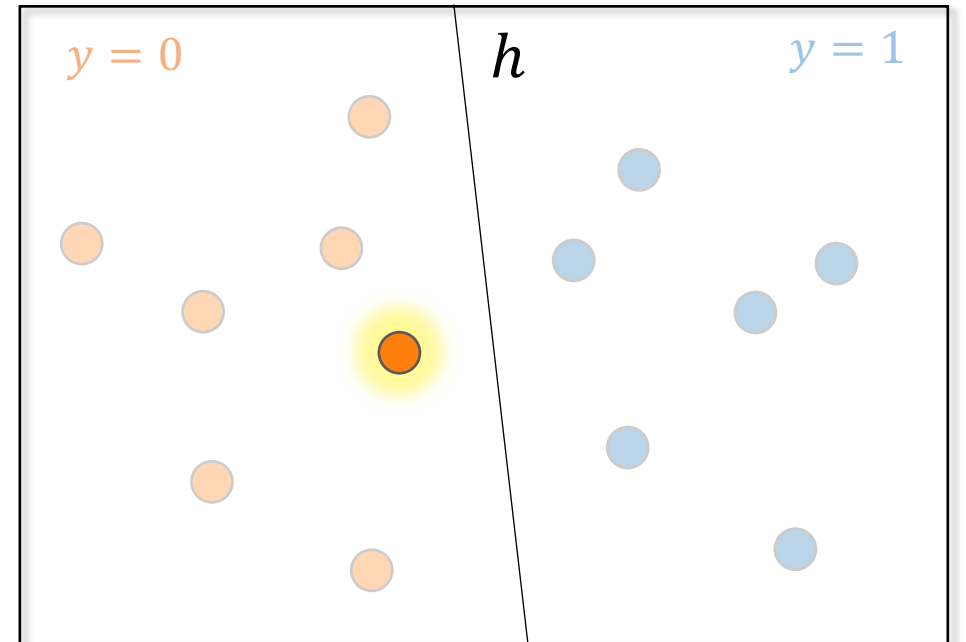
input features



standard classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$
learned model
input features

test: $h(x) = \hat{y} \approx y$
prediction *ground truth*



strategic

classification:

[BS'2011, HMPW'2016]

learned model

train:

$$\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$$

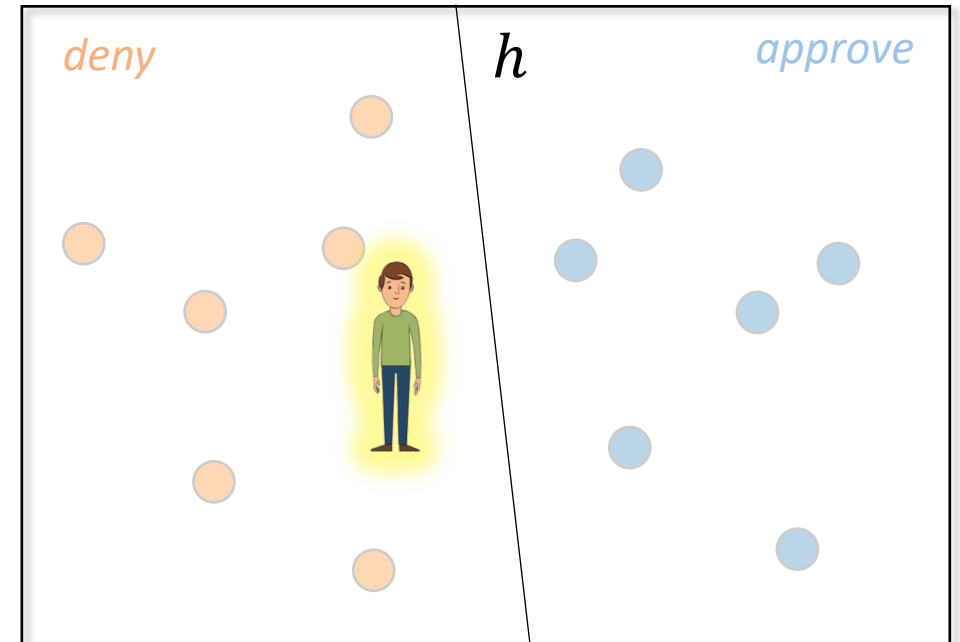
*representation of
human agent*

test:

$$h(x) = \hat{y} \approx y$$

prediction

ground truth

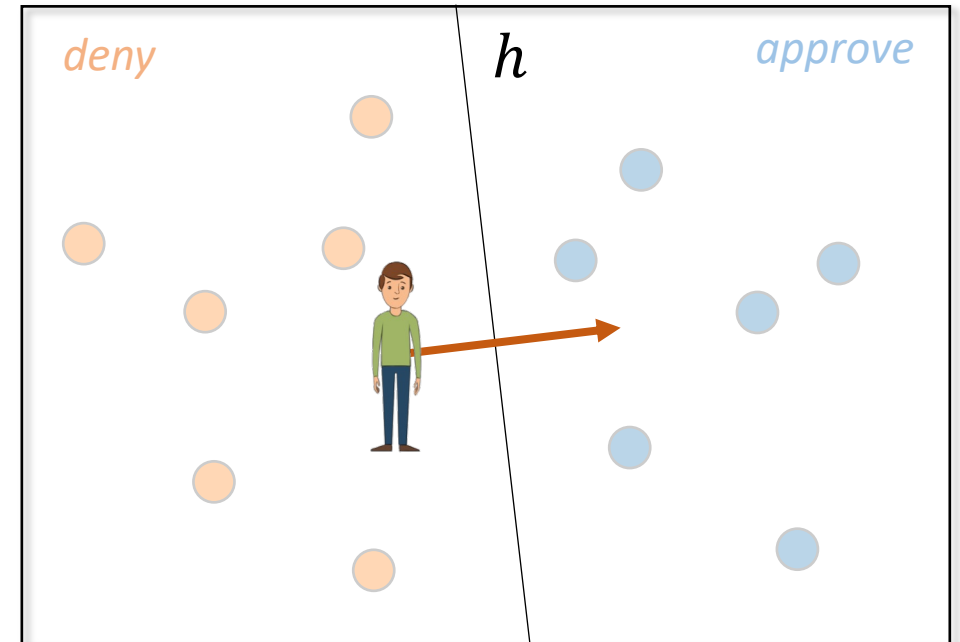


strategic
classification:



train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(x) = \hat{y} \approx y$

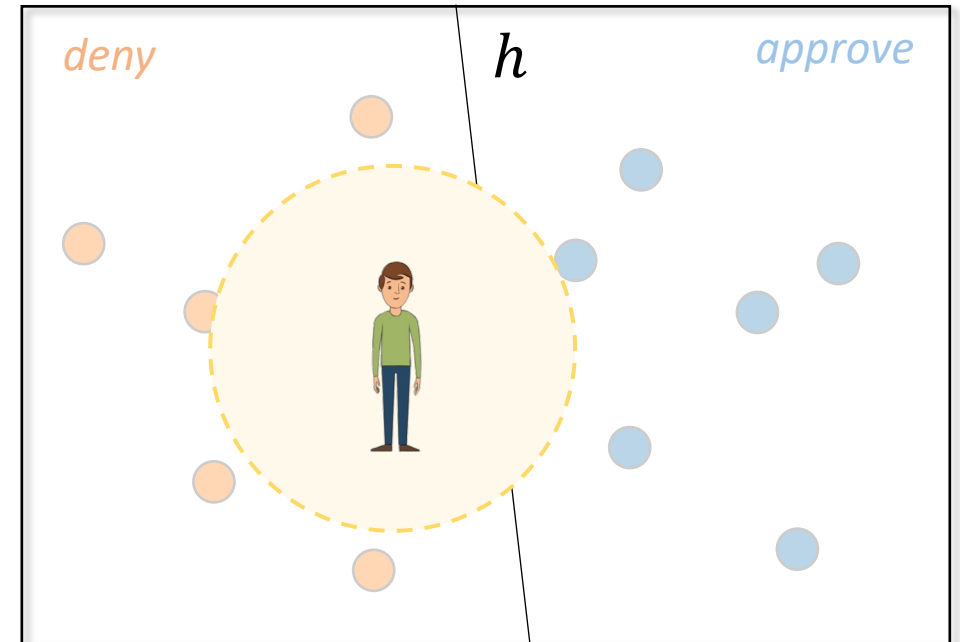


1. **want:** $\hat{y} = 1$ (get the loan)

strategic
classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(x) = \hat{y} \approx y$

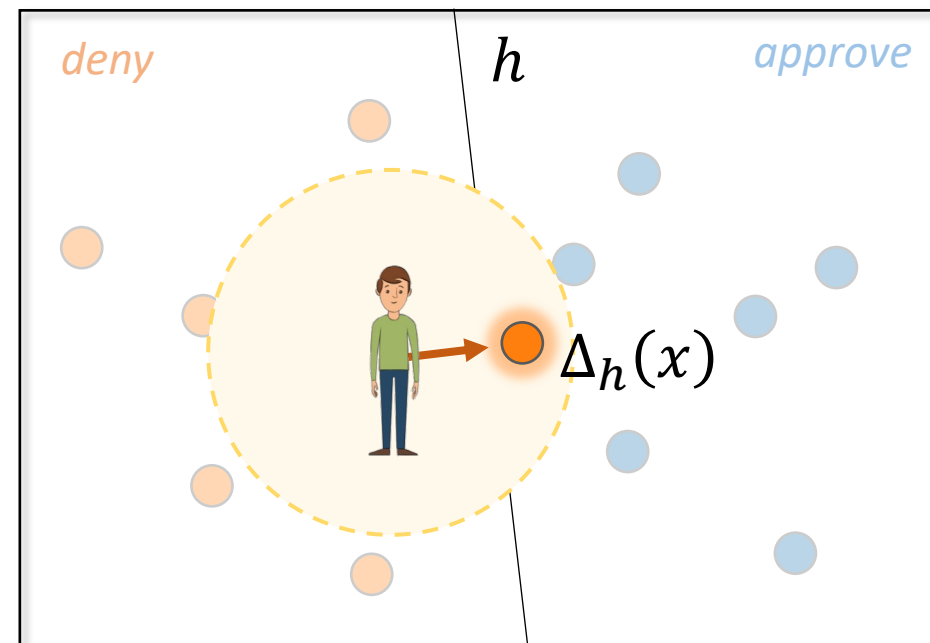


- 1. **want:** $\hat{y} = 1$ (get the loan)
- 2. **do:** modify features (at cost)

strategic
classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(x) = \hat{y} \approx y$



1. **want:** $\hat{y} = 1$ (get the loan)
2. **do:** modify features (at cost)
3. **know:** h (and cost function)

strategic
classification:



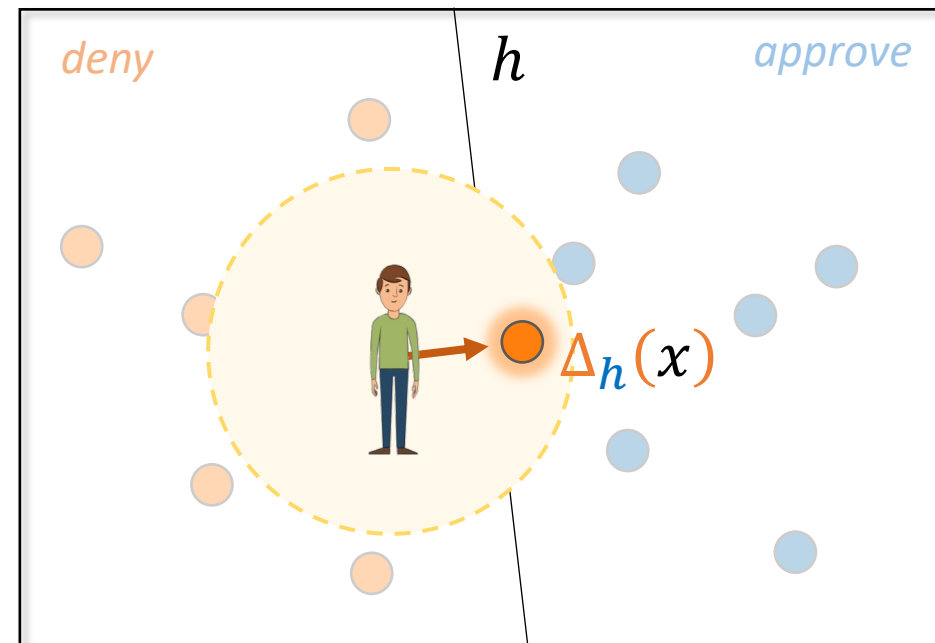
train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(x) = \hat{y} \approx y$

response: $x \mapsto x^h \triangleq \Delta_h(x)$



behavior



- 1. want:** $\hat{y} = 1$ (get the loan)
- 2. do:** modify features (at cost)
- 3. know:** h (and cost function)

strategic

classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

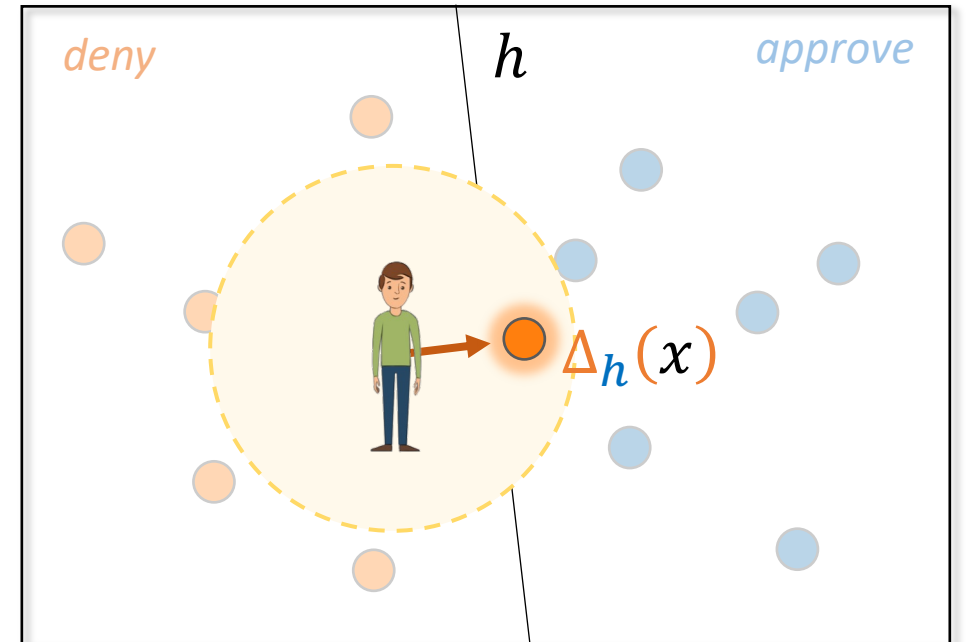
test: $h(x) = \hat{y} \approx y$

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

rational

utility
= prediction

cost
(e.g., norm)



rational \Rightarrow most cost-effective

\Rightarrow move on decision boundary

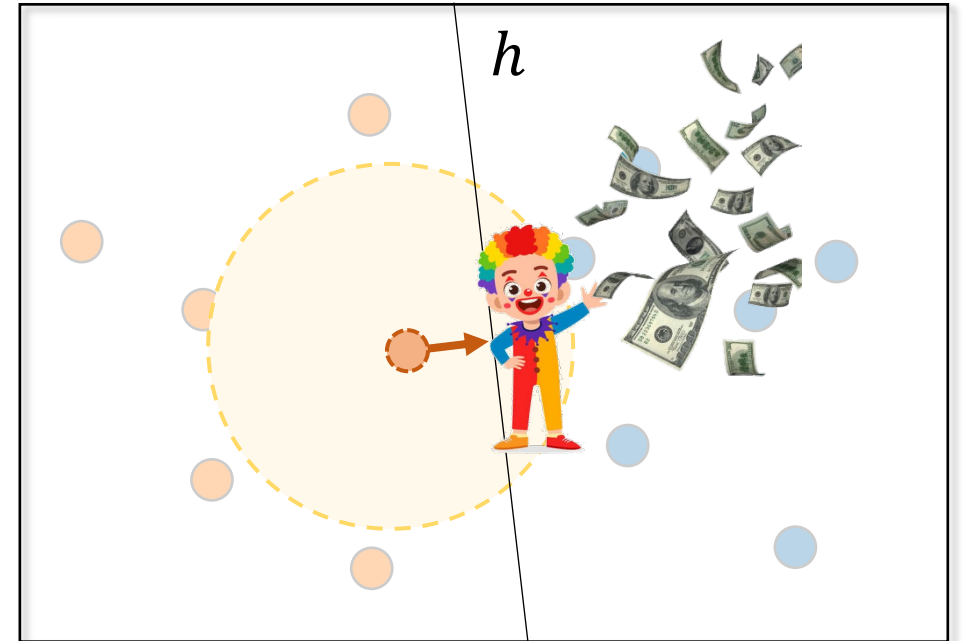
strategic

classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(\Delta_h(x)) = \hat{y} \neq y$

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$



goal: learning that is robust to strategic “gaming” behavior

strategic

classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(\Delta_h(x)) = \hat{y} \neq y$

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

Goodhart's law:

“If a measure becomes the public's goal, it is no longer a good measure.”



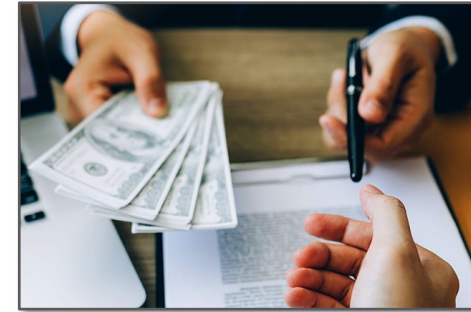
➤ Common examples:

**strategic
classification:**

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(\Delta_h(x)) = \hat{y} \neq y$

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$



strategic
classificat

train:

test:

respons



Holy grail: a realistically practical, well-understood, plug-n-play framework for strategic learning

- SC is great, by **frustrating**
- **Culprit** – lots (and lots) of **assumptions:**

- *outcomes are binary*
- *users always want positive outcomes*
- *costs are fixed, uniform, and known to all*
- *classifier is made public*
- *modifying x does not affect y*
- *changes to x are real (no mis-reporting)*
- *user actions = modify features*
- *users are rational (best-respond)*
- *users respond independently*
- *input data are 'clean' (=unmodified)*
- *playing order is fixed*
- *only single playing round*
- *system cares only for accuracy*
- *....*

- ongoing community effort to relax, extend, scrutinize, and generalize

strategic classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

test: $h(\Delta_h(x)) = \hat{y} \approx y$

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

- standard setup has **lots (and lots) of assumptions:** (implicit/explicit)
- *modifying x does not affect y*
 - *outcomes are binary*
 - *input data are 'clean' (=unmodified)*
 - *changes to x are real (no mis-reporting)*
 - *users always want positive outcomes*
 - *costs are fixed, uniform, and known to all*
 - *classifier is made public*
 - *user actions = modify features*
 - *users are rational (best-respond)*
 - *users respond independently*
 - *playing order is fixed*
 - *only single playing round*
 - *system cares only for accuracy*
- ongoing community effort to relax, extend, scrutinize, and generalize

strategic classification:

train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(\Delta_h(x))\}]$

test: $h(\Delta_h(x)) = \hat{y} \approx y$

consistent

response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

key point:

far from trivial! minor change \Rightarrow major implications

➤ standard setup has
lots (and lots) of assumptions: (implicit/explicit)

- *modifying x does not affect y*
- *outcomes are binary*
- *input data are 'clean' (=unmodified)*
- *changes to x are real (no mis-reporting)*
- *users always want positive outcomes*
- *costs are fixed, uniform, and known to all*
- *classifier is made public*
- *user actions = modify features*
- *users are rational (best-respond)*
- *users respond independently*
- *playing order is fixed*
- *only single playing round*
- *system cares only for accuracy*

ongoing community effort to
relax, extend, scrutinize, and generalize

strategic

classification as a Stackelberg game: [HMPW '16]

- Players: [1st] **Learner** [2nd] **Users** (dist.)
- Actions: classifier h modify $x \mapsto x^h$
- Payoffs: $\mathbb{E}[\mathbb{1}\{h(x^h) = y\}]$ $\mathbb{E}[\mathbb{1}\{h(x^h) = 1\}]$
- Best response: $x^h = \Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

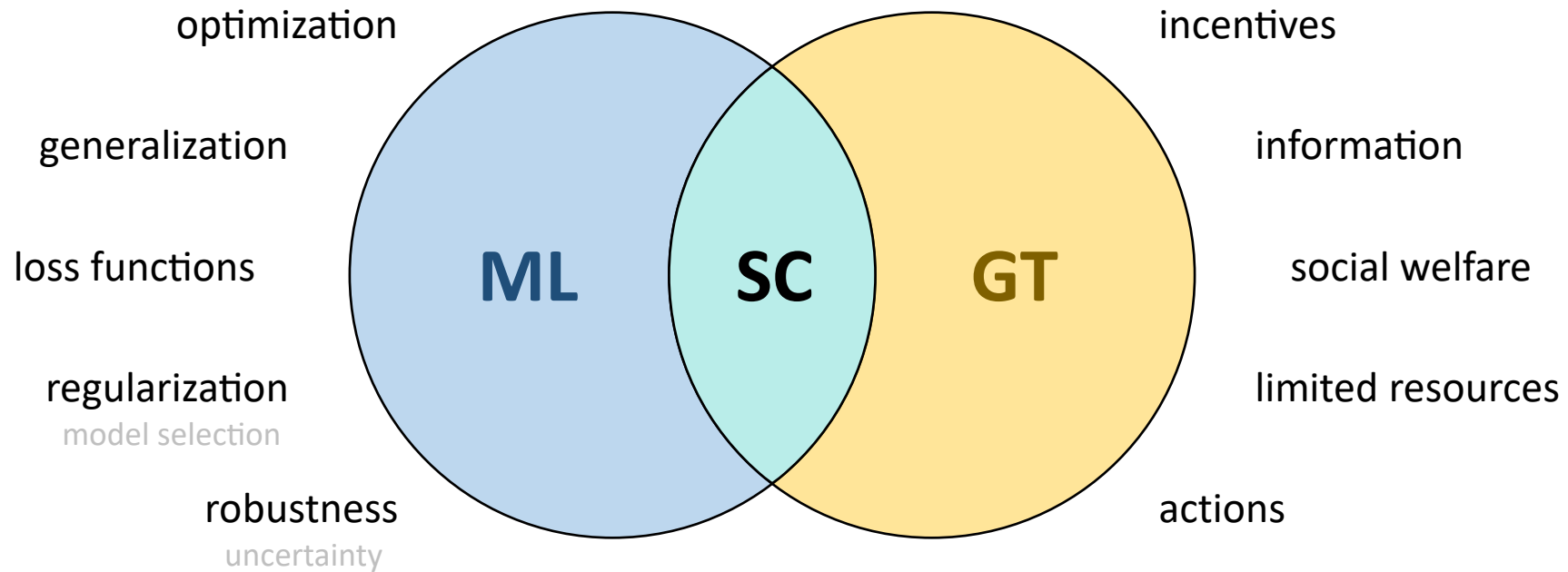
- **Solve equilibrium** \Leftrightarrow **solve learning**
- Holds in idealized setting; trickier as becomes more realistic (finite data, partial information, weaker assumptions, ...)
- **Still:** SC = fundamental ML task + basic economic questions



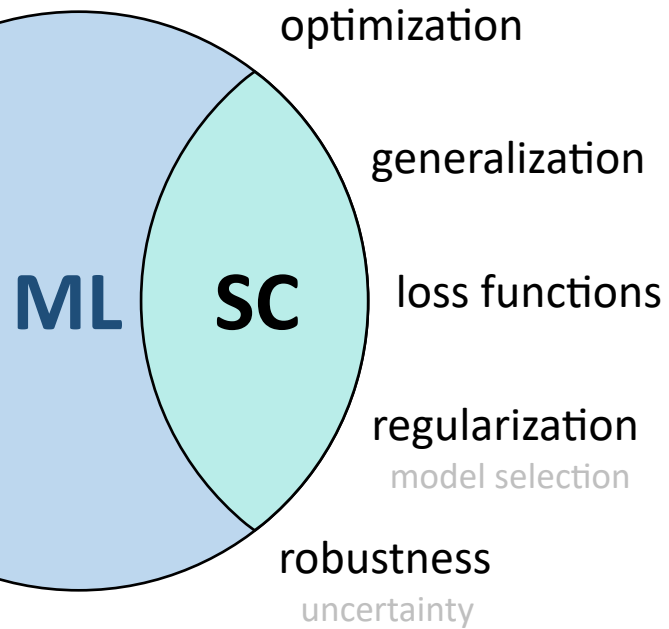
Users

play order is **crucial modeling choice** – choose with care! [NGTR'21, [ZMSJ'21]

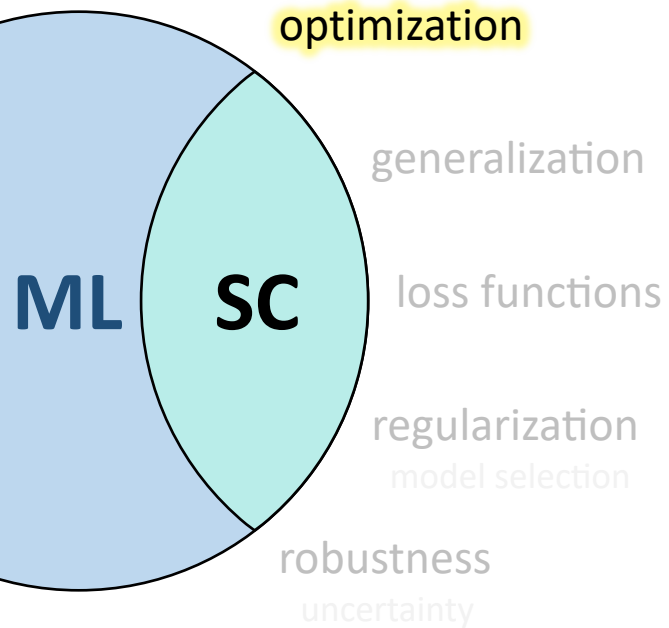
strategic classification as an **interface**
between **machine learning** and **game theory**:



revisit **old questions** + tackle **new ones**



Learning aspects of strategic classification



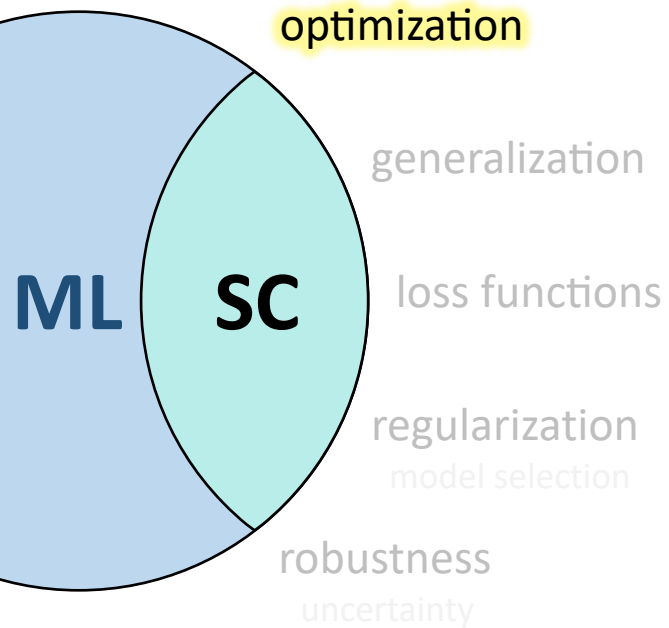
learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

$$\text{s.t. } \Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

nasty nested min-argmax problem!

ask: how to optimize objective?

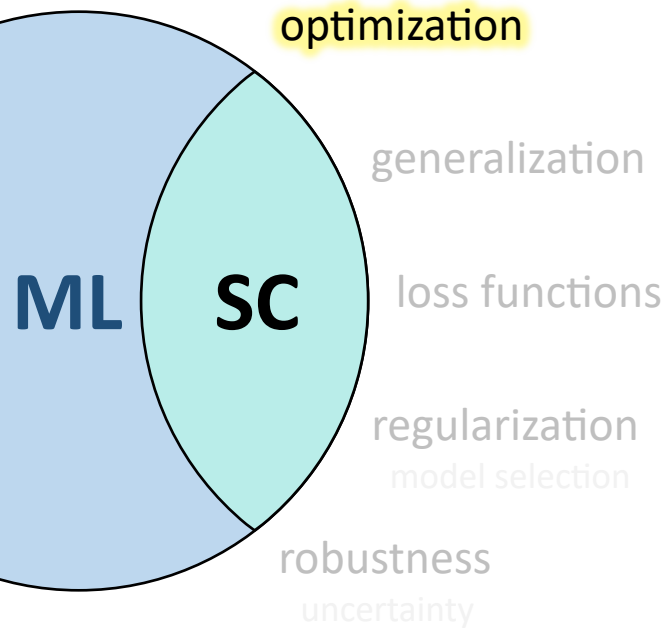


learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

ask: how to optimize objective?

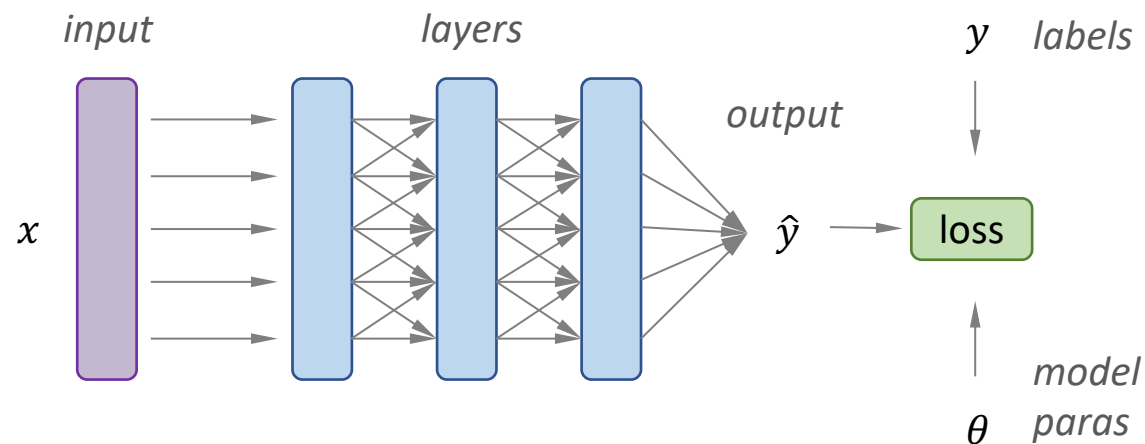


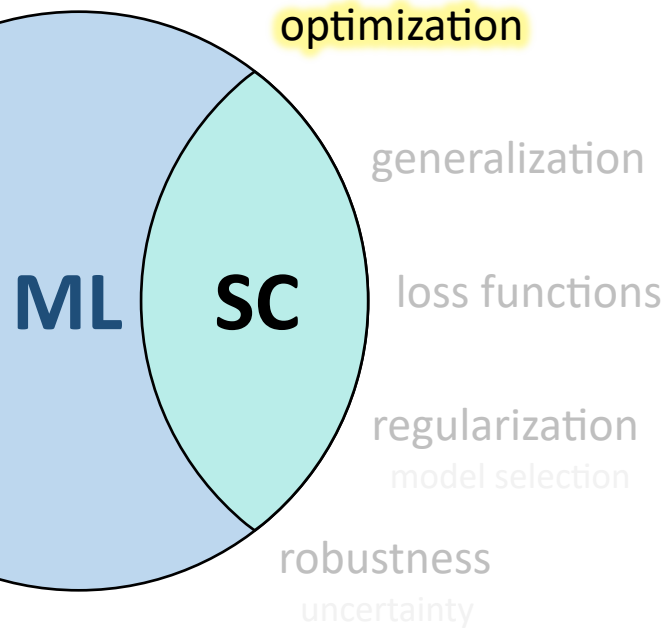
learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

ask: how to optimize objective?



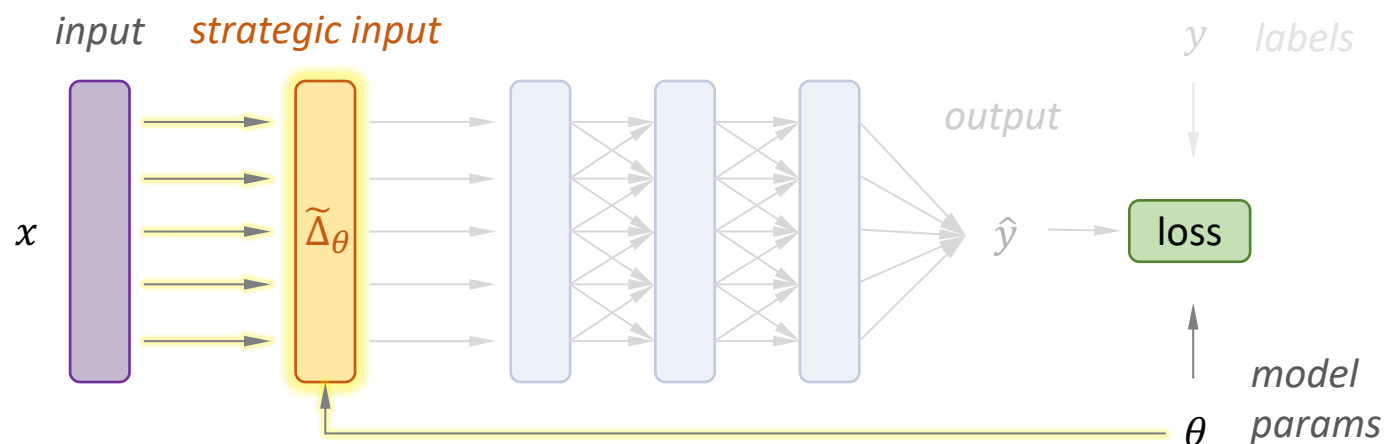


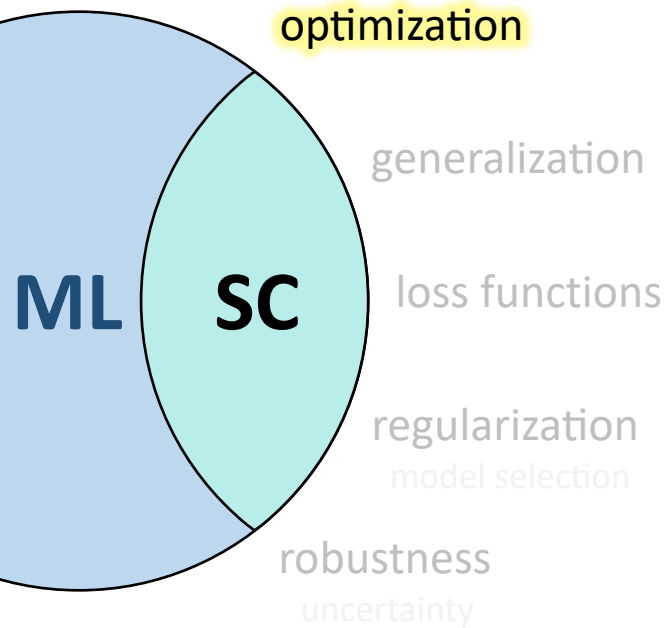
learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

ask: how to optimize objective?





learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

- For **common case** where:

➤ $h(x) = w^\top x$

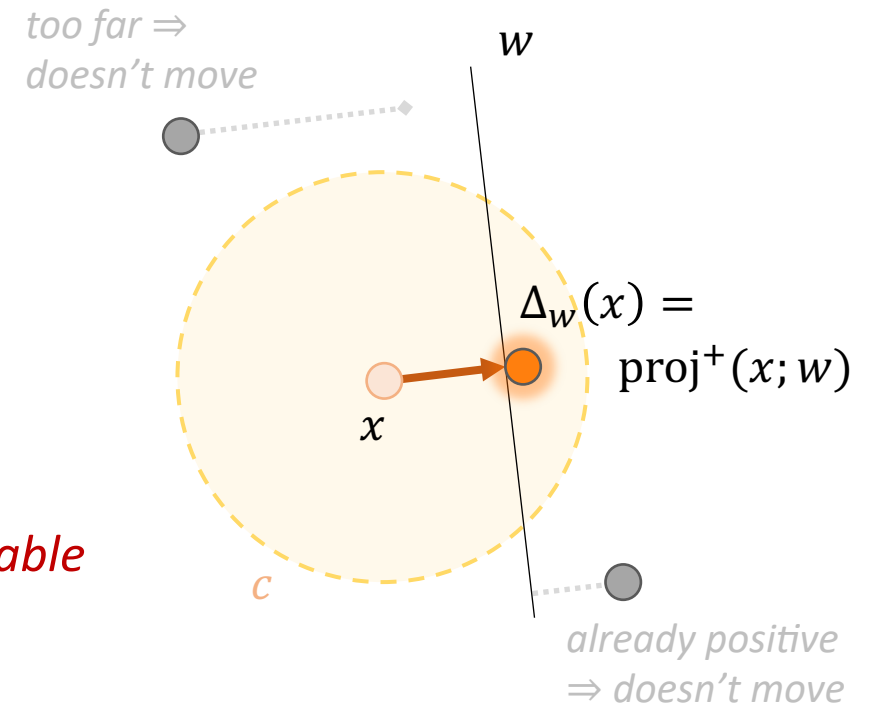
➤ $c(x, x') = \|x' - x\|_2$ (or squared, or PSD)

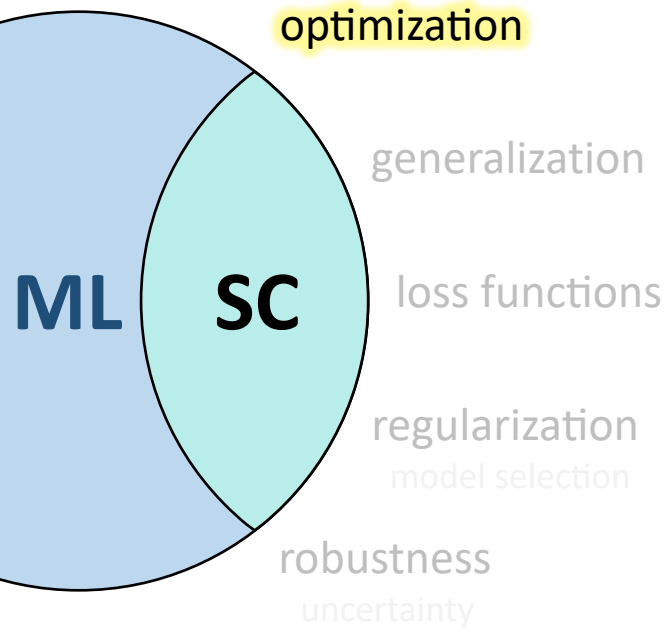
- Admits **simple closed-form solution**:

$$\Delta_w(x) = \begin{cases} x & w^\top x \geq 0 \text{ or } \operatorname{dist}(x; w) > 2 \\ \operatorname{proj}^+(x; w) & \text{o. w.} \end{cases}$$

$= x - \min \left\{ 0, \frac{w^\top x + b}{\|w\|_2^2} \right\}$ *differentiable!*

- Just replace hard-if with soft-if (e.g., sigmoid)





learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

• **Otherwise, when:**

➤ $h(x) = w^\top \phi(x) + \psi(x)$
(for some non-linear ϕ, ψ)

➤ Δ applies to $z = \phi(x)$

➤ c is convex (in z)

• Then can use **plato**: [LR ICML21]

implements Δ as **concave optimization layer** [AABBDK'19]

• **Code:** <https://plato.codes/>

e.g., if Δ is LP:

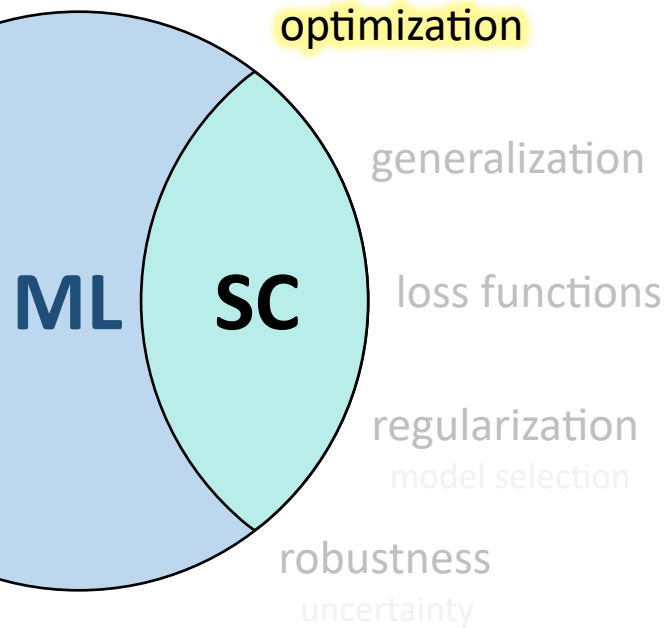
$$z^* = \operatorname{argmax}_z x^\top A z$$

$$\text{s.t. } B z \leq 0$$



$$z^* = h(x; \underbrace{A, B}_{=\theta})$$

can differentiate!



varied costs:
(concave in x, z)

learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

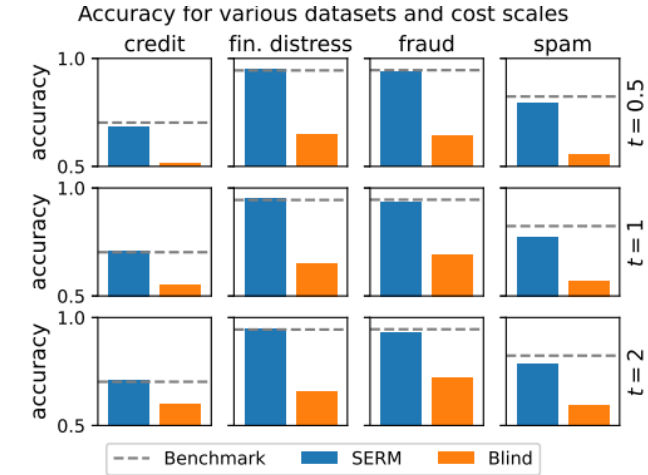
• **Otherwise, when:**

- $h(x) = w^\top \phi(x) + \psi(x)$
(for some non-linear ϕ, ψ)
- Δ applies to $z = \phi(x)$
- c is convex (in z)

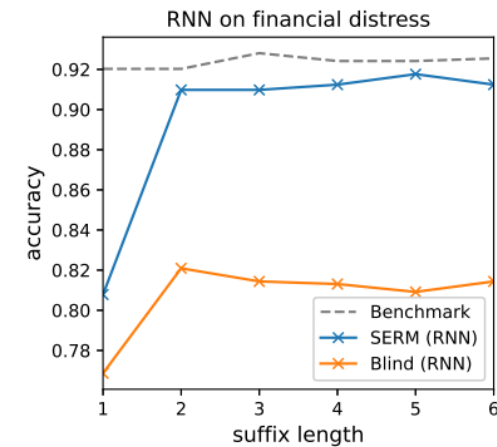
• Then can use **plato**: [LR ICML21]

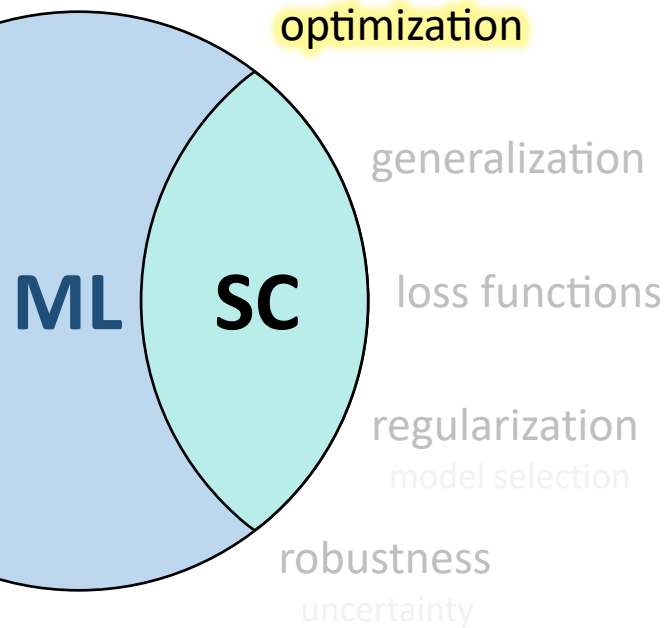
implements Δ as **concave optimization layer** [AABBDK'19]

• **Code:** <https://plato.codes/>



flexible models:
(concave in z)
e.g., certain RNNs





varied costs:
(concave in x, z)

learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

current state of affairs:

s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is differentiable

experiments = semi-syntetic

• Otherwise, when:

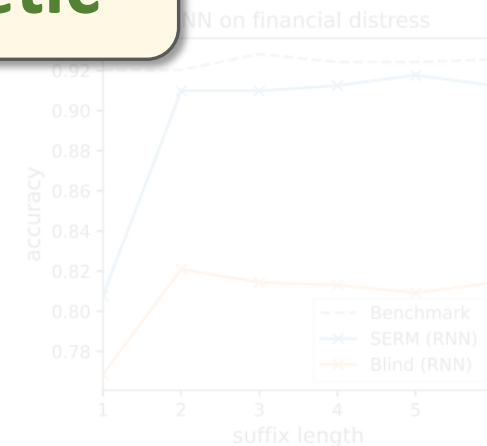
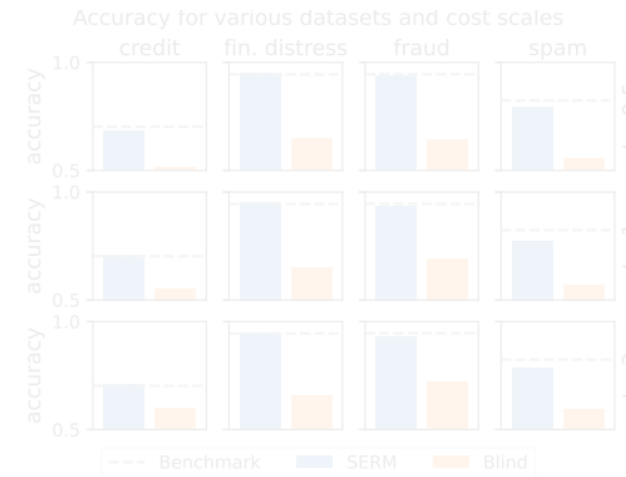
- $h(x) = w^T \phi(x) + \psi(x)$
(for some non-linear ϕ, ψ)
- Δ applies to $z = \phi(x)$
- c is convex (in z)

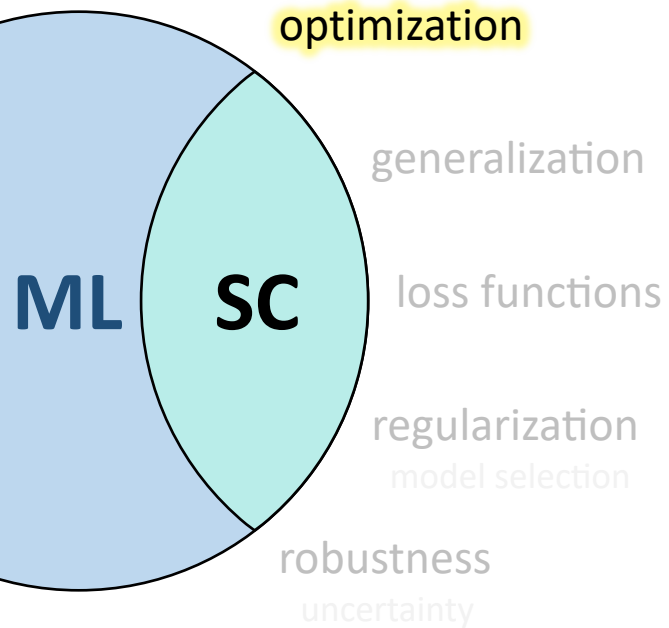
(concave in z)
e.g., certain RNNs

• Then can use **plato**: [LR ICML21]

implements Δ as **concave optimization layer** [AABBDK'19]

• Code: <https://plato.codes/>



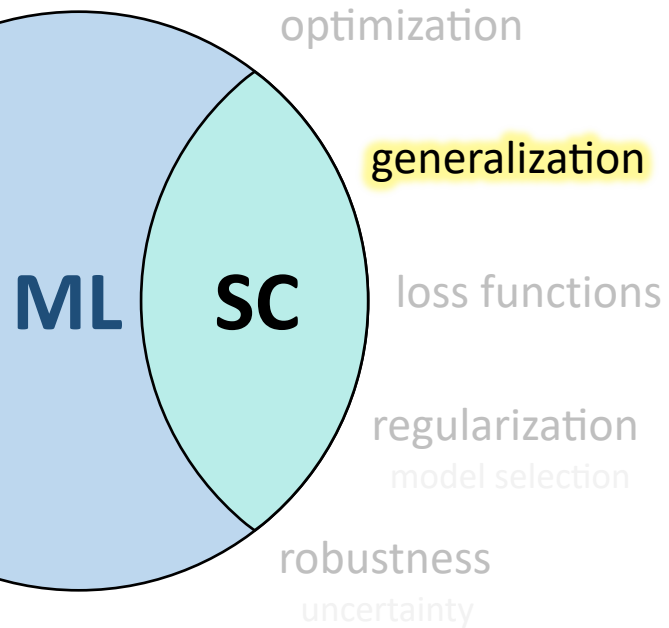


learning objective:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\tilde{\Delta}_h(x_i)))$$

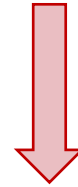
s.t. $\tilde{\Delta}_h(x) \approx \Delta_h(x)$ and is *differentiable*

- **Otherwise – uncharted territory**
- **Idea:** borrow methods from **adversarial learning** literature (e.g., FGSM [GSS'15] or PGD [MMSTV'18])
- Essentially, optimize objective by alternating between:
 - fixing features x^h and updating θ
 - fixing parameters θ and updating x^h
- Technically possible – but hasn't been done yet in strategic learning
- More on **strategic** ↔ **adversarial** connection to follow!



empirical loss:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

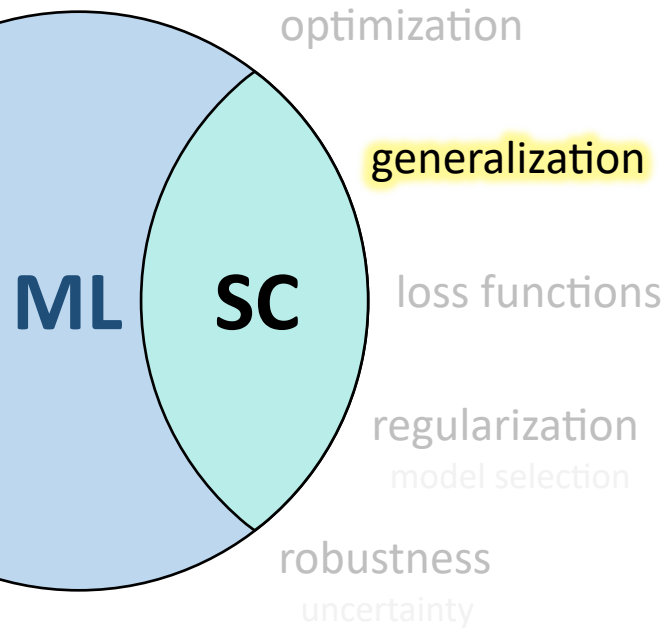


generalization

expected loss:

$$\operatorname{argmin}_h \mathbb{E} \left[\ell(y, h(\Delta_h(x))) \right]$$

ask: how does behavior affect generalization?

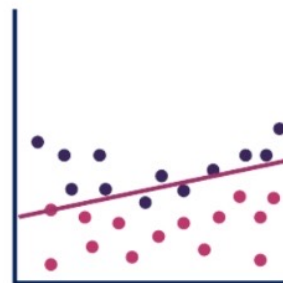


empirical loss:

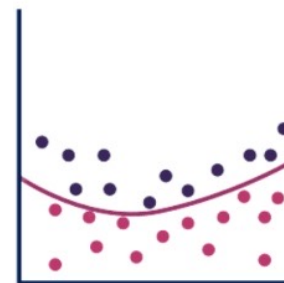
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

- **Q** – will strategic behavior:
 1. *increase* overfitting?
 2. *reduce* overfitting?
 3. *make no difference*?
- **Rephrase:** how does behavior affect sample complexity?

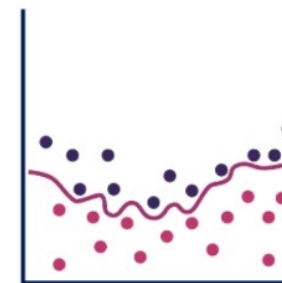
ask: how does behavior affect generalization?



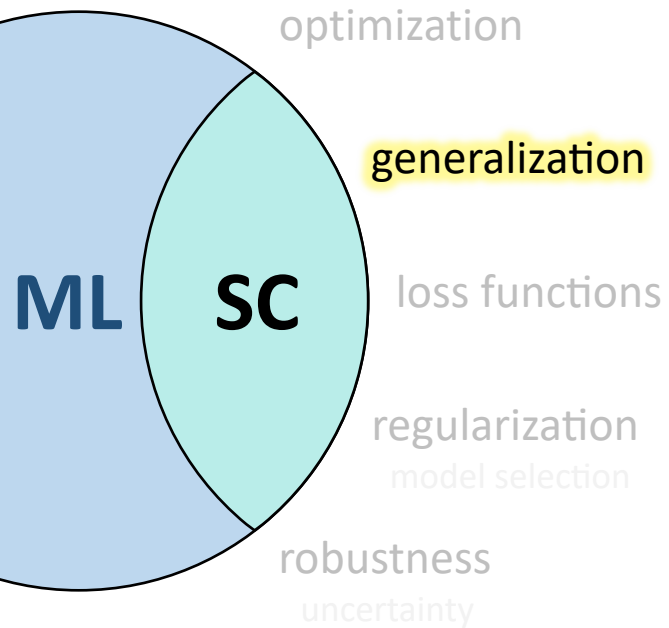
Underfitting



Balanced



Overfitting

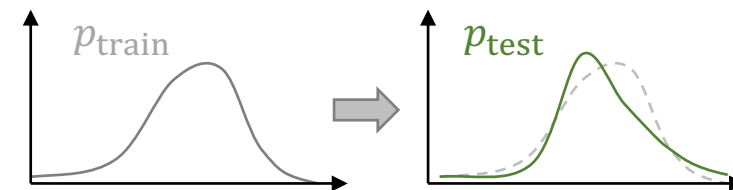


empirical loss:

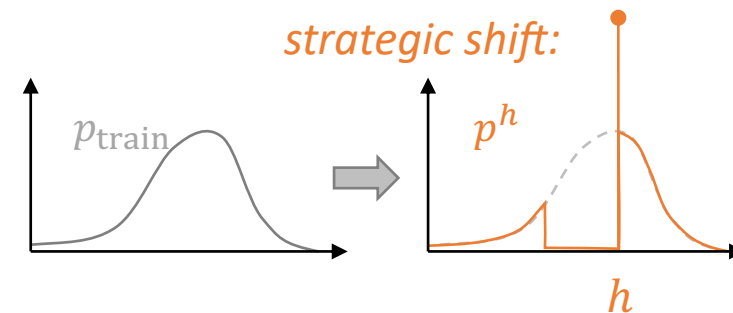
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

- SC = **model-dependent distribution shift**
- In **typical distribution shift**, p_{test} is assumed to be “close” to p_{train} (e.g., in ball)
- Contrarily, in **strategic shift**:
 1. only points in “band” before h move
 2. entire region moves on decision boundary
 3. moving region determined by choice of h

typical shift:



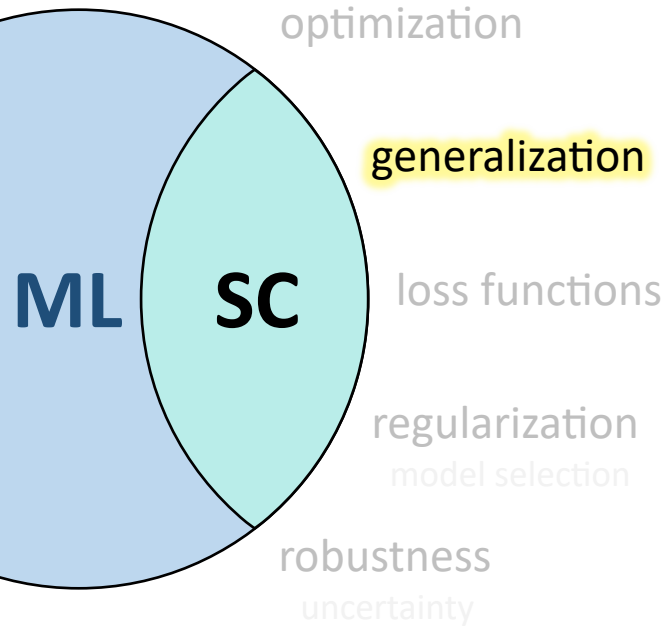
strategic shift:



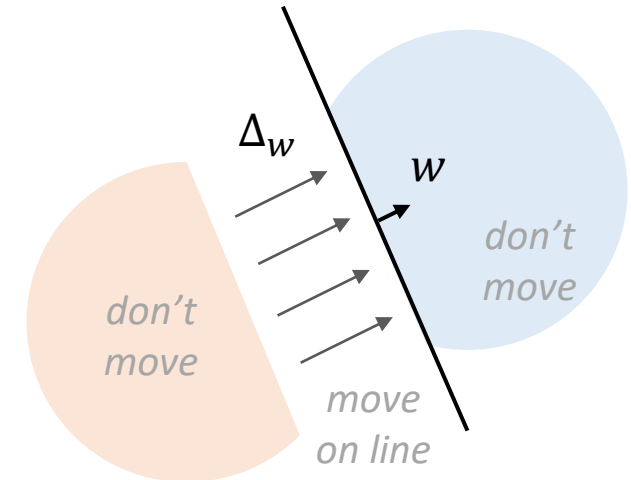
ask: how does behavior affect generalization?

empirical loss:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$



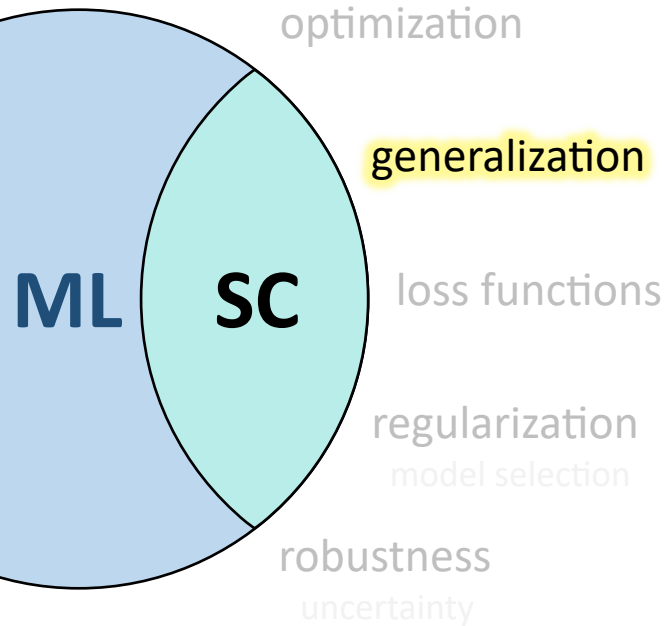
- SC = **model-dependent distribution shift**
- In **typical distribution shift**, p_{test} is assumed to be “close” to p_{train} (e.g., in ball)
- Contrarily, in **strategic shift**:
 1. only points in “band” before h move
 2. entire region moves on decision boundary
 3. moving region determined by choice of h



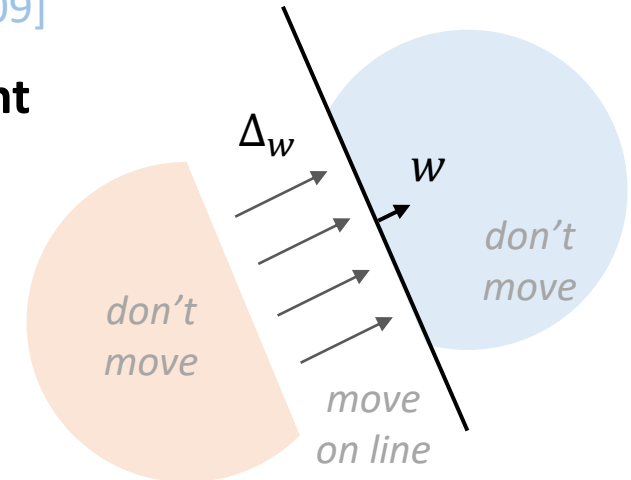
ask: how does behavior affect generalization?

empirical loss:

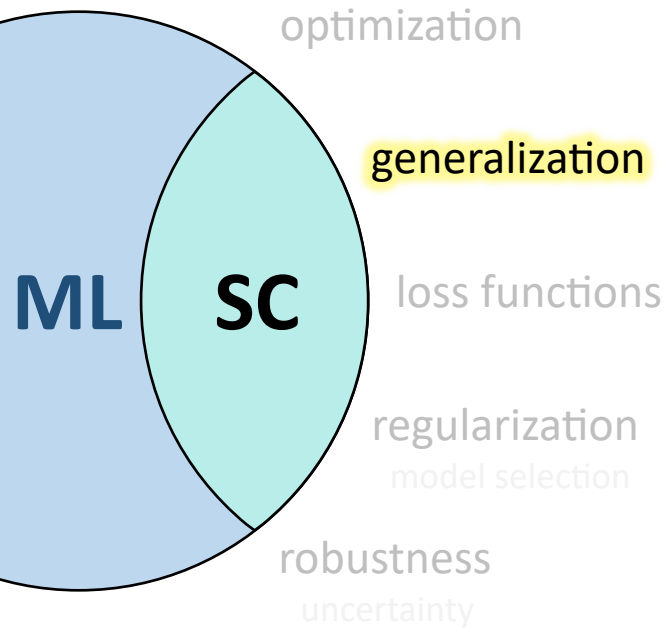
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$



- Generalization theory typically relies on **discrepancy measures** $d(p_{\text{train}}, p_{\text{test}})$ [MMR ICML09]
- \Rightarrow bounds are **shift** (and so **dsitribution**) **dependent**
- Interestingly, strategic shifts admit **distribution-independent** generalization bounds



ask: how does behavior affect generalization?



empirical loss:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

- **Induced class:** $H_\Delta = \{h(\Delta_h(x)) : h \in H\}$
- **Strategic VC:** $SVC(H) = VC(H_\Delta)$

• **Result:** for standard setting, recover non-strategic bounds (almost!)

- But – **cost form matters!** [SVXY'23] show:

- **instance-invariant costs:**

$$c(x - x') \Rightarrow SVC \approx VC \text{ (for linear } h)$$

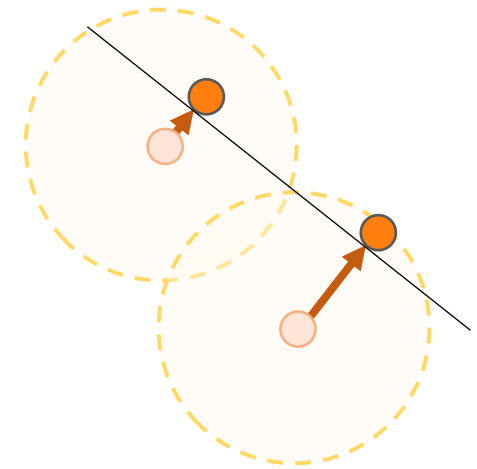
i.e., learning is **not harder**

- **instance-wise costs:** *=individualized*

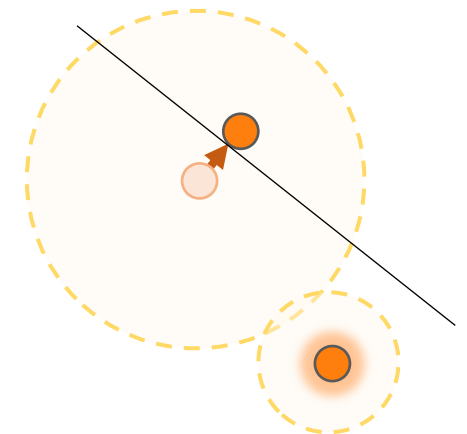
$$c_x(x') \Rightarrow \text{unlearnable! (in the worst case)}$$

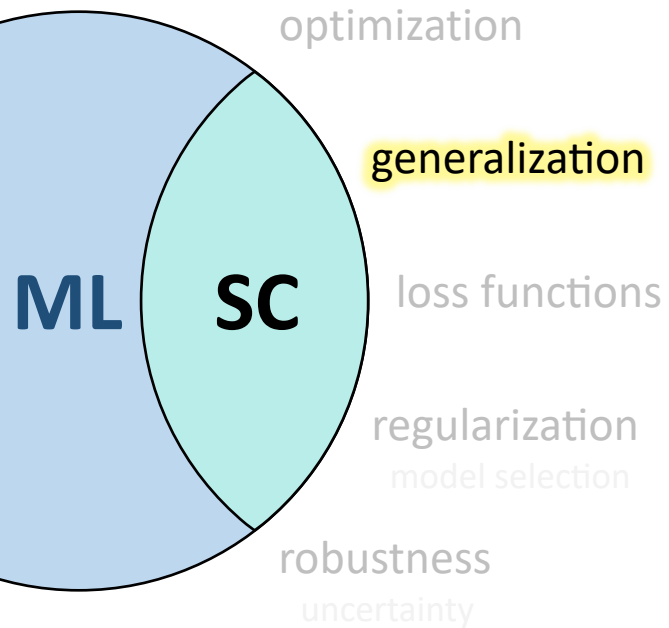
i.e., learning is **impossible**

instance-invariant:



instance-wise:





empirical loss:

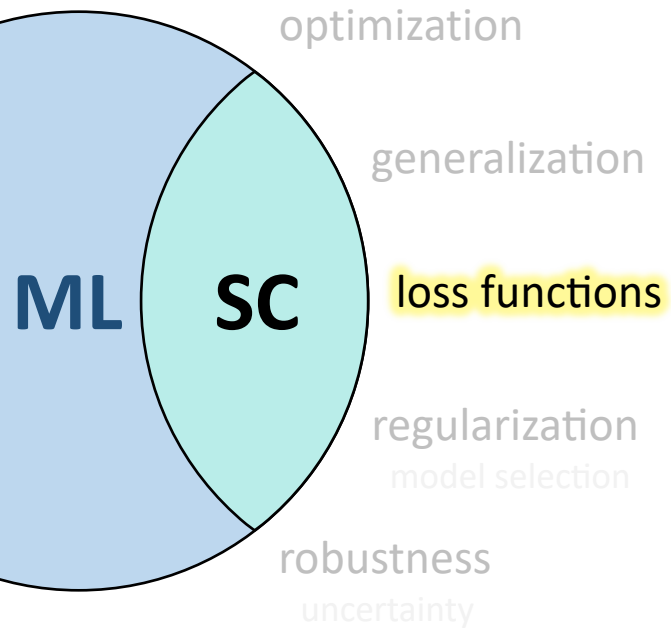
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

- Also have data-dependent Rademacher bounds:

$$\mathcal{L}_{0/1} \leq \mathcal{L}_{s\text{-hinge}} + \frac{4r\|\hat{w}\|}{\sqrt{m}} + (1 + 2(\overset{= \max\|x\|}{r+2})\|w\|) \sqrt{\frac{2 \ln(4\|\hat{w}\|/\delta)}{m}}$$

only difference from non-strategic

- Notice behavior **just adds constant** to scale
- Applies to more general settings (to follow)
- **Also:** regret analysis for online strategic classification



tractable proxy loss: (e.g., hinge, log-loss, ...)

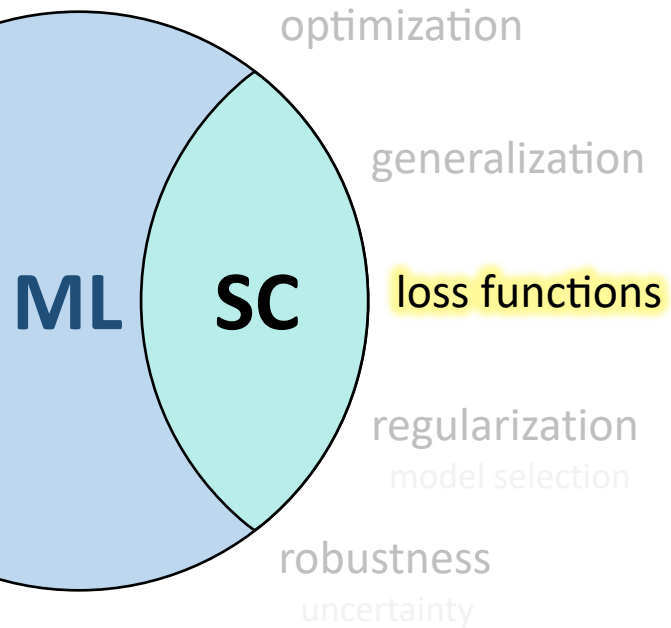
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

↓ *surrogate*

true 0-1 loss:

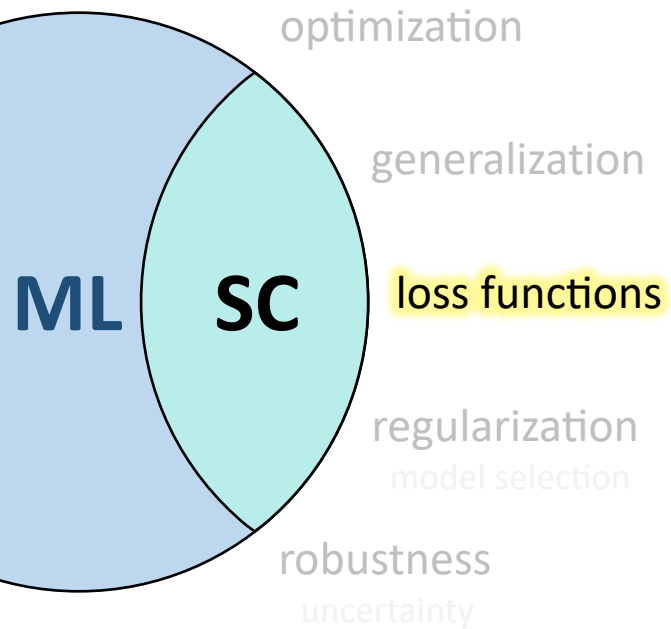
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y_i \neq h(\Delta_h(x_i))\}$$

ask: can we just use conventional proxies?

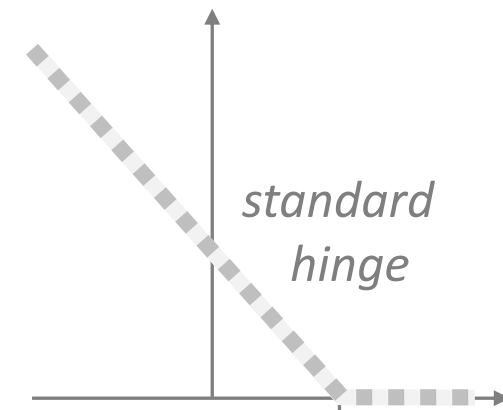
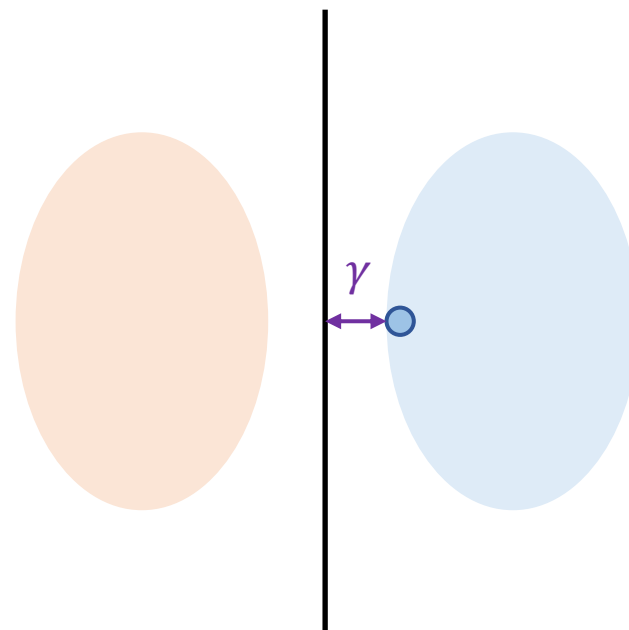


standard hinge: $\max\{0, 1 - yw^T x\}$

ask: can we just use conventional proxies?



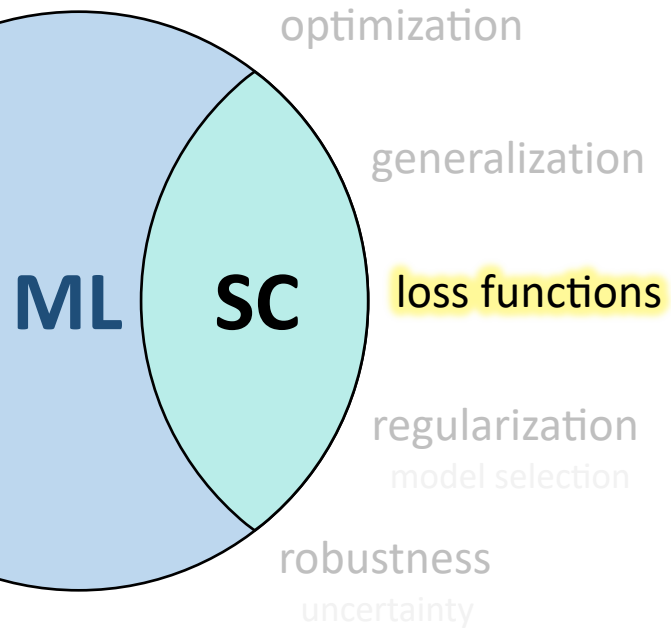
standard hinge: $\max\{0, 1 - yw^T x\}$



max-margin classifier

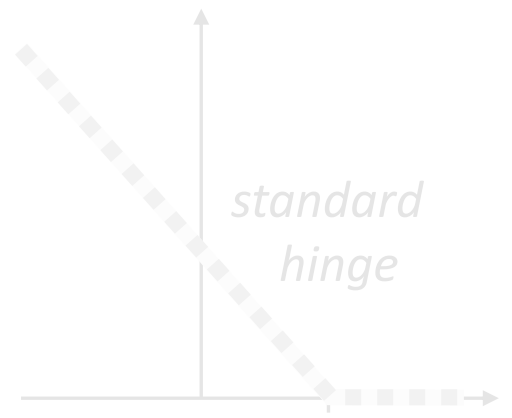
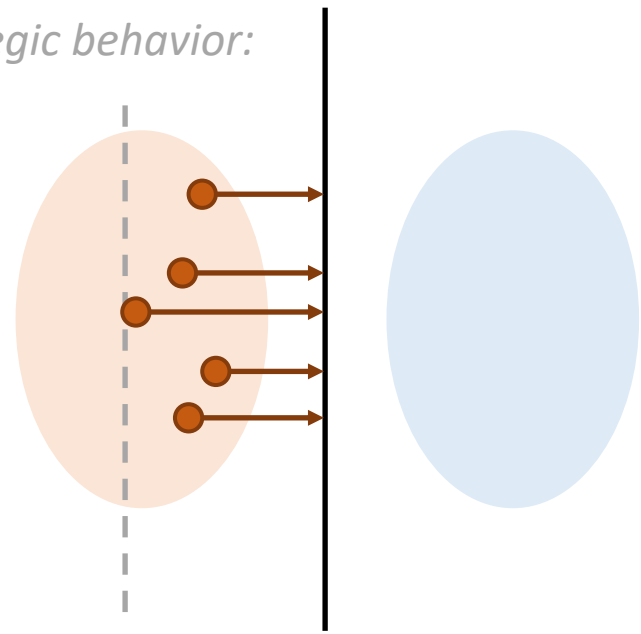
- selection criterion
- good generalization
- tractable

ask: can we just use conventional proxies?



naïve hinge: $\max\{0, 1 - yw^T \Delta_h(x)\}$

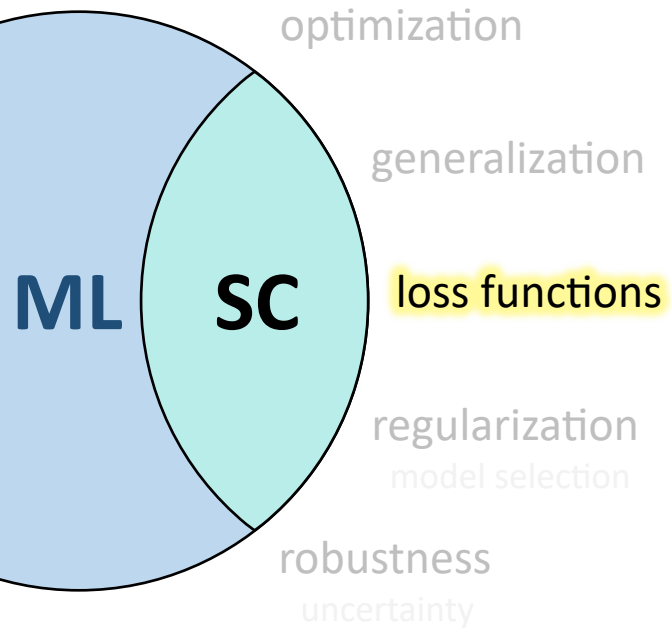
- strategic behavior:



max-margin classifier

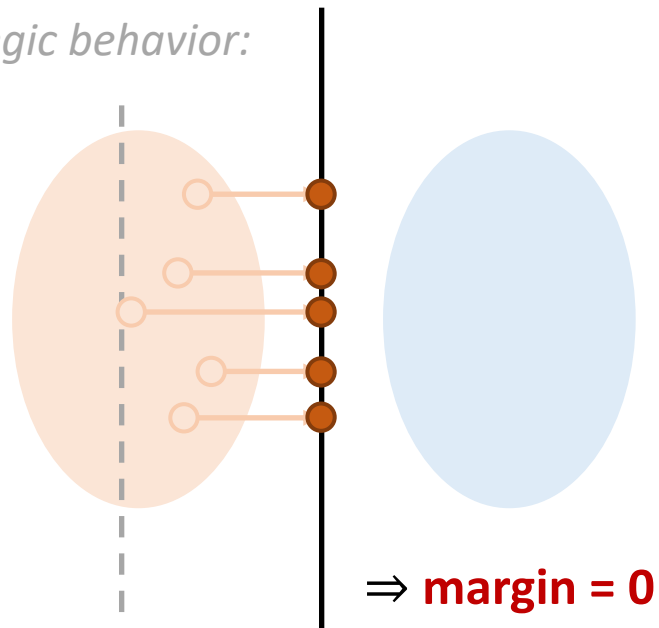
- selection criterion
- good generalization

ask: can we just use conventional proxies?



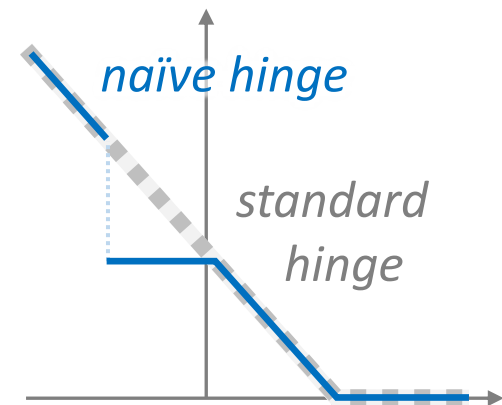
naïve hinge: $\max\{0, 1 - yw^T \Delta_h(x)\}$

- strategic behavior:

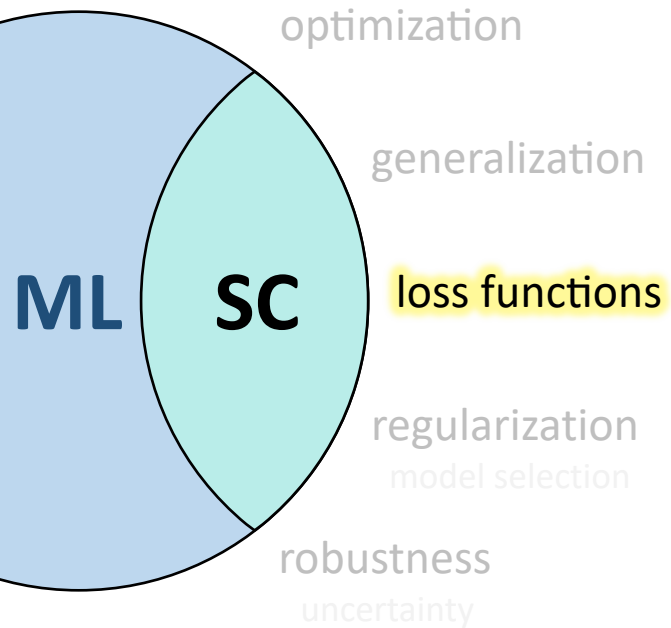


naïve max-margin classifier

- vacuous criterion
- unclear if generalizes

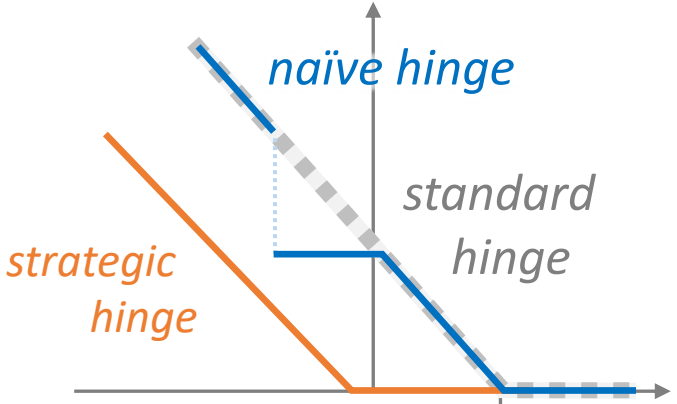
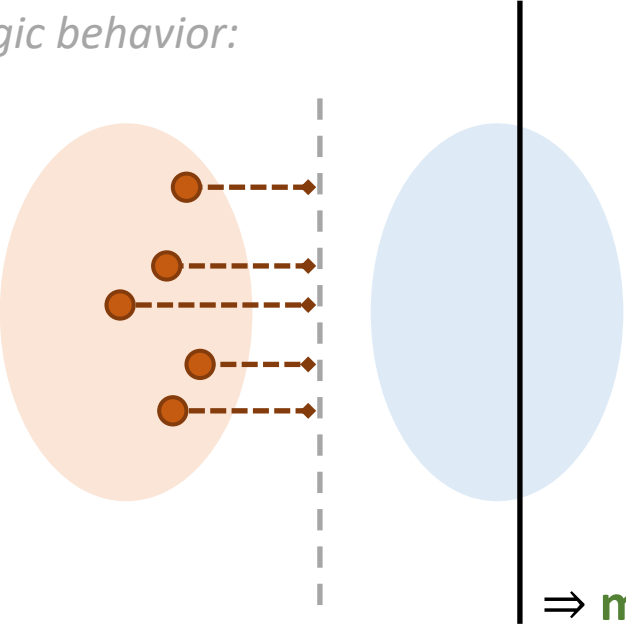


ask: can we just use conventional proxies?



strategic hinge: $\max\{0, 1 - yw^T x - 2\|w\|\}$

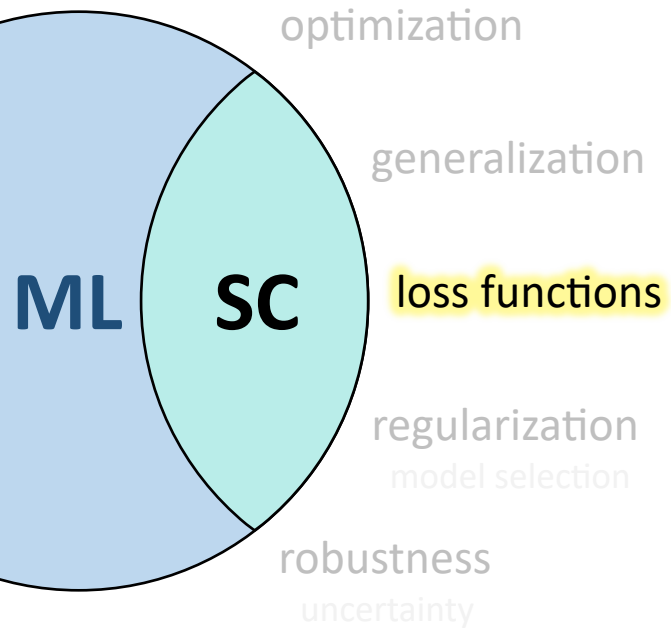
- *strategic behavior:*



strategic max-margin classifier

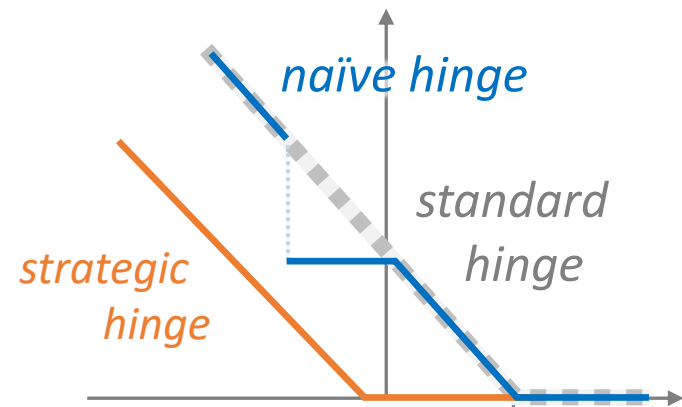
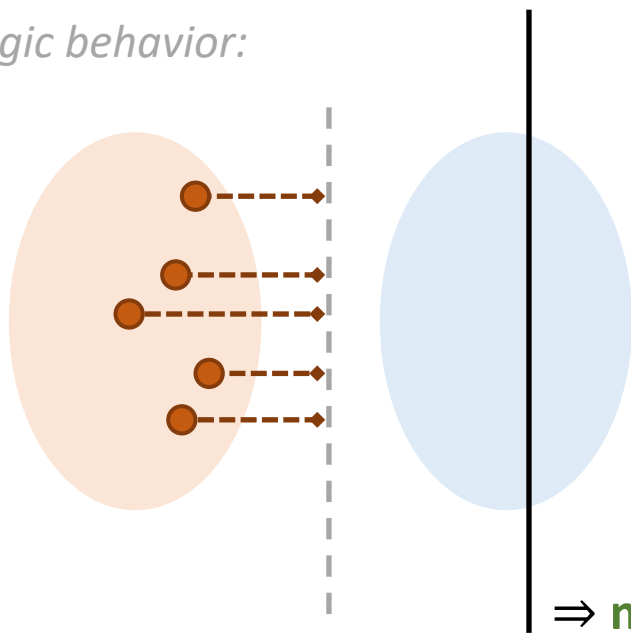
- *regain selection criterion*
- *comparable generalization*
- *reasonably tractable*

ask: can we just use conventional proxies?



strategic hinge: $\max\{0, 1 - yw^T x - 2\|w\|\}$

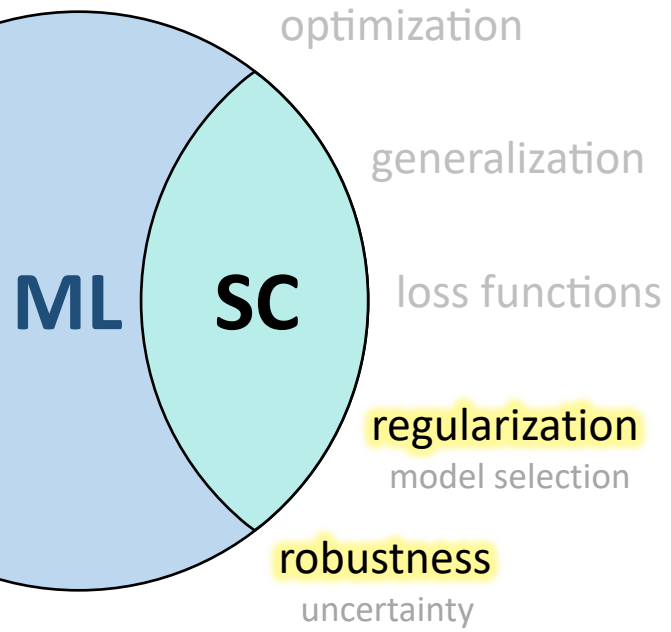
- *strategic behavior:*



strategic max-margin classifier

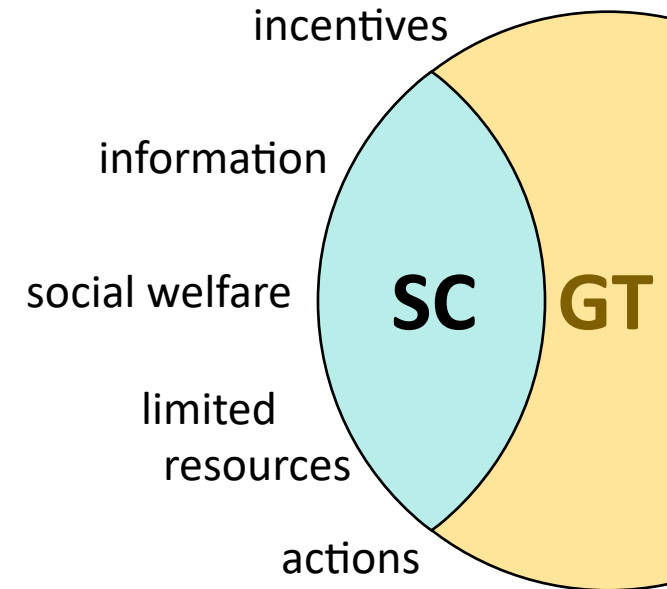
- *regain selection criterion*
- *comparable generalization*
- *reasonably tractable*

conclusion: strategic robustness requires **rethinking fundamental learning concepts**



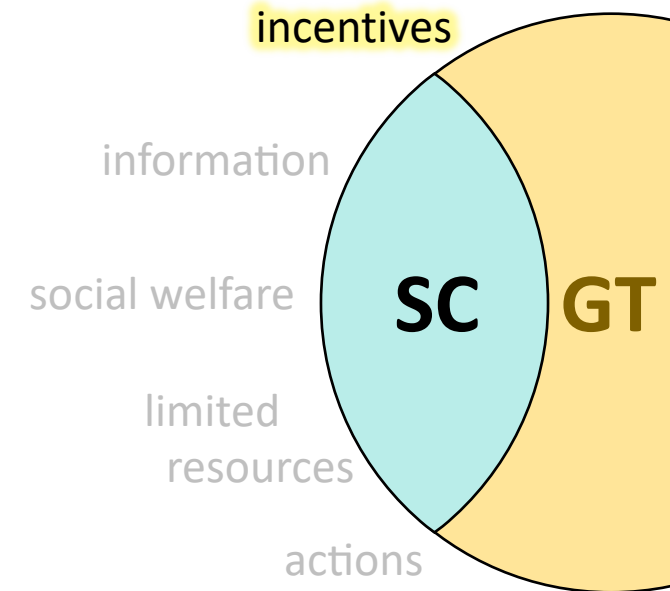
} *to follow!*

Economic aspects of strategic classification



standard SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

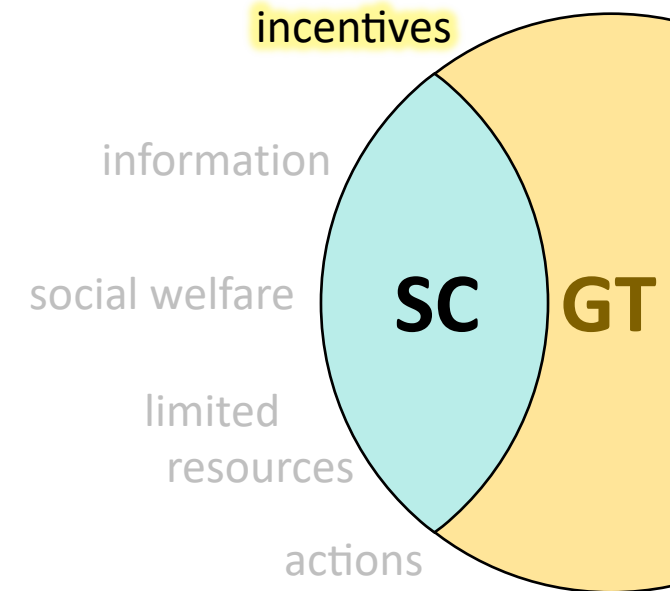


ask: what else could users want?

standard SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x') - c(x, x')$$

└──────────┬──────────→ = $\begin{cases} +1 & \hat{y} = +1 \\ -1 & \hat{y} = -1 \end{cases}$



ask: what else could users want?

standard SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x') - c(x, x')$$

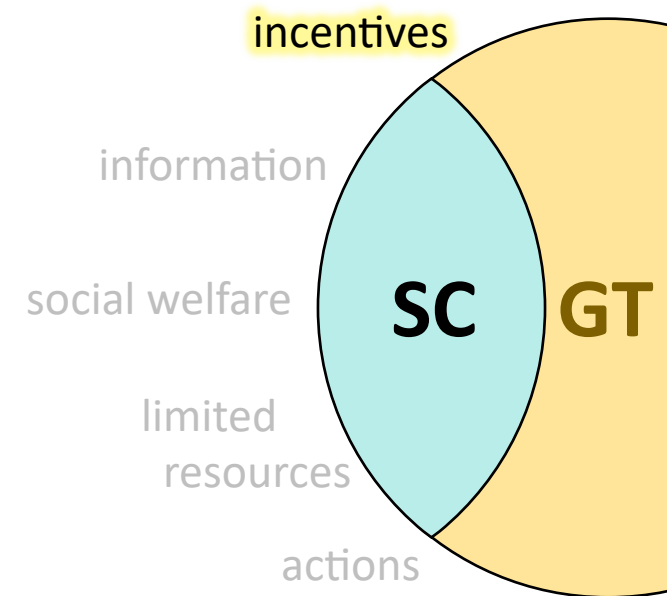
└──────────┬──────────┘

$$= \begin{cases} +1 & \hat{y} = +1 \\ -1 & \hat{y} = -1 \end{cases}$$

generalized SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x') - c(x, x')$$

1. arbitrary **utility function**



ask: what else could users want?

standard SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x') - c(x, x')$$

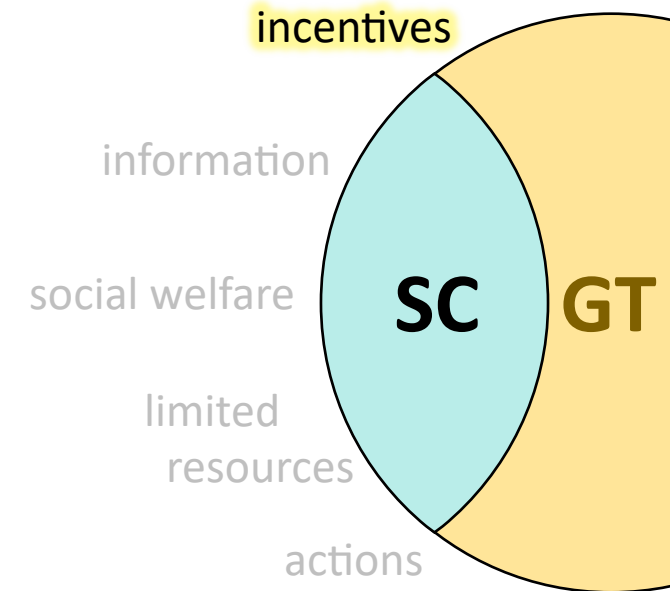
↘

$$= \begin{cases} +1 & \hat{y} = +1 \\ -1 & \hat{y} = -1 \end{cases}$$

generalized SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x'; z) - c(x, x')$$

1. arbitrary **utility function**
2. can depend on **private information**



ask: what else could users want?

standard SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} u(x') - c(x, x')$$

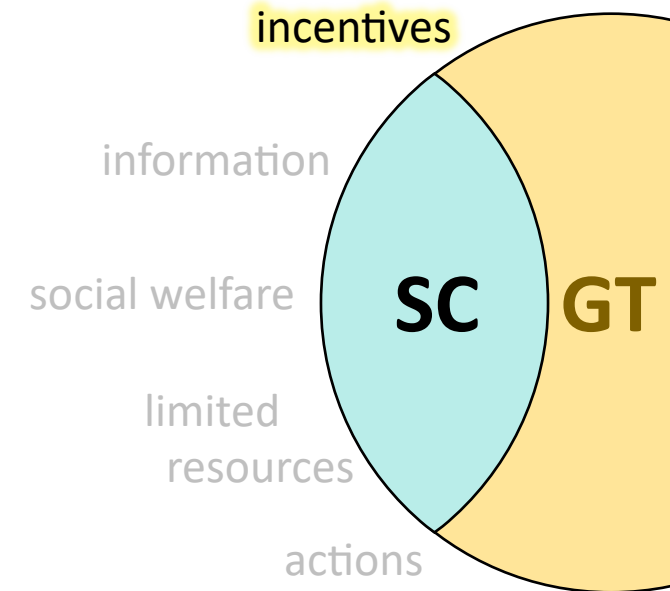
↘

$$= \begin{cases} +1 & \hat{y} = +1 \\ -1 & \hat{y} = -1 \end{cases}$$

generalized SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} \tilde{u}(x'; z) - c(x, x')$$

1. arbitrary **utility function**
2. can depend on **private information**
3. act on **perceived utility** (\neq true utility)



ask: what else could users want?

generalized SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} \tilde{u}(x'; z) - c(x, x')$$

- **Q:** how to learn?
- **A:** generalize strategic margins and hinge!

➤ standard hinge:

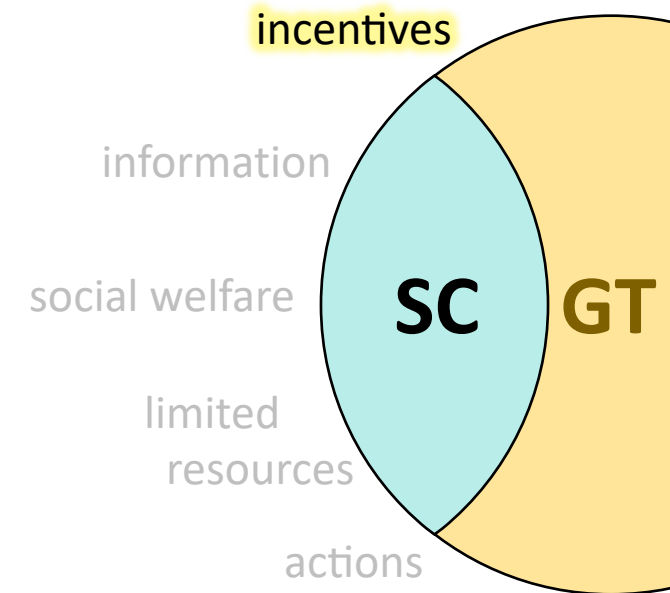
$$\begin{aligned} & \max\{0, 1 - yw^\top x\} \\ & = \max\{0, 1 - \underbrace{\operatorname{sign}(yw^\top x)}_{\text{correctness}} |\underbrace{w^\top x}_{\text{distance}}|\} \end{aligned}$$

➤ naïve hinge:

$$\max\{0, 1 - \operatorname{sign}(yw^\top \Delta_h(x, z)) |w^\top \Delta_h(x, z)|\}$$

➤ generalized strategic hinge: (gs-hinge)

$$\max\{0, 1 - \operatorname{sign}(yw^\top \Delta_h(x, z)) d_\Delta(x, z; w)\}$$



ask: what else could users want?

generalized SC:

$$\Delta_h(x) = \operatorname{argmax}_{x'} \tilde{u}(x'; z) - c(x, x')$$

reinterpretation of “margin”:

➤ distance to nearest x' that *flips label*:

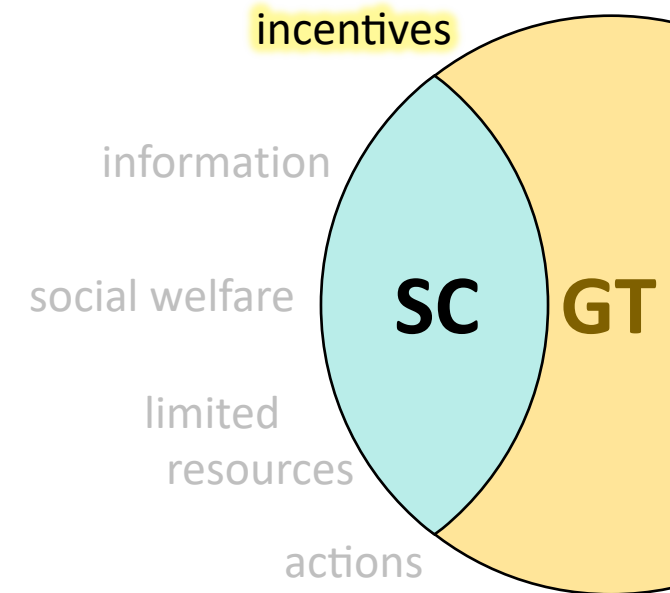
$$d_{\Delta}(x, z; w) = \min_{x'} \frac{\|x - x'\|}{\|w\|} \quad \begin{array}{l} \text{minimal distance} \\ \text{between points} \\ \text{(normalized)} \end{array}$$

flip label (after movement) s.t. $h(\Delta_h(x, z)) \neq h(\Delta_h(x', z))$
subsumes non-strategic case

- admits convenient tractable form for several known special cases

➤ **generalized strategic hinge:** (gs-hinge)

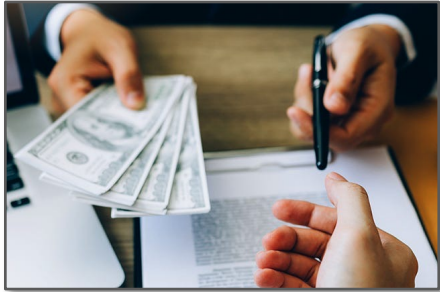
$$\max\{0, 1 - \operatorname{sign}(yw^{\top} \Delta_h(x, z)) d_{\Delta}(x, z; w)\}$$



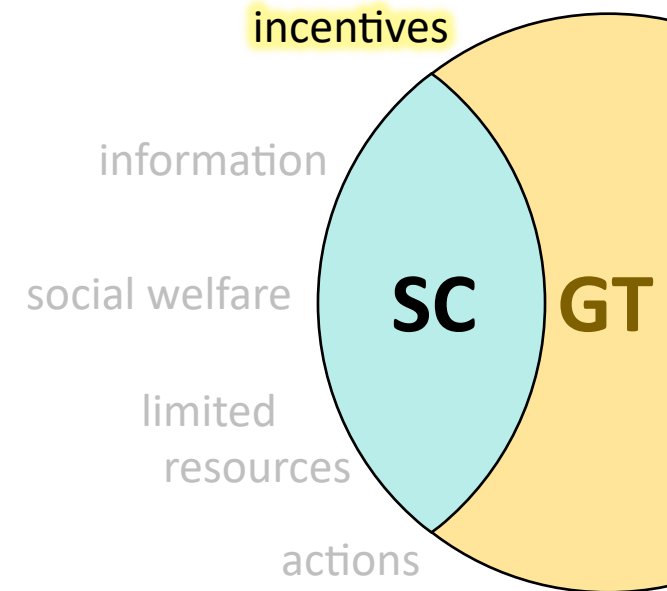
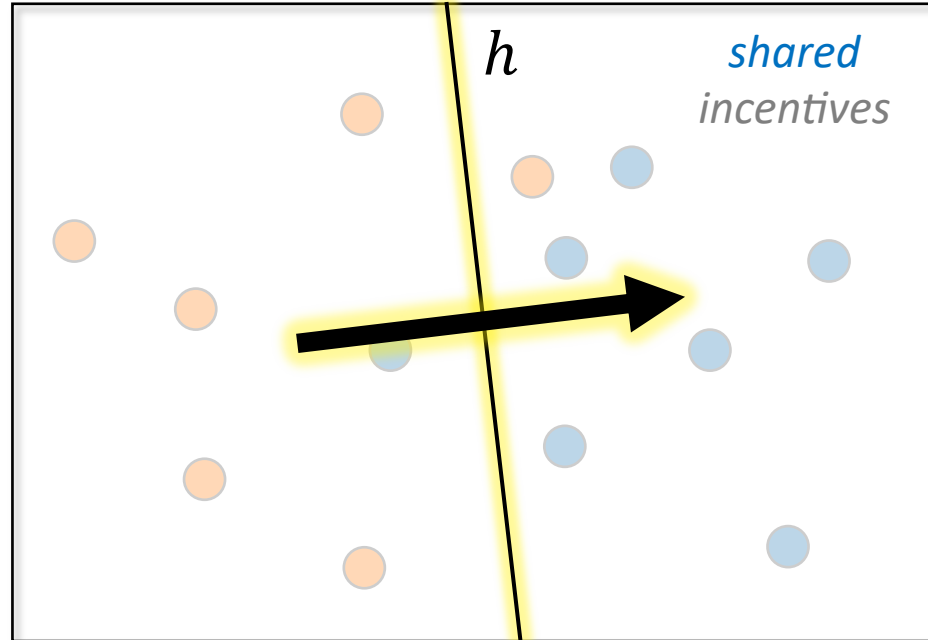
ask: what else could users want?

standard SC:

classification **about** humans

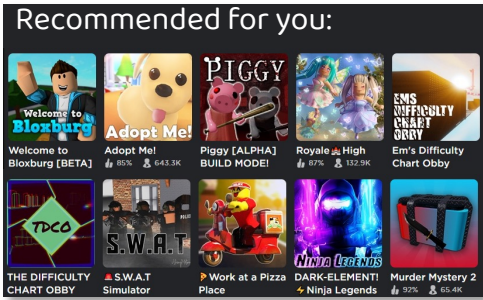


system wants: correct predictions
users want: positive predictions



ask: what else could users want?

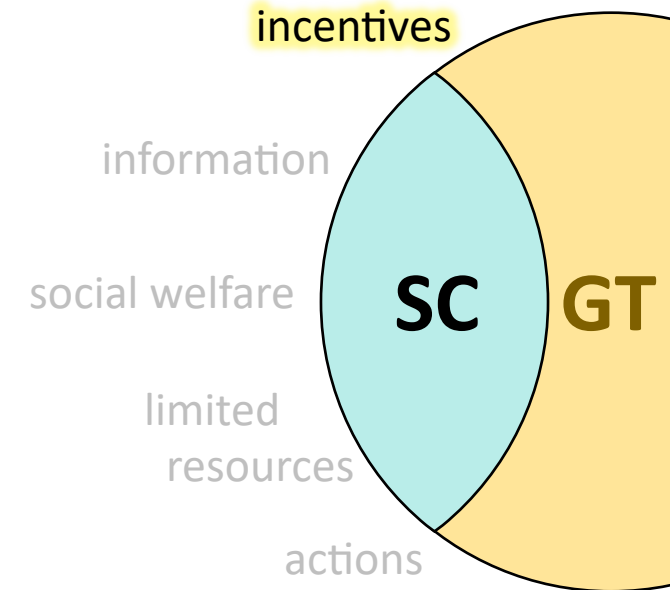
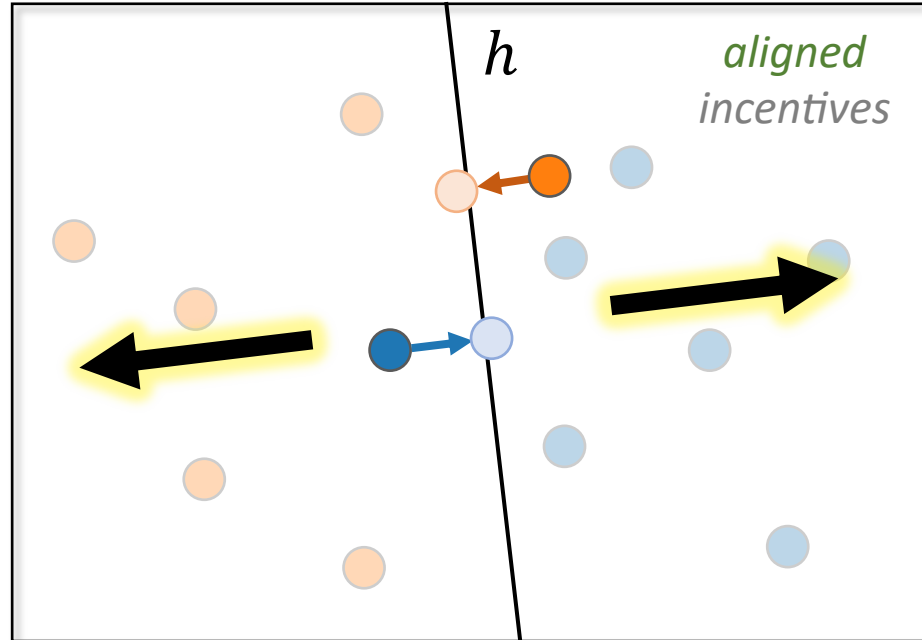
classification *for* humans (as a service)



system wants: *correct predictions*

users want: *correct predictions*

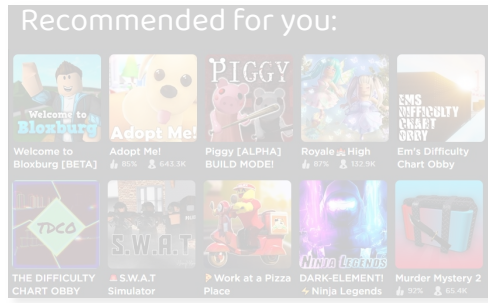
incentive-aligned:



ask: can learning (implicitly) coordinate cooperation?

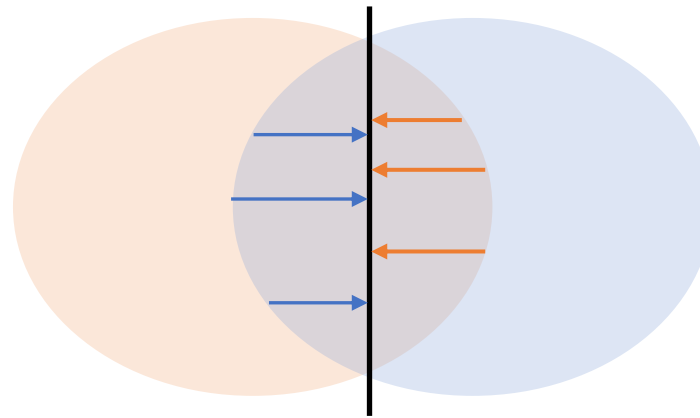
incentive-aligned:

classification *for* humans (as a *service*)

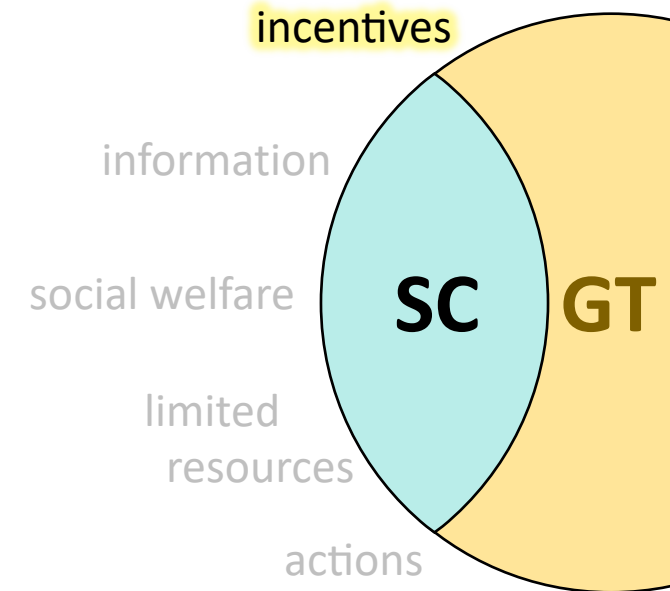


system wants: correct predictions

users want: correct predictions

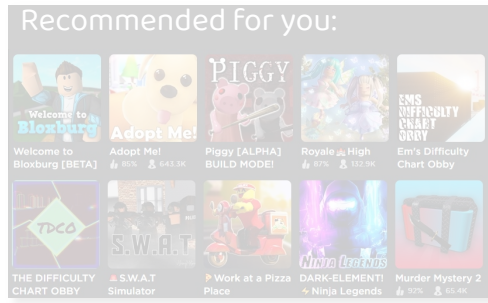


not linearly separable



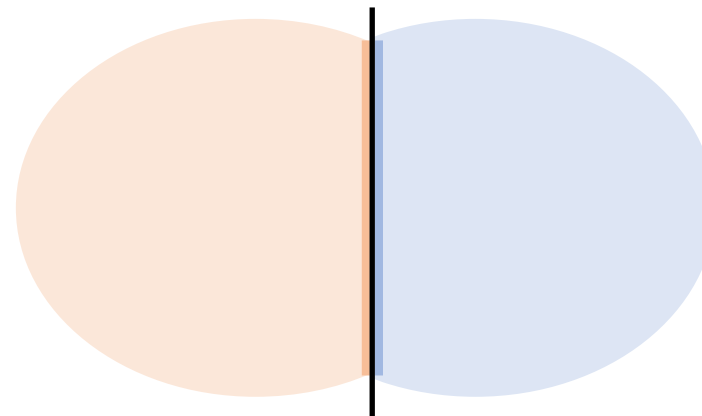
incentive-aligned:

classification *for* humans (as a service)

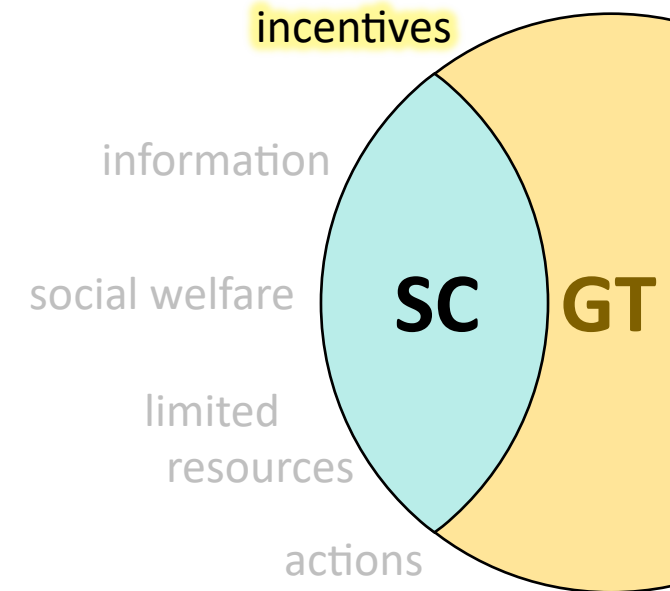


system wants: correct predictions

users want: correct predictions

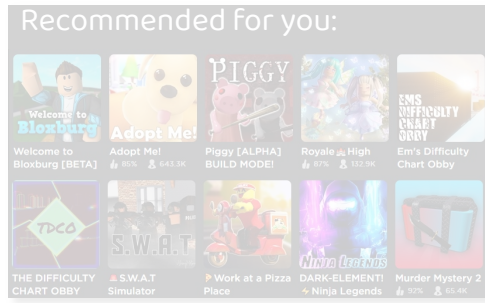


strategically linearly separable



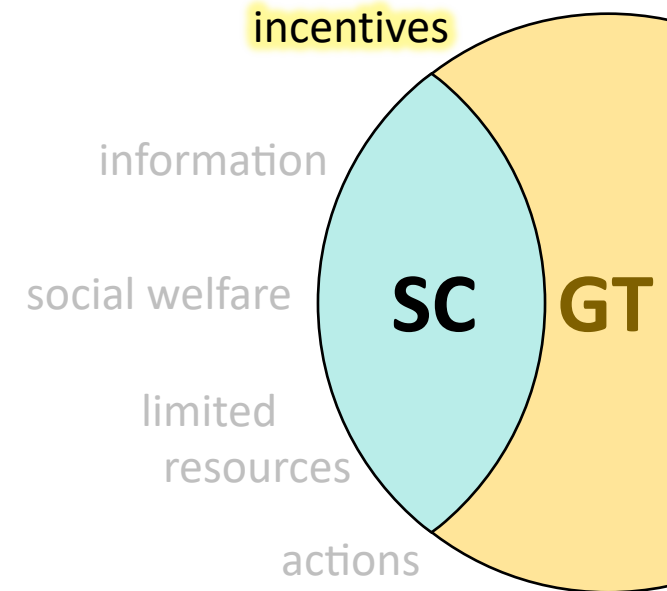
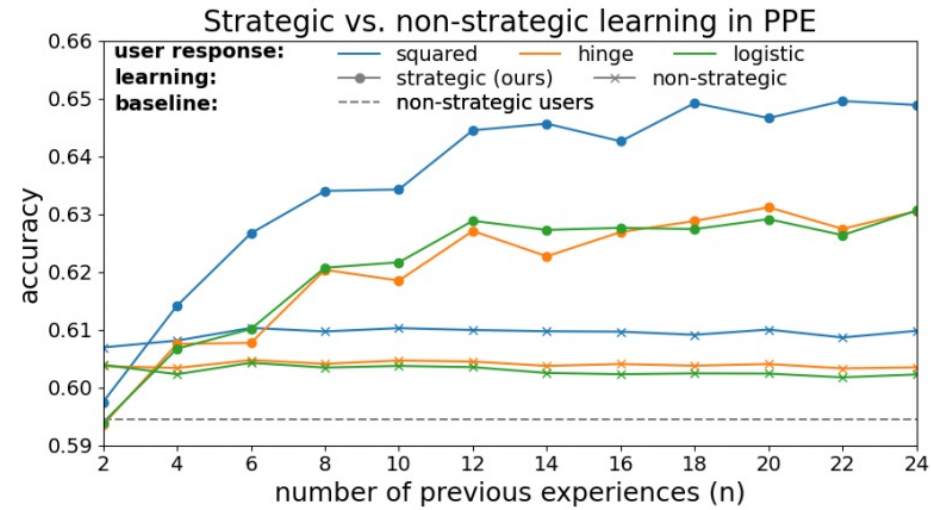
incentive-aligned:

classification *for* humans (as a service)



system wants: correct predictions

users want: correct predictions



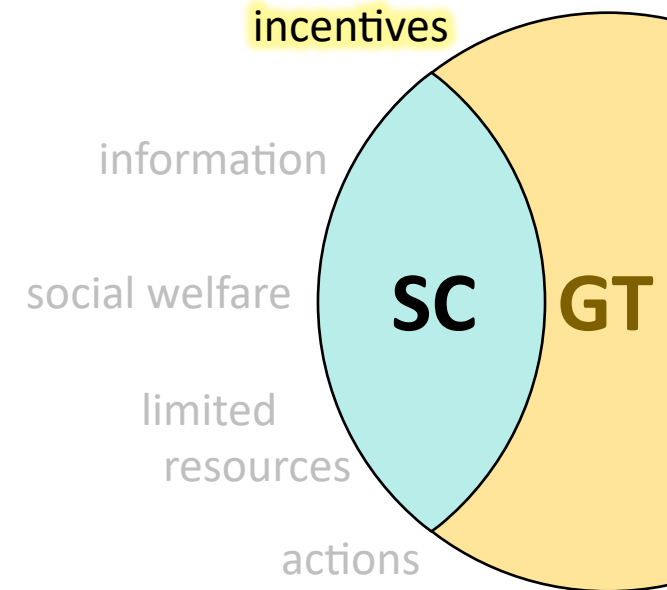
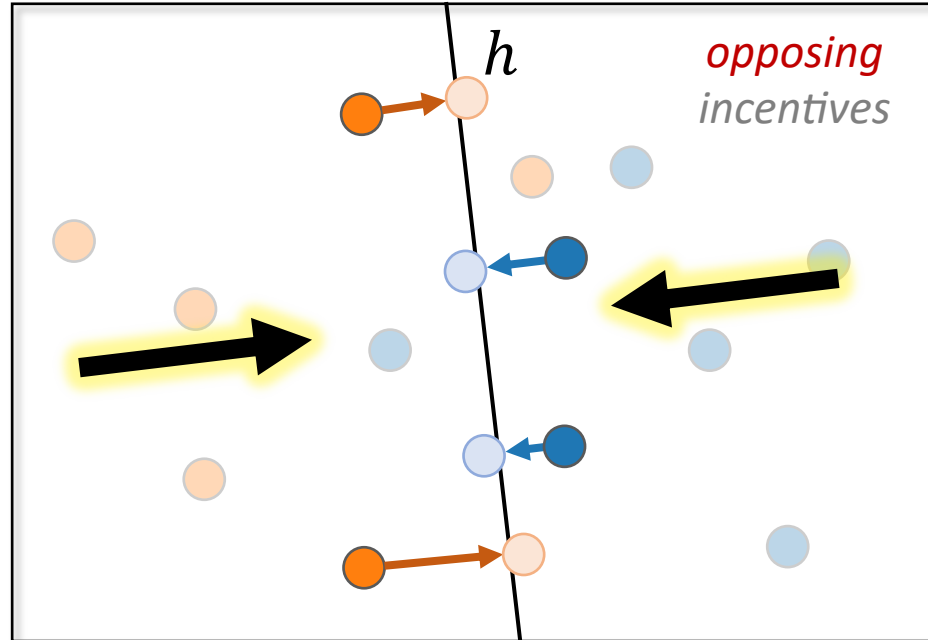
classification **against** humans (?)



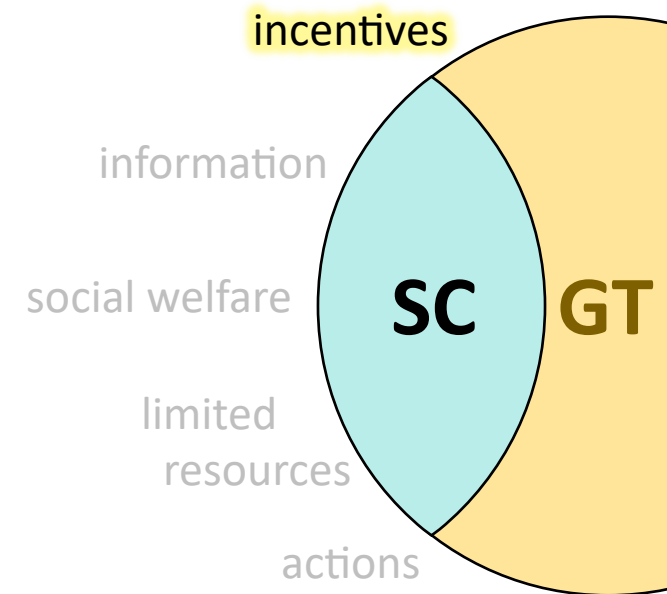
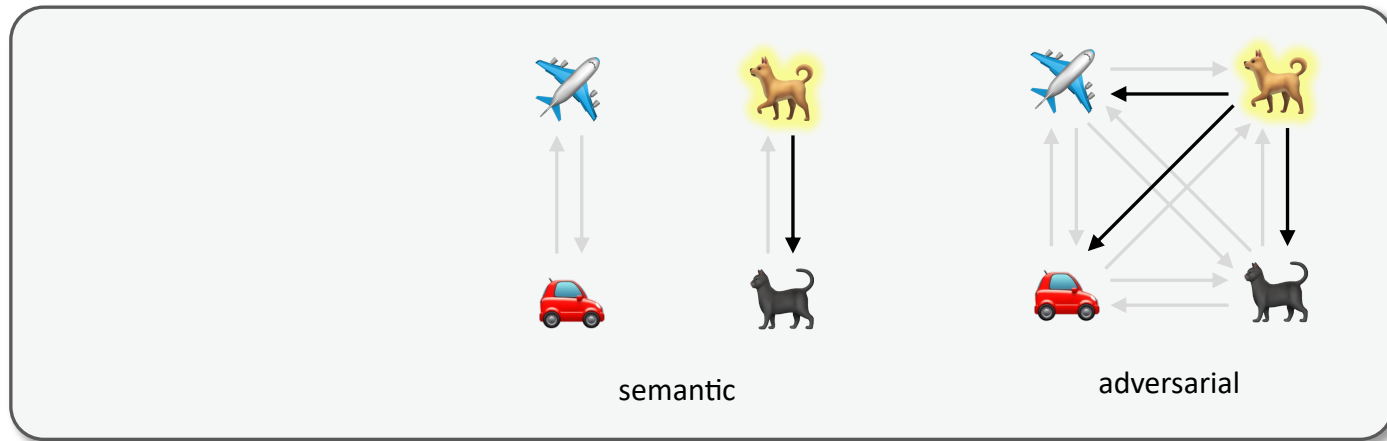
system wants: *correct* predictions

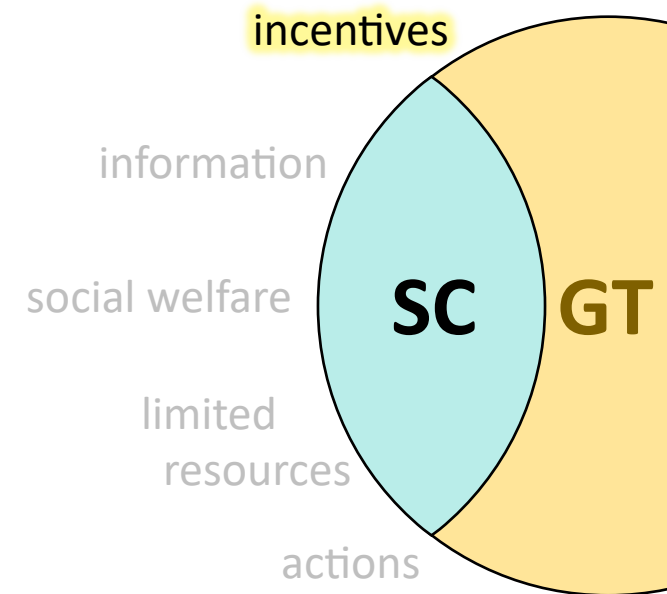
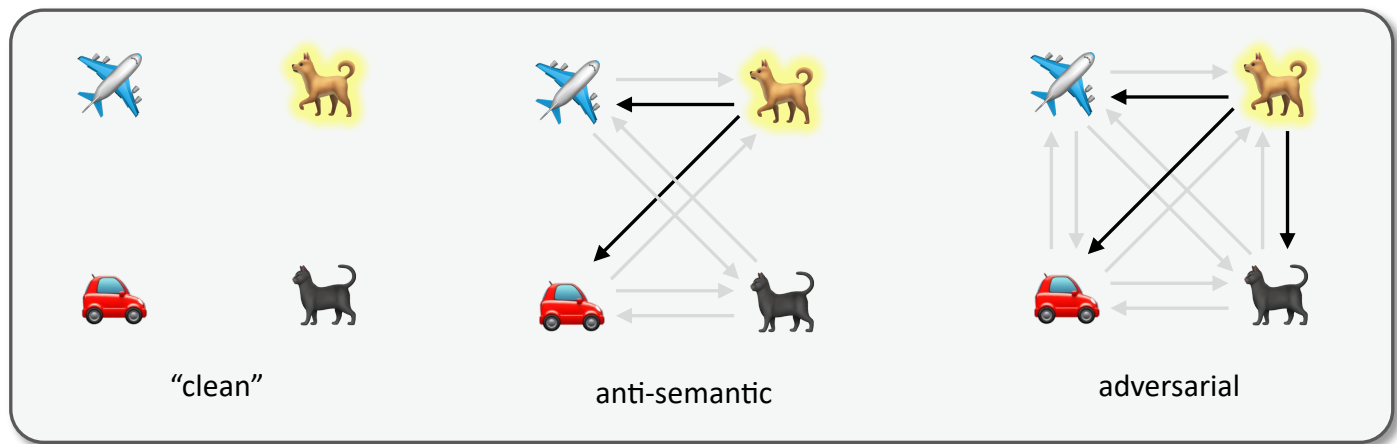
users want: *wrong* predictions

adversarial:



ask: can strategic modeling help make adversarial training *less conservative*?





“continuum”



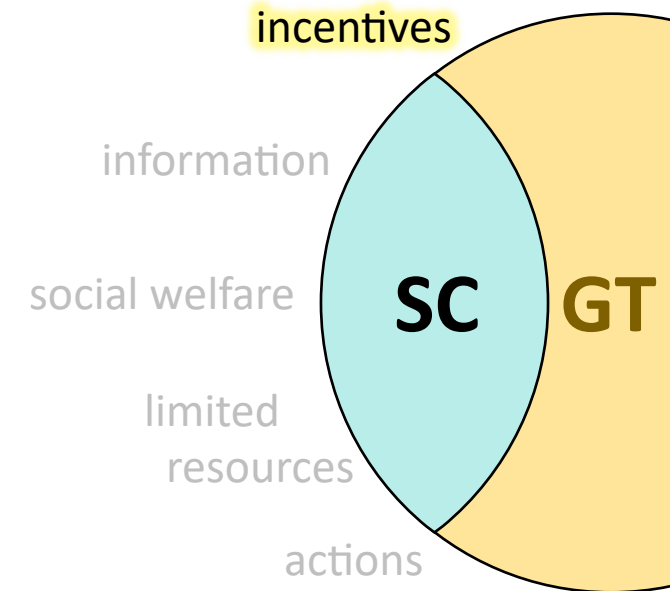
		CIFAR-10									GTSRB								
		VGG			ResNet18			ViT			VGG			ResNet18			ViT		
	test	clean	strat	adv	clean	strat	adv	clean	strat	adv	clean	strat	adv	clean	strat	adv	clean	strat	adv
a-sem.	train																		
	clean	89.2	0.1	0.6	93.1	0.0	0.0	77.5	0.1	0.0	95.1	5.8	15.4	96.6	1.8	1.6	90.7	2.8	1.8
	strategic	72.2	50.9	33.7	78.9	52.5	40.7	58.1	44.0	23.6	82.4	58.3	47.0	92.0	64.9	61.7	80.2	49.9	42.4
adversarial	67.7	46.4	39.3	79.8	49.6	45.5	53.1	39.8	33.3	80.5	47.2	43.8	90.9	55.4	53.9	79.4	49.6	47.1	
a-sem.	train																		
	clean	89.2	0.2	0.6	93.1	0.0	0.0	77.5	0.2	0.0	95.1	14.6	15.4	96.6	4.3	1.6	90.7	6.4	1.8
	strategic	80.3	69.8	25.0	85.1	73.4	33.4	65.3	58.7	13.1	88.2	71.2	36.7	92.0	79.1	49.0	83.5	74.0	34.0
adversarial	67.7	56.8	39.3	79.8	67.4	45.5	53.1	45.8	33.3	80.5	63.2	43.8	90.9	73.0	53.9	79.4	64.6	47.1	

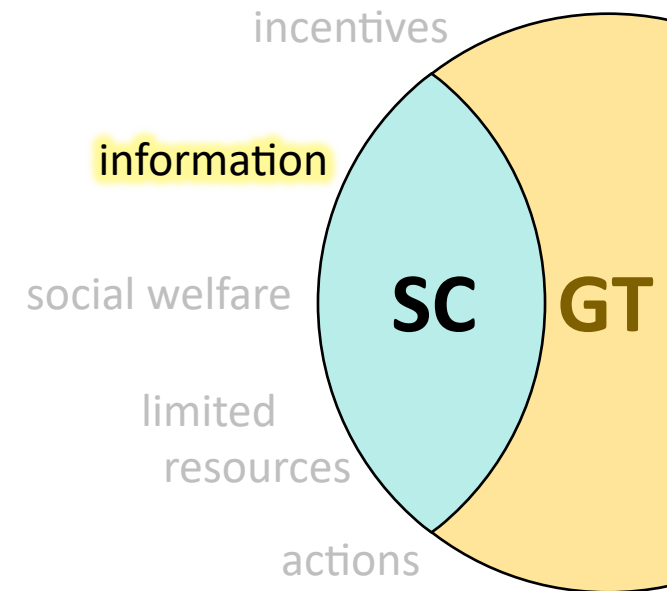
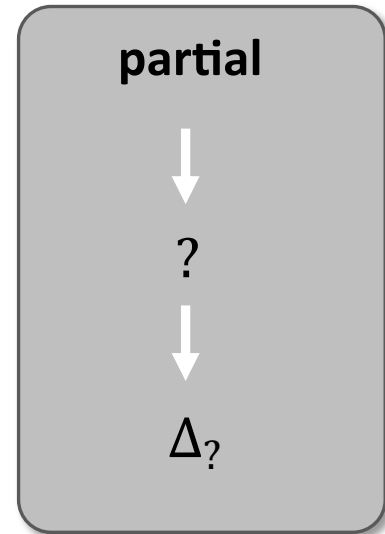
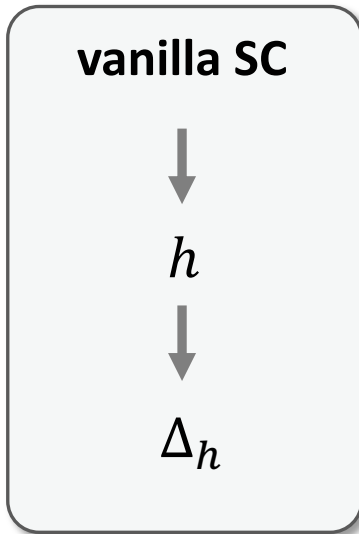
much potential for synergy!
will return to this

A note on strategic vs. adversarial learning:

- From SC perspective, **adversarial** is “**special case**”
- But only in a narrow sense – many distinctions in practice
- E.g., in **adversarial learning** (vs. **strategic learning**):

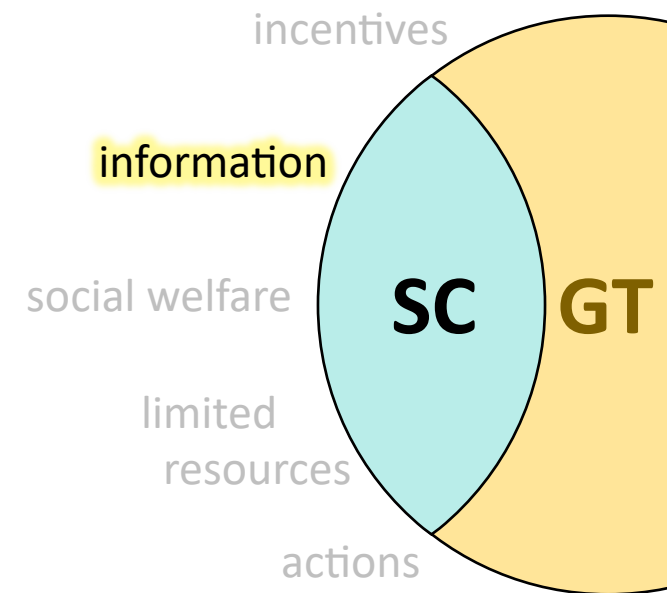
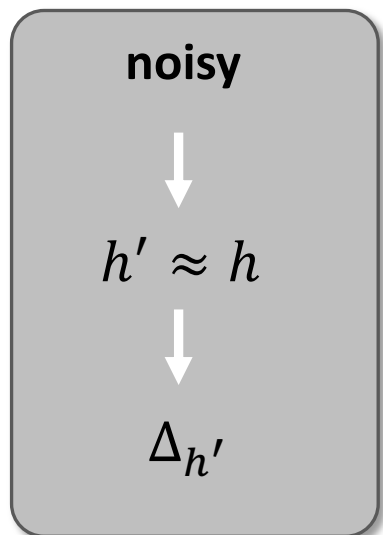
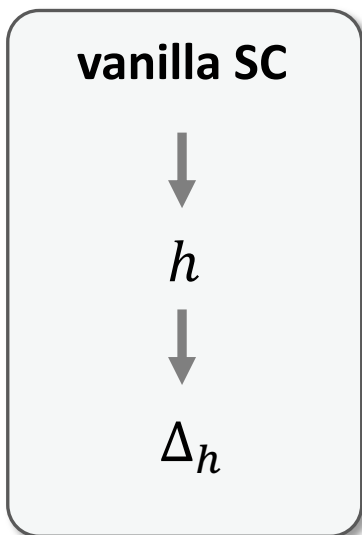
- attack **proxy loss** (e.g. log-loss) → vs. 0-1
- focus on **non-linear models** → vs. mostly linear
- focus on **complex modalities** (e.g. images) → vs. mostly tabular
- ⇒ best-responses are **approximate** → vs. exact best-responses
- vulnerabilities mostly in **latent space** → vs. no latent space
- maximize utility under **budget constraints** → vs. minimize cost
- ⇒ features **always modified** and **to the max** → vs. modify minimally – and only if needed
- ⇒ optimize **minimax objective** → vs. nested min-argmax



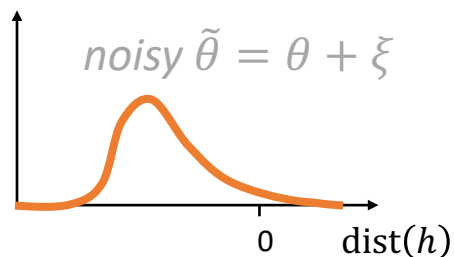
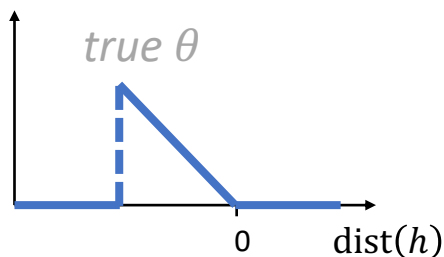


ask: what happens when users have partial knowledge of h ?

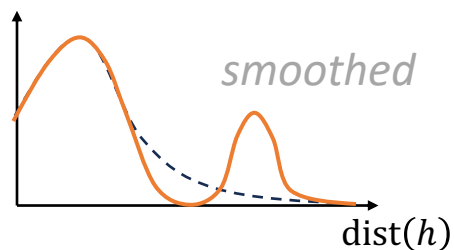
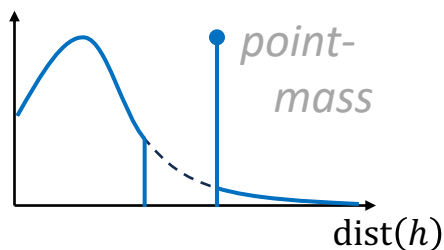
[JMH ICML21]



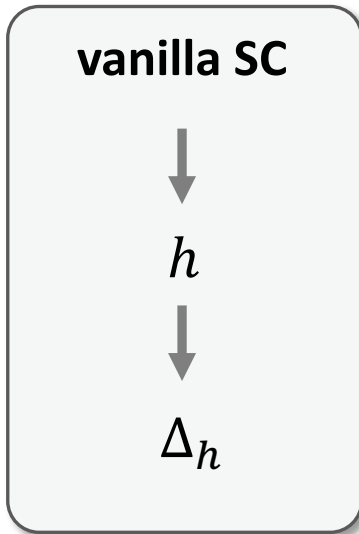
response curve:



induced distribution:



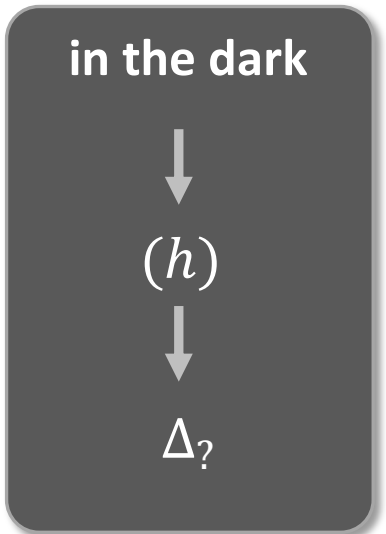
ask: what happens when users have partial knowledge of h ?



transparent



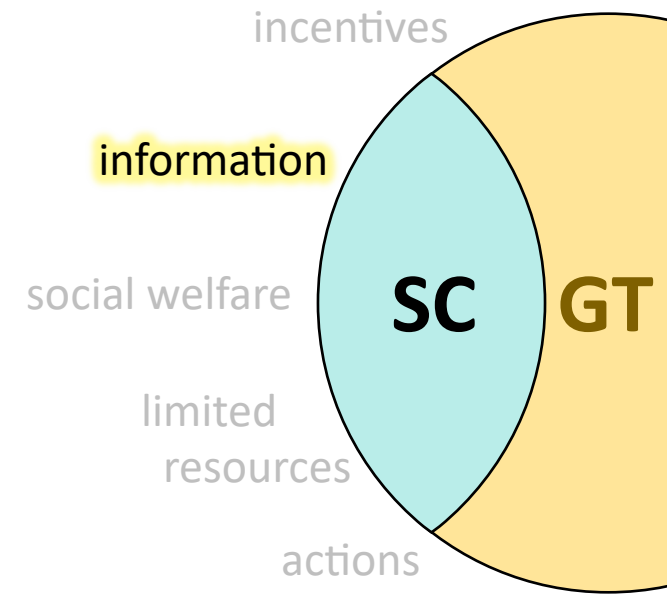
- informed users:**
- have **more** power
 - **easy** to anticipate



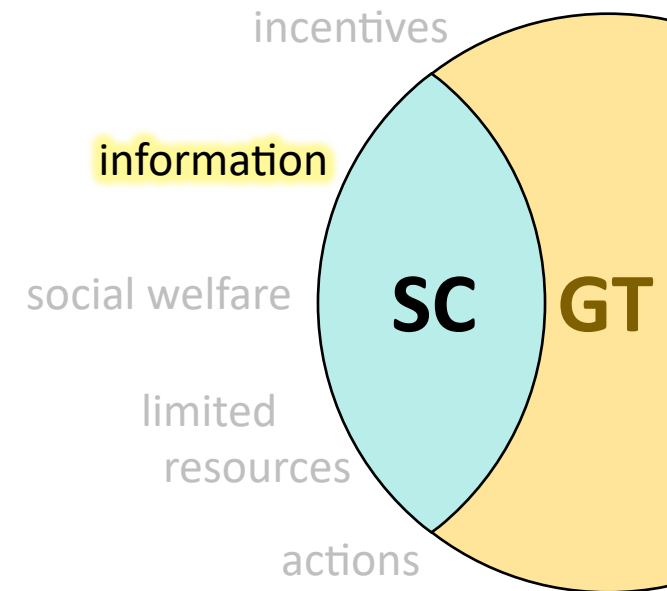
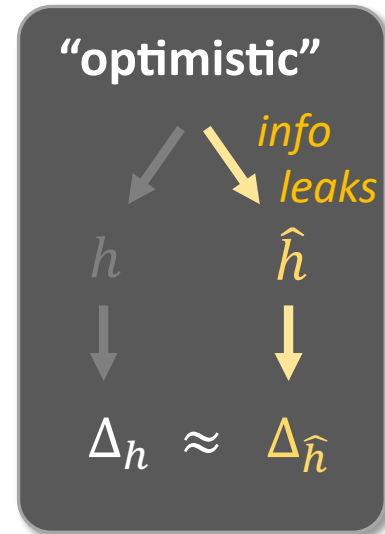
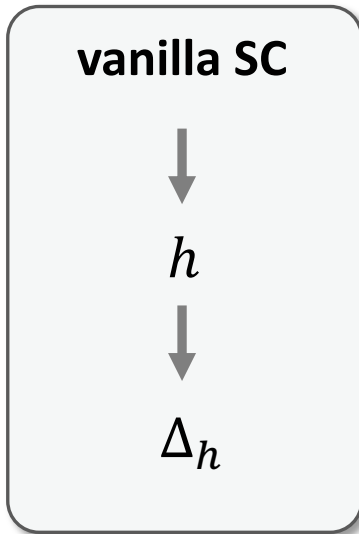
opaque



- uninformed users:**
- have **less** power
 - **harder** to anticipate

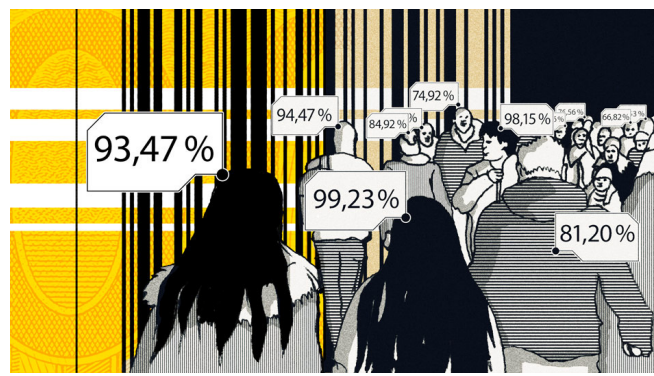
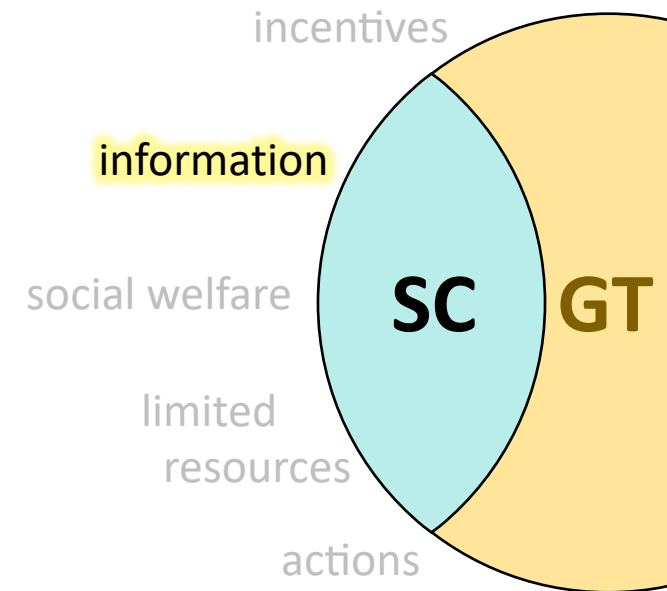
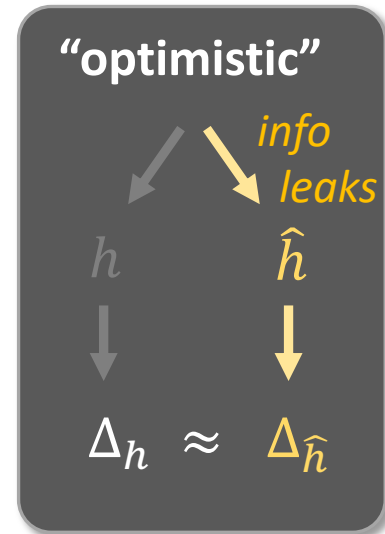
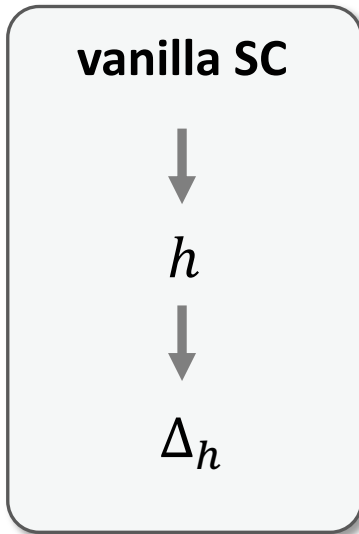


ask: what happens when users have partial knowledge of h ?

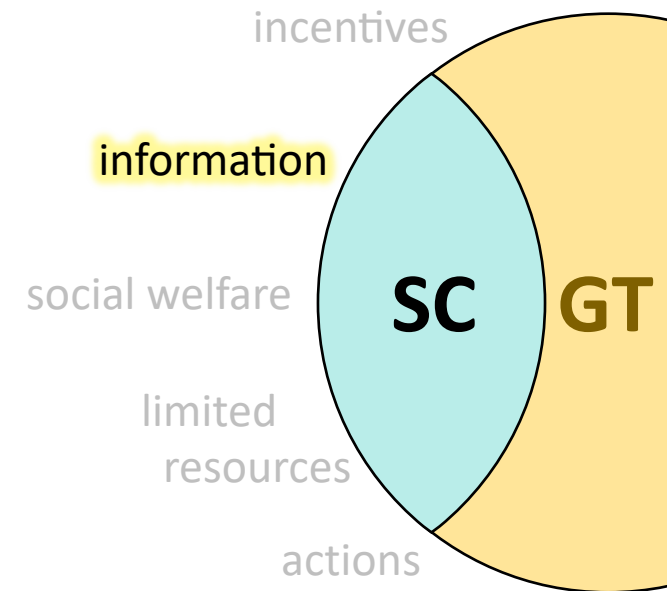
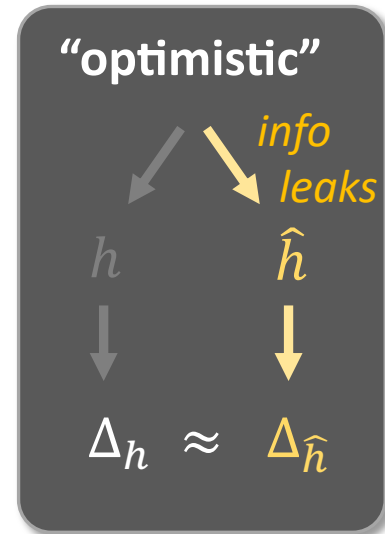
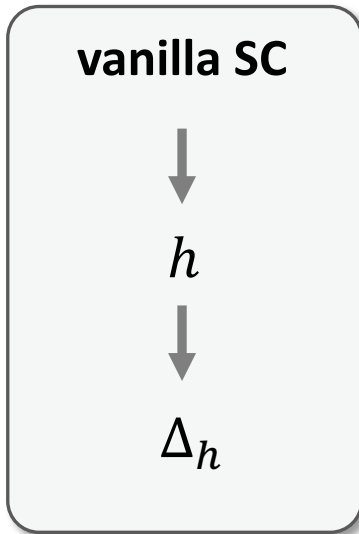


SCHUFA

ask: what happens when users have partial knowledge of h ?

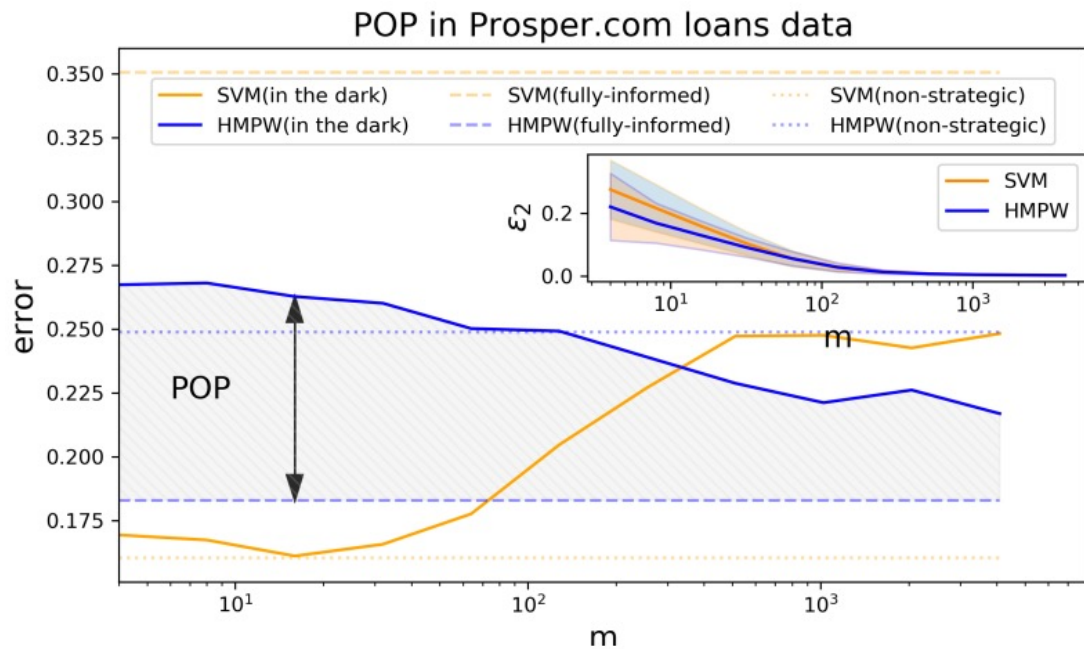


ask: what happens when users have partial knowledge of h ?



- **Price of OPacity: (POP)**
 $err(h, \hat{h}) - err(h, h)$
- **Main result:** can be **arbitrarily bad**
 \Rightarrow **transparency** is often in **best interest of system!**

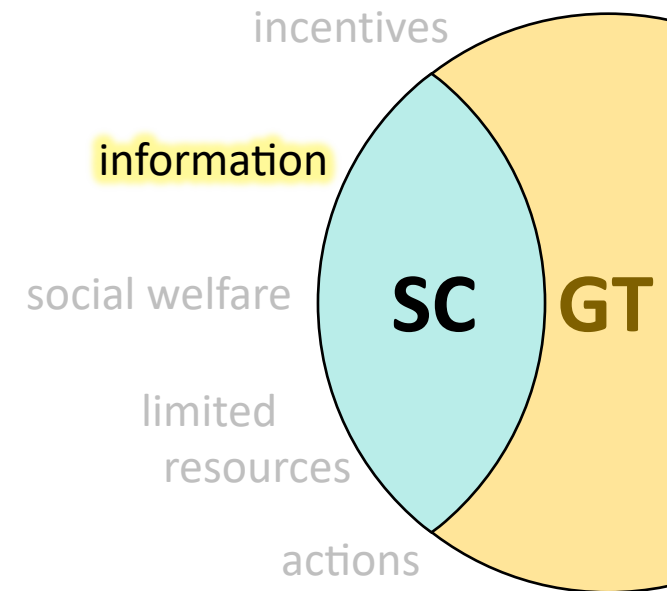
ask: what happens when users have partial knowledge of h ?



- **Price of OPacity: (POP)**

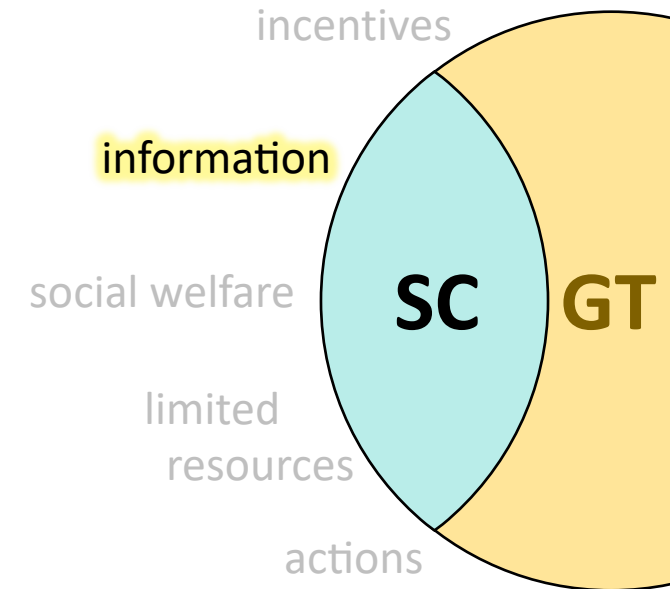
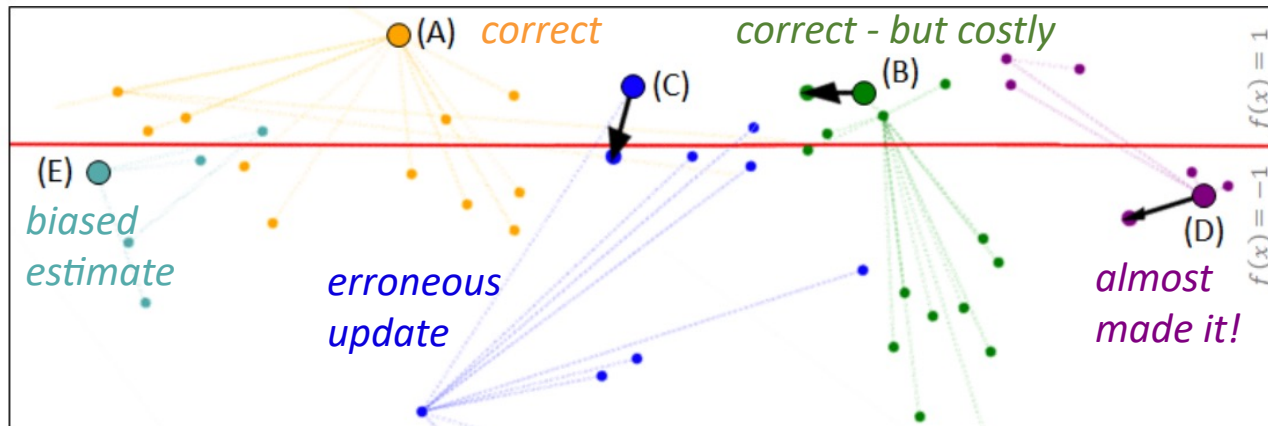
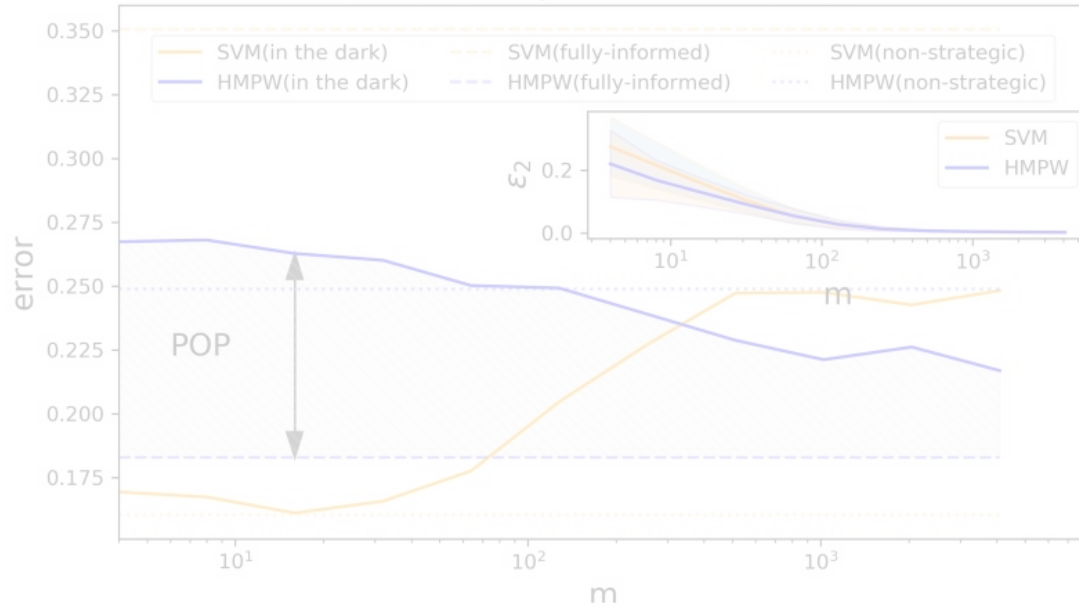
$$err(h, \hat{h}) - err(h, h)$$

- **Main result:** can be **arbitrarily bad**
 \Rightarrow **transparency** is often in **best interest of system!**



ask: what happens when users have partial knowledge of h ?

POP in Prosper.com loans data



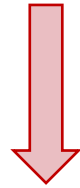
ask: what happens when users have partial knowledge of h ?

known user response:

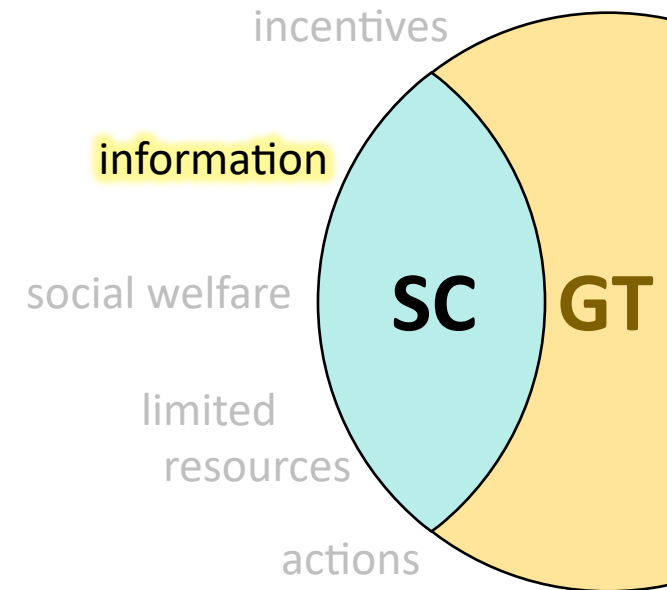
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

unknown user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_?(x_i)))$$



uncertainty



ask: how can learning contend with uncertain user behavior?

known user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

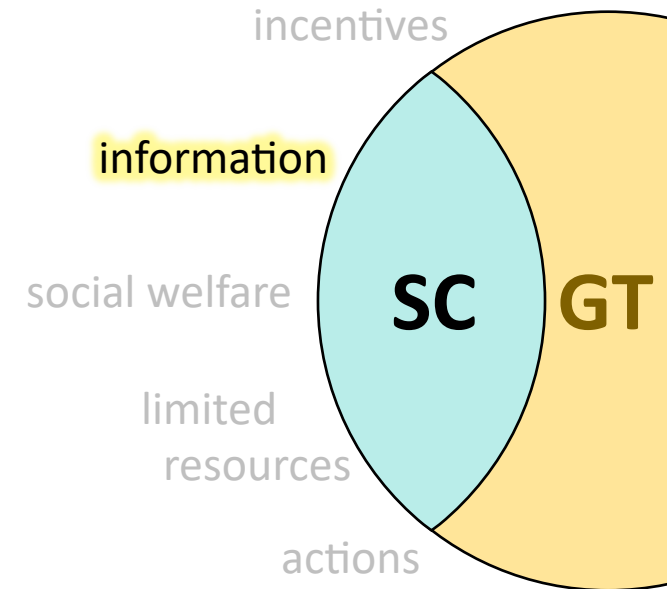
unknown user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_?(x_i)))$$



uncertainty

1) infer Δ over time (more on this later)



ask: how can learning contend with uncertain user behavior?

known user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$



uncertainty

unknown user response:

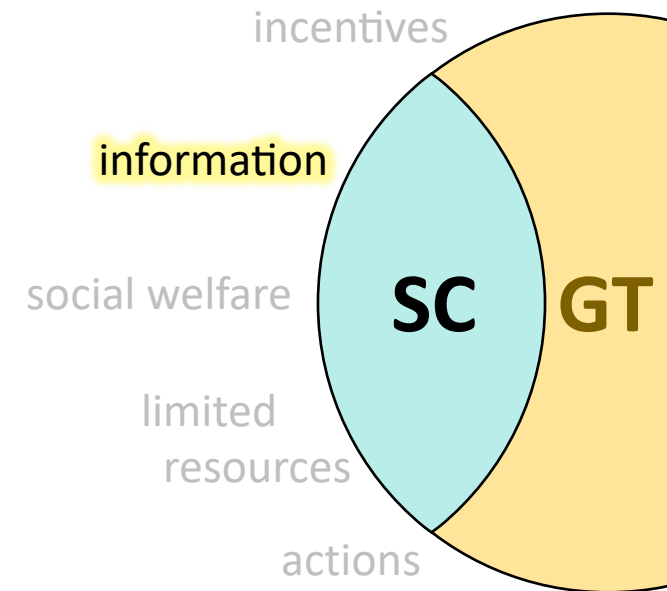
$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_?(x_i)))$$

1) infer Δ over time (more on this later)

2) **robust learning:**

$$\operatorname{argmin}_h \max_{\Delta \in \mathcal{U}} \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

↑
uncertainty set



ask: how can learning contend with uncertain user behavior?

known user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

unknown user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_?(x_i)))$$

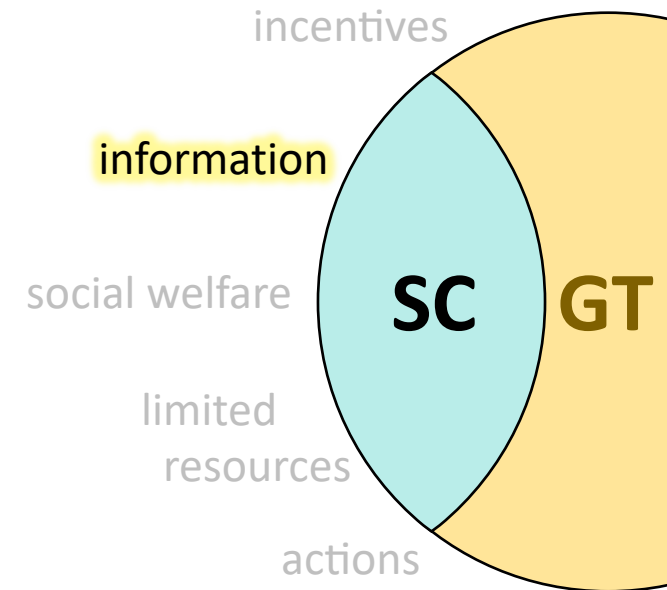
↓ uncertainty

1) infer Δ over time (more on this later)

2) **robust learning** – unknown costs:

$$\operatorname{argmin}_h \max_{c \in \mathcal{C}} \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h^c(x_i)))$$

↑ uncertainty set



ask: how can learning contend with uncertain user behavior?

known user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

unknown user response:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_?(x_i)))$$

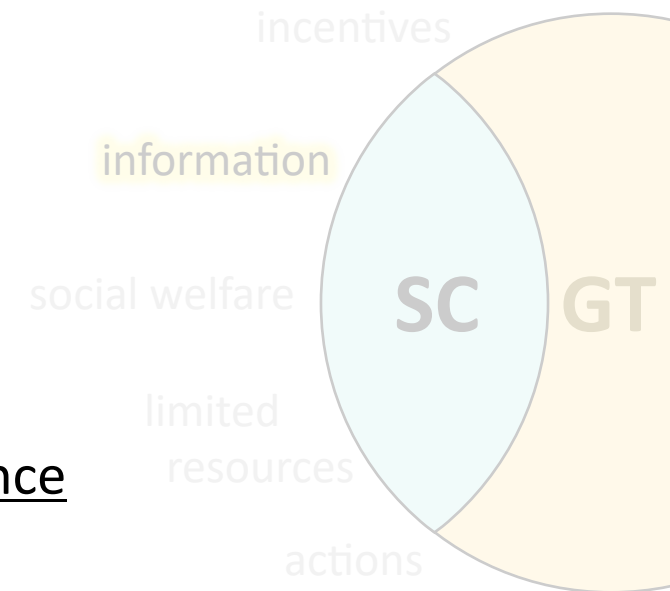
1) infer Δ over time (more on this later)

2) **robust learning** – unknown costs:

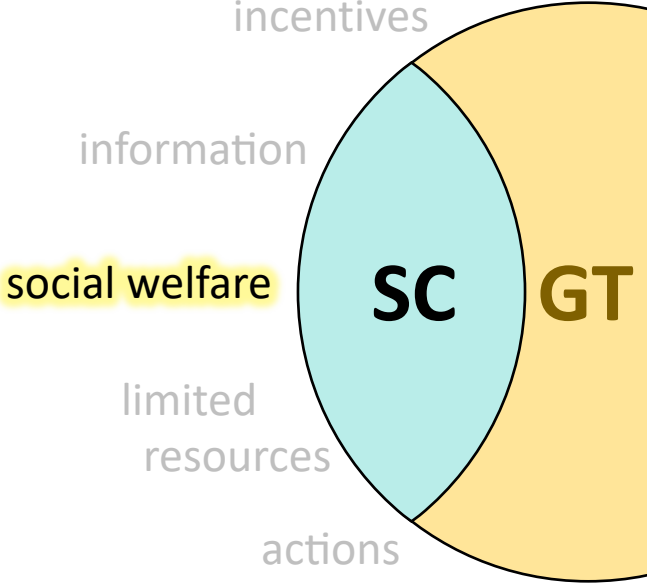
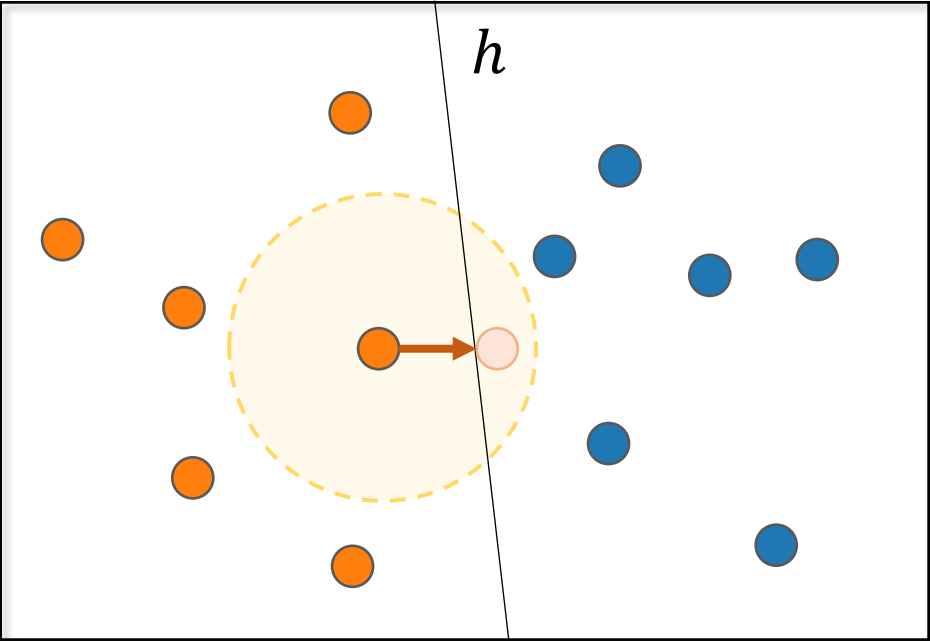
$$\operatorname{argmin}_h \max_{c \in \mathcal{C}} \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h^c(x_i)))$$

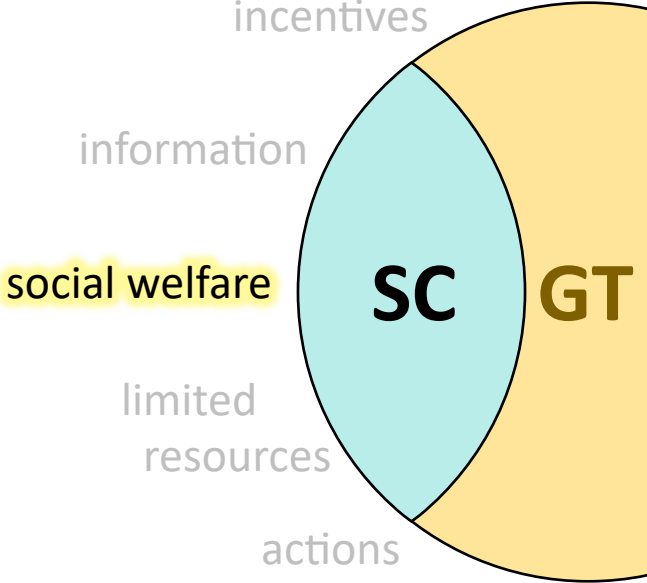
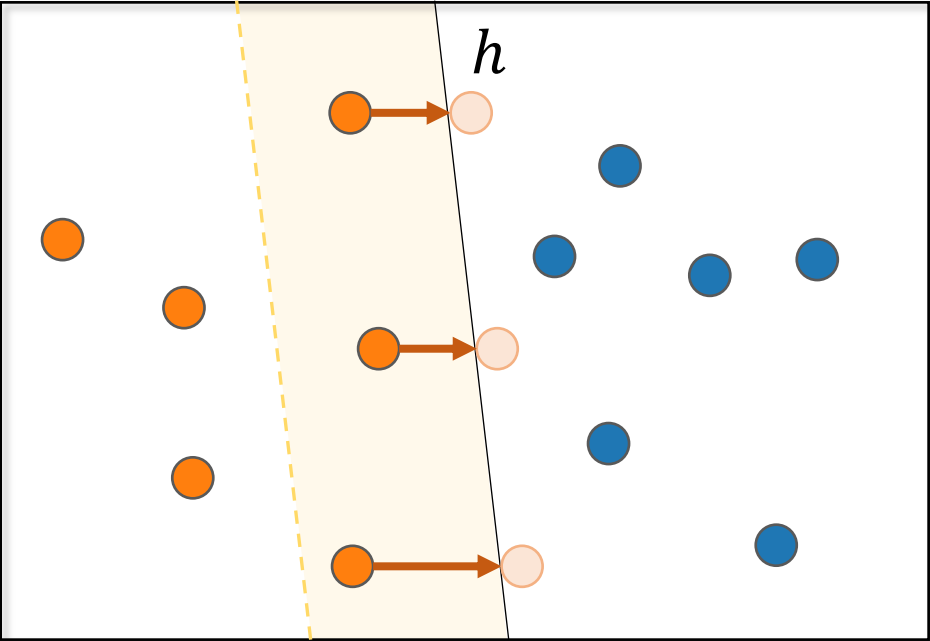
↑
uncertainty set

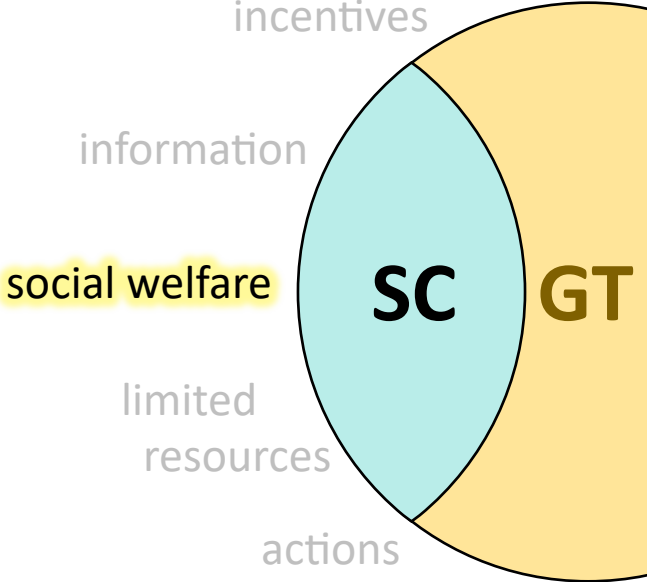
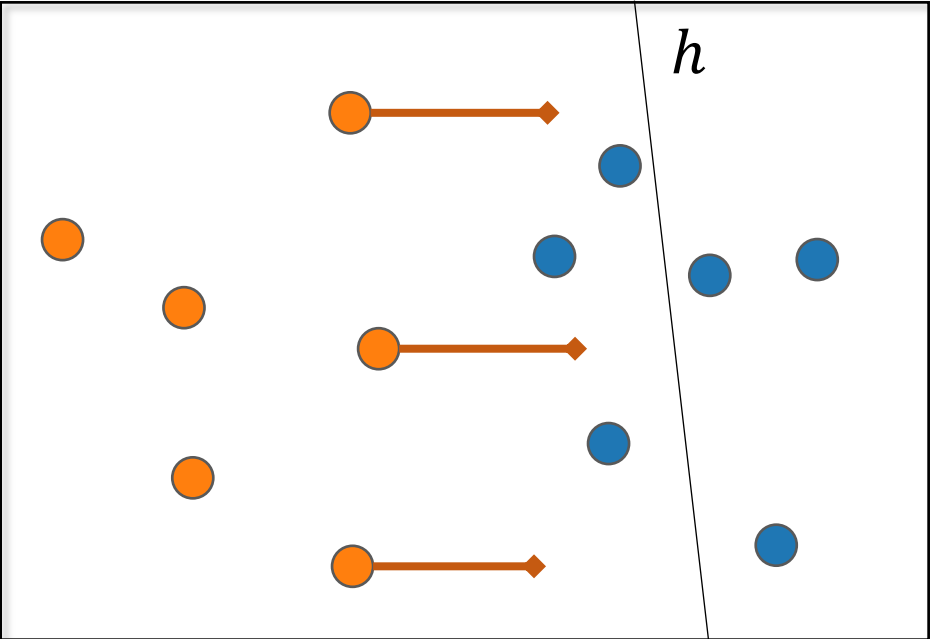
(public) policy problems:



- **"One shot"** – can deploy only once
- **Goal:** learn to be doubly-robust:
 - vs. strategic behavior
 - vs. worst-case cost $c \in \mathcal{C}$
- **Hardness:** not knowing c can be catastrophic
- **Convexification:** updated ad-hoc s-hinge
- **Algorithm:** effective, converge to opt. min-max

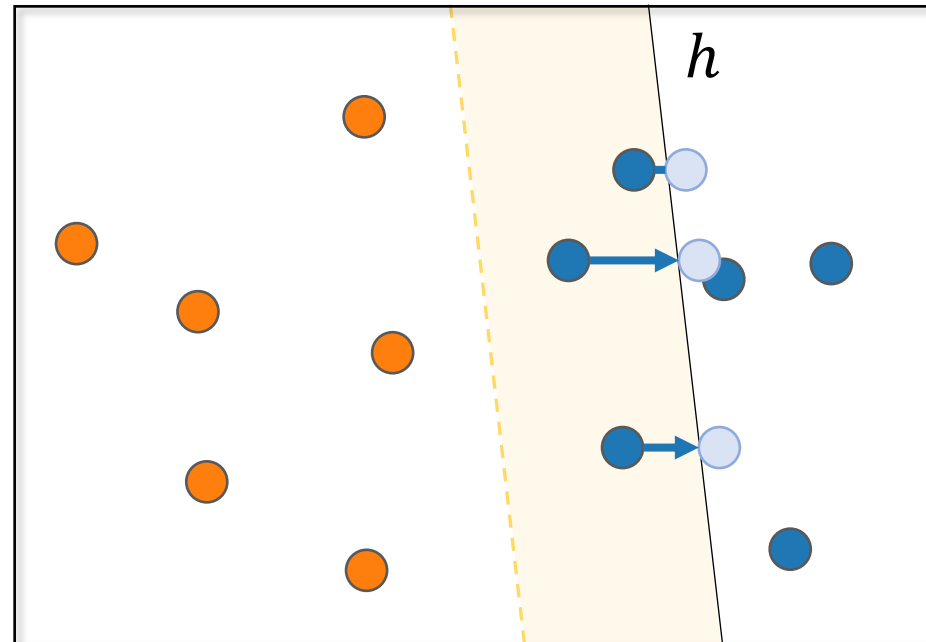






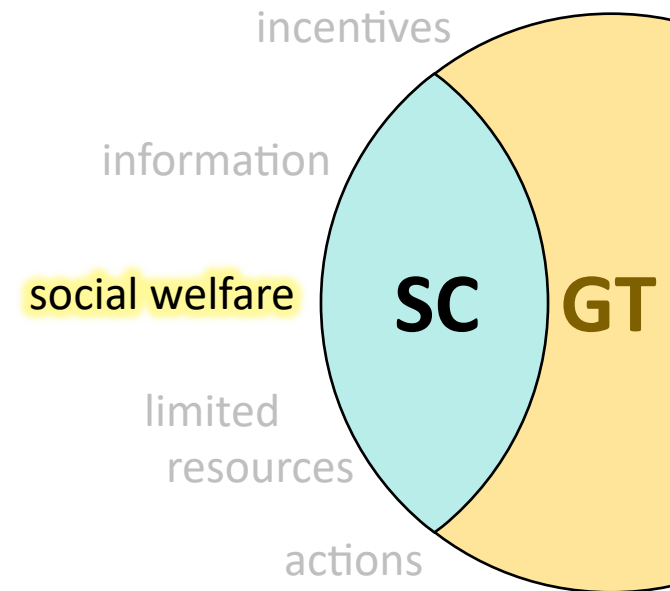
- Robustness via penalizing deserving sub-population
- Main result is **negative**:
increased accuracy \Rightarrow
increased social burden
- However, results apply to certain **monotone** setting
- In more general settings, there is **reason for optimism!**

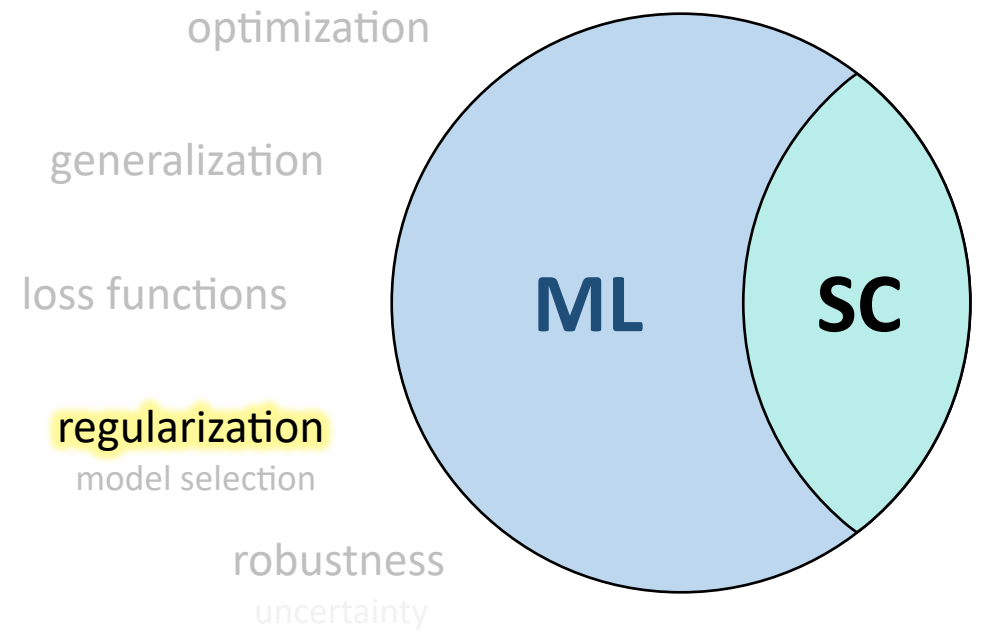
ask: when and how can we reduce social harm?



= “social burden” [MMDH’19]

$$\text{burden}(h) = \mathbb{E} \left[\min_{x': h(x')=1} c(x, x') \mid y = 1 \right]$$

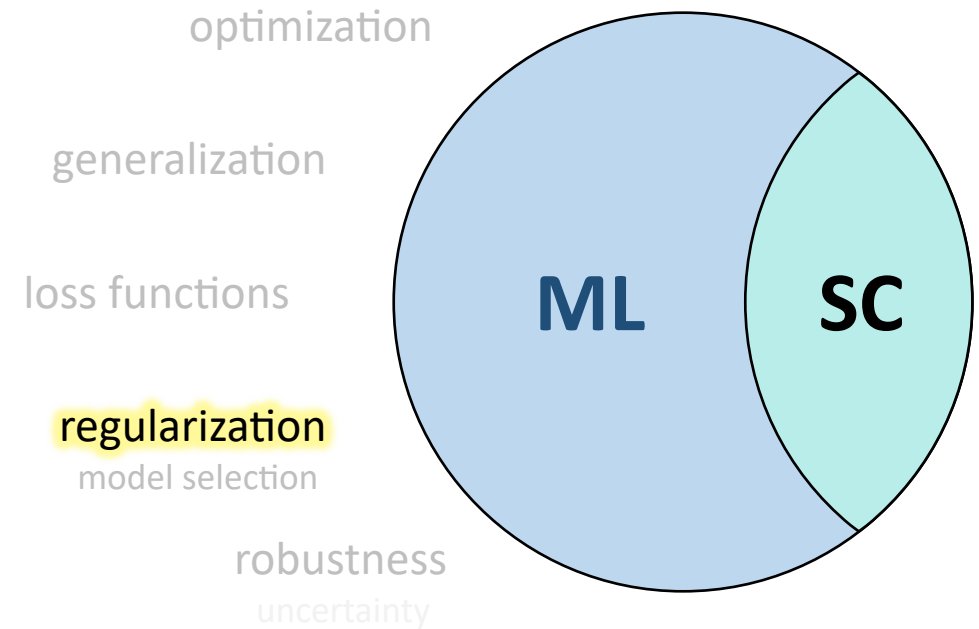
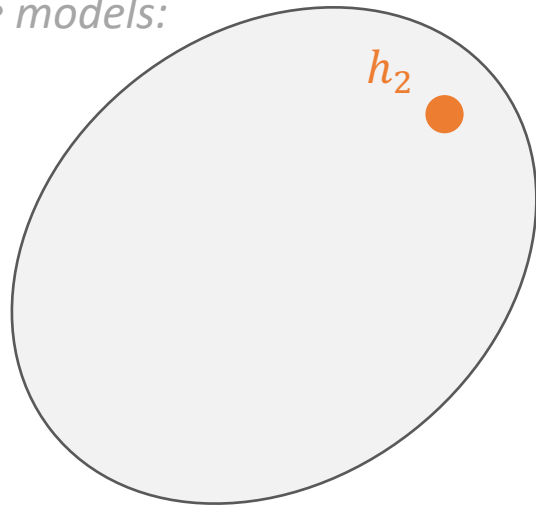




- **Conjecture:** many good models, **vary in burden**
- Learning objective underspecified – **can exploit!**
- Regularize for **generalization**:

$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda \|h\|_2$$

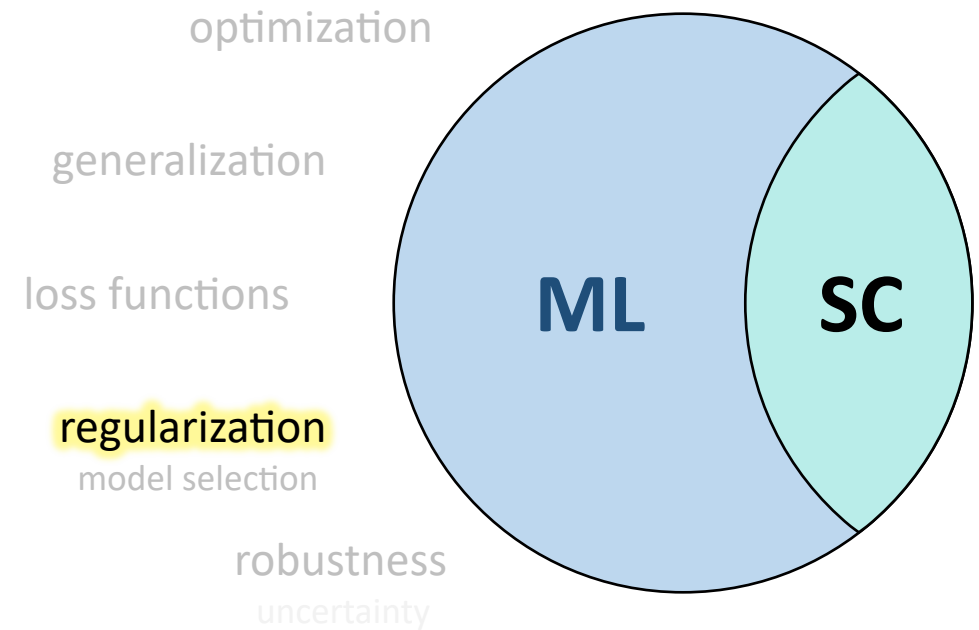
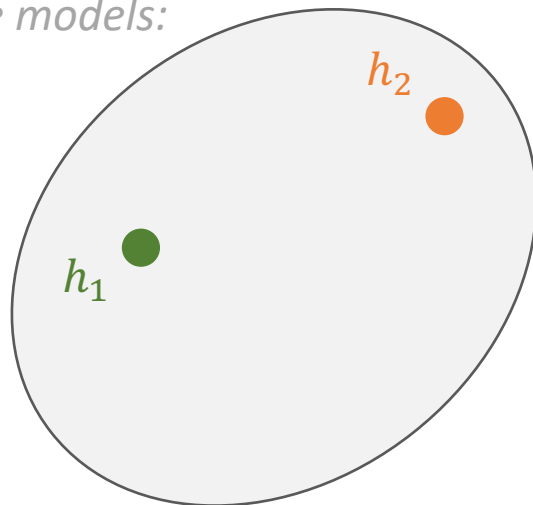
set of accurate models:



- **Conjecture:** many good models, vary in induced burden
- Learning objective underspecified – **can exploit!**
- Regularize for **sparsity**:

$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda \|h\|_1$$

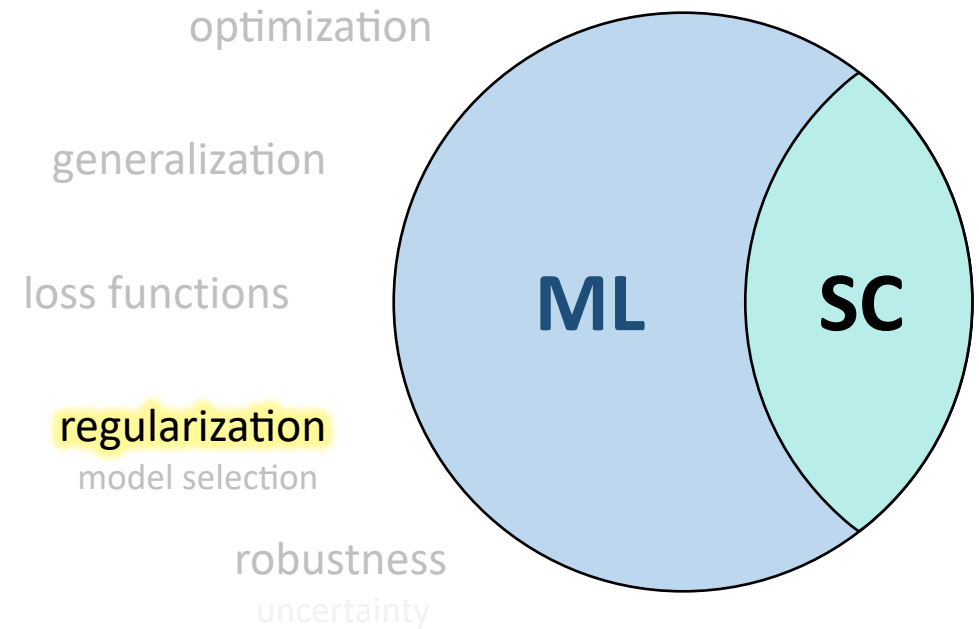
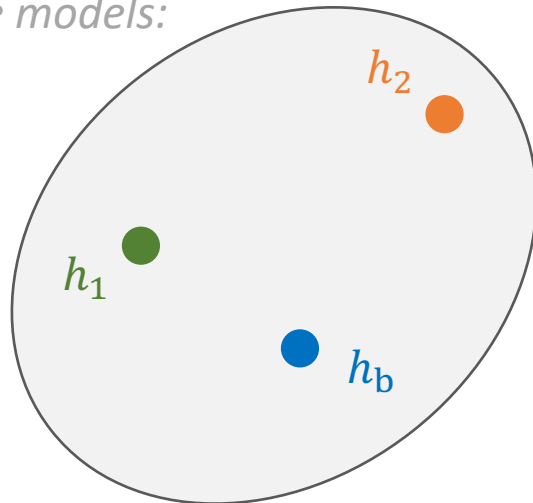
set of accurate models:



- **Conjecture:** many good models, vary in induced burden
- Learning objective underspecified – **can exploit!**
- Regularize for... **social good?**

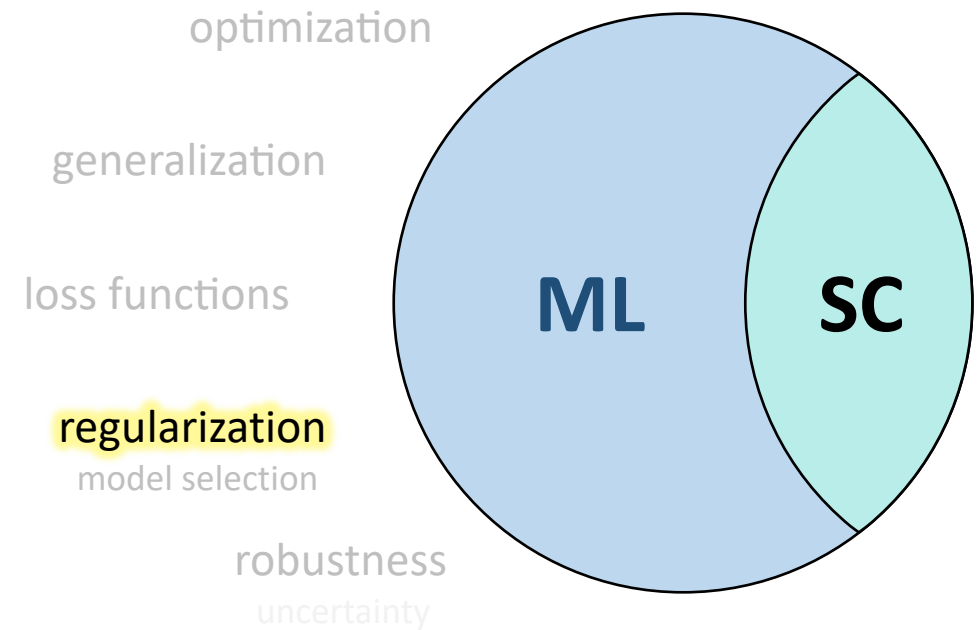
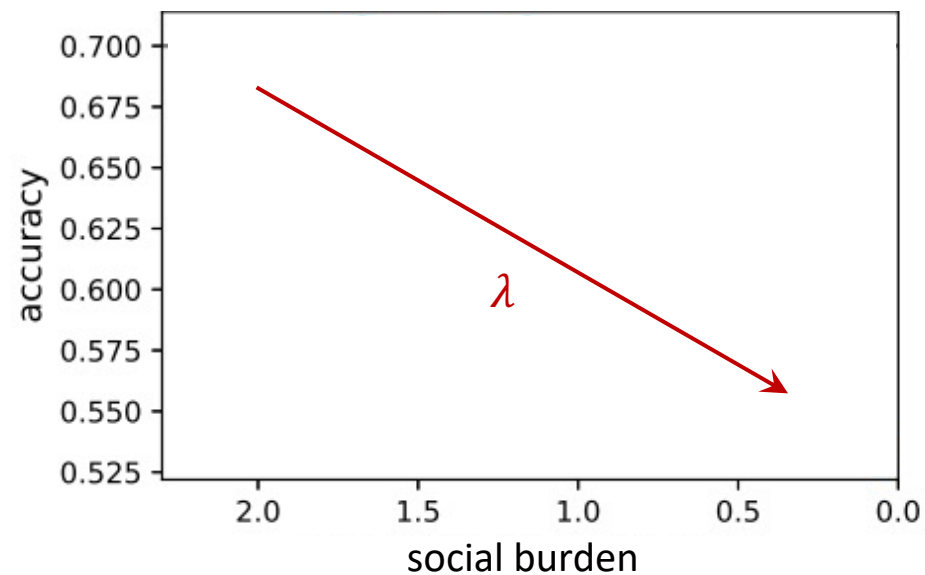
$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda R_{\text{burden}}(h)$$

set of accurate models:



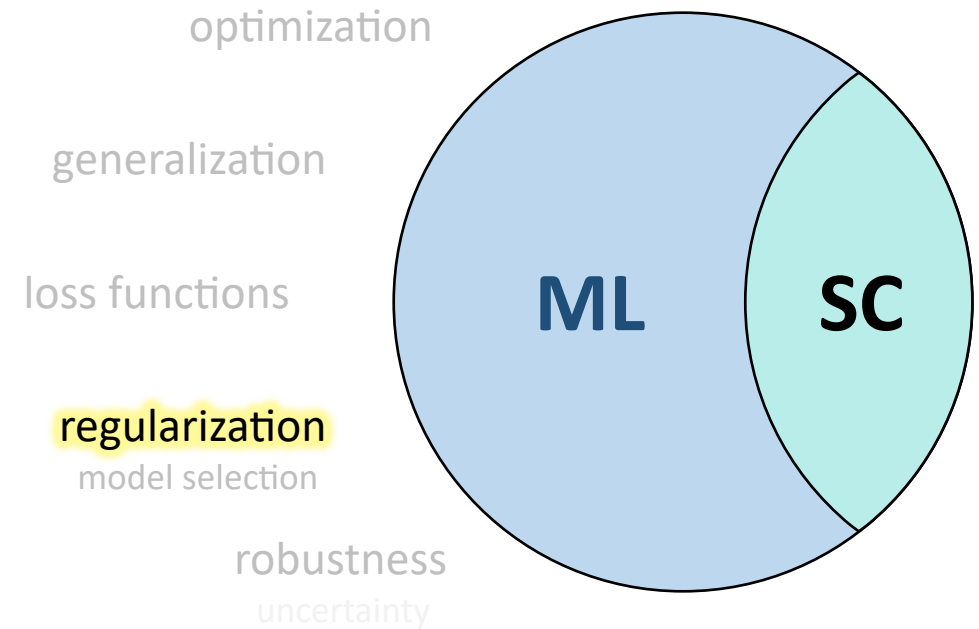
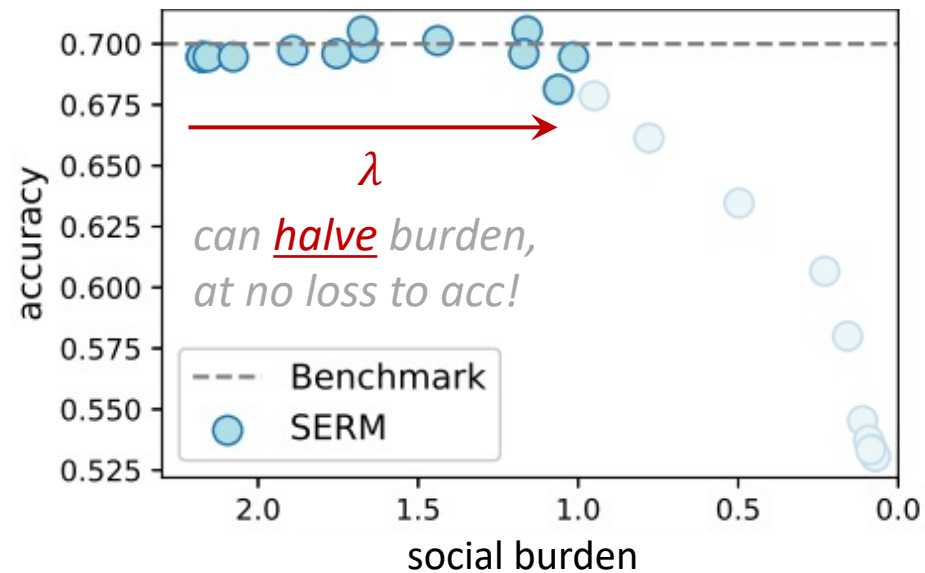
- **Conjecture:** many good models, vary in induced burden
- Learning objective underspecified – **can exploit!**
- Regularize for... **social good?**

$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda R_{\text{burden}}(h)$$



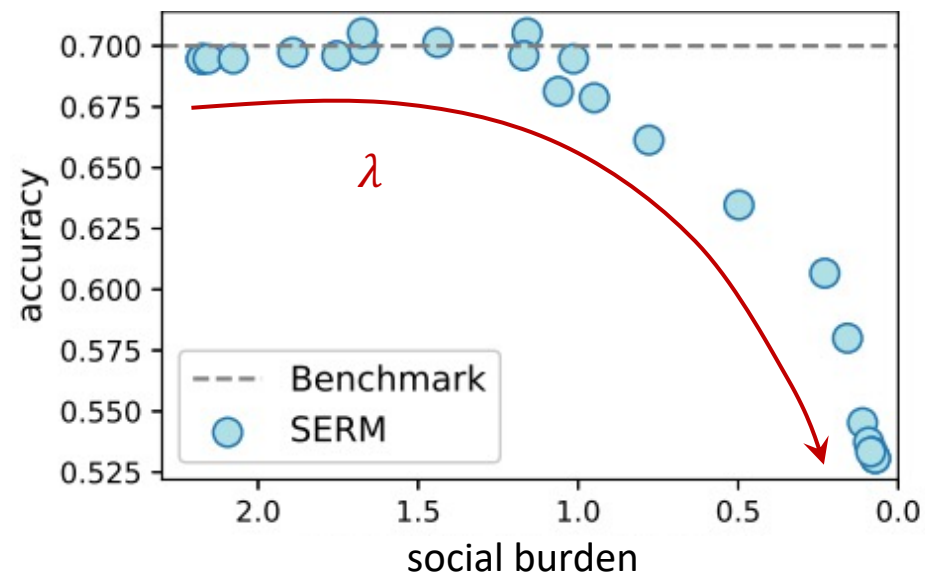
- **Conjecture:** many good models, vary in induced burden
- Learning objective underspecified – **can exploit!**
- Regularize for... **social good!**

$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda R_{\text{burden}}(h)$$



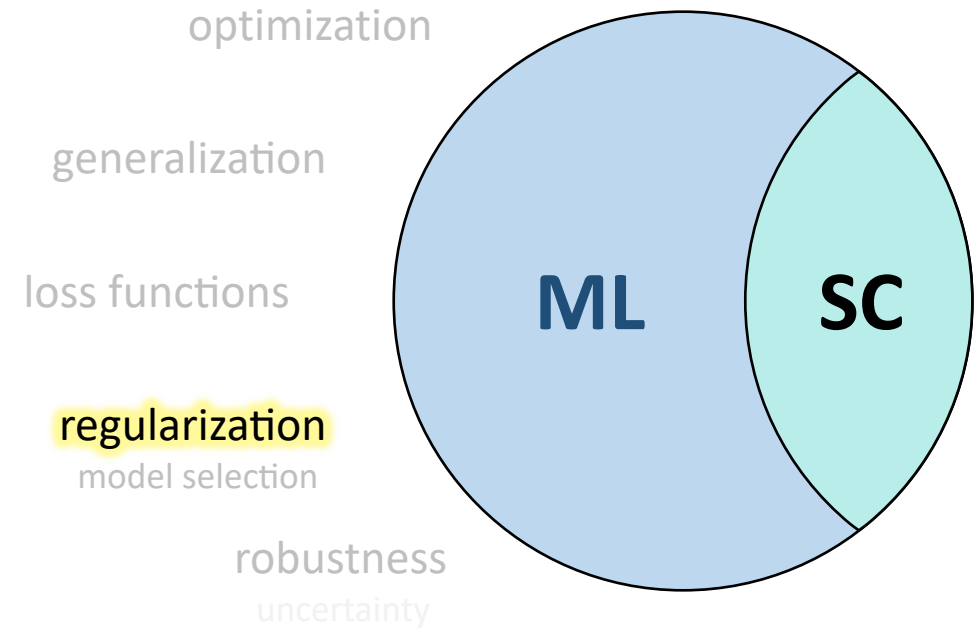
- **Conjecture:** many good models, vary in induced burden
- Learning objective underspecified – **can exploit!**
- Regularize for... **social good!**

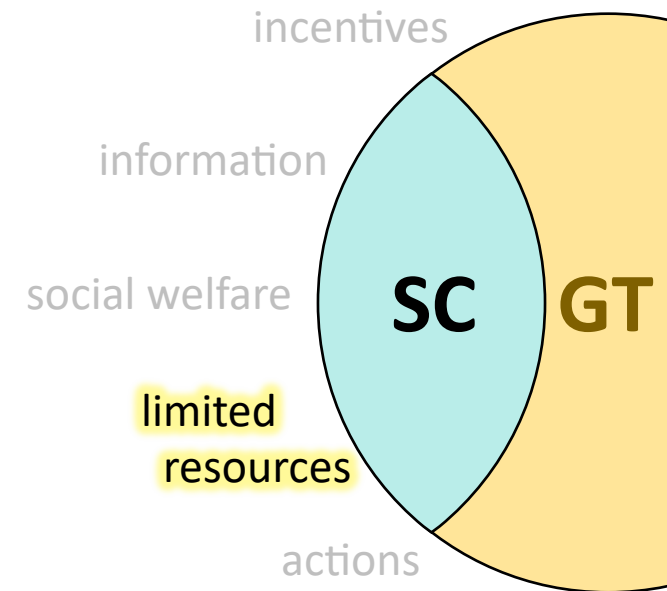
$$\operatorname{argmin}_{h \in H} \sum_{i=1}^m \ell(y, h(\Delta_h(x))) + \lambda R_{\text{burden}}(h)$$

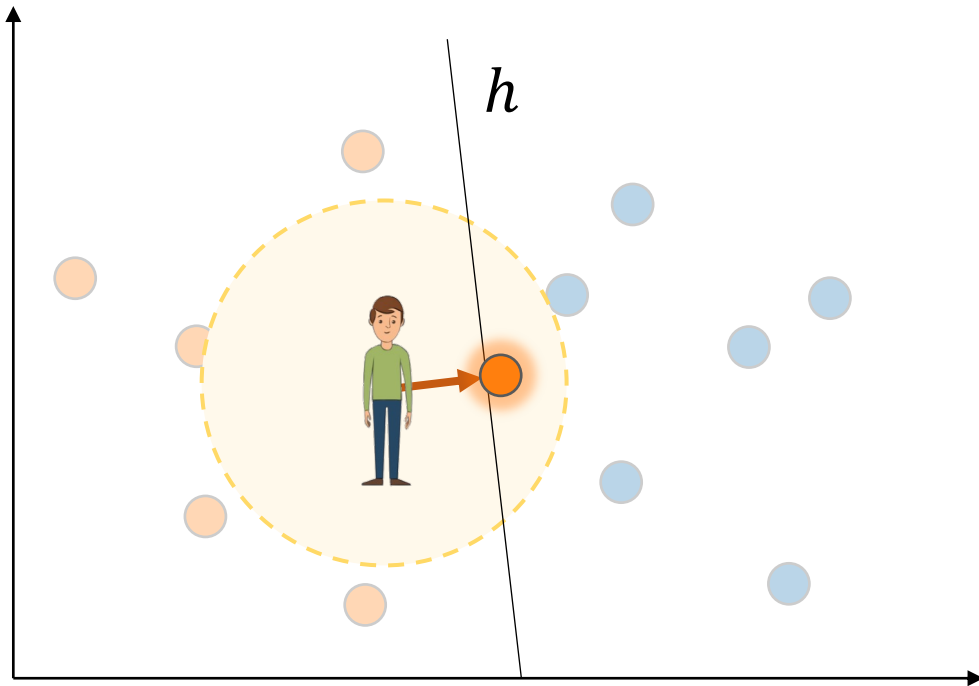


- Applies to other social good metrics (utility, recourse, ...)
- Similarly underspecified – similar pareto fronts!

additional social good metrics:

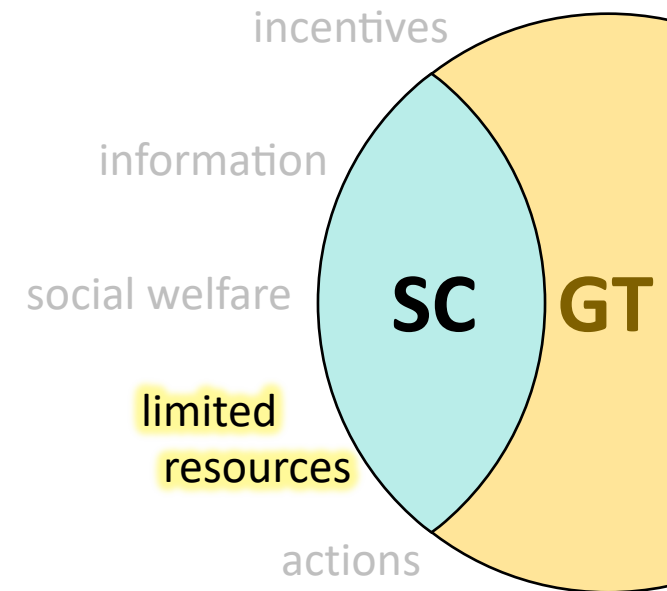






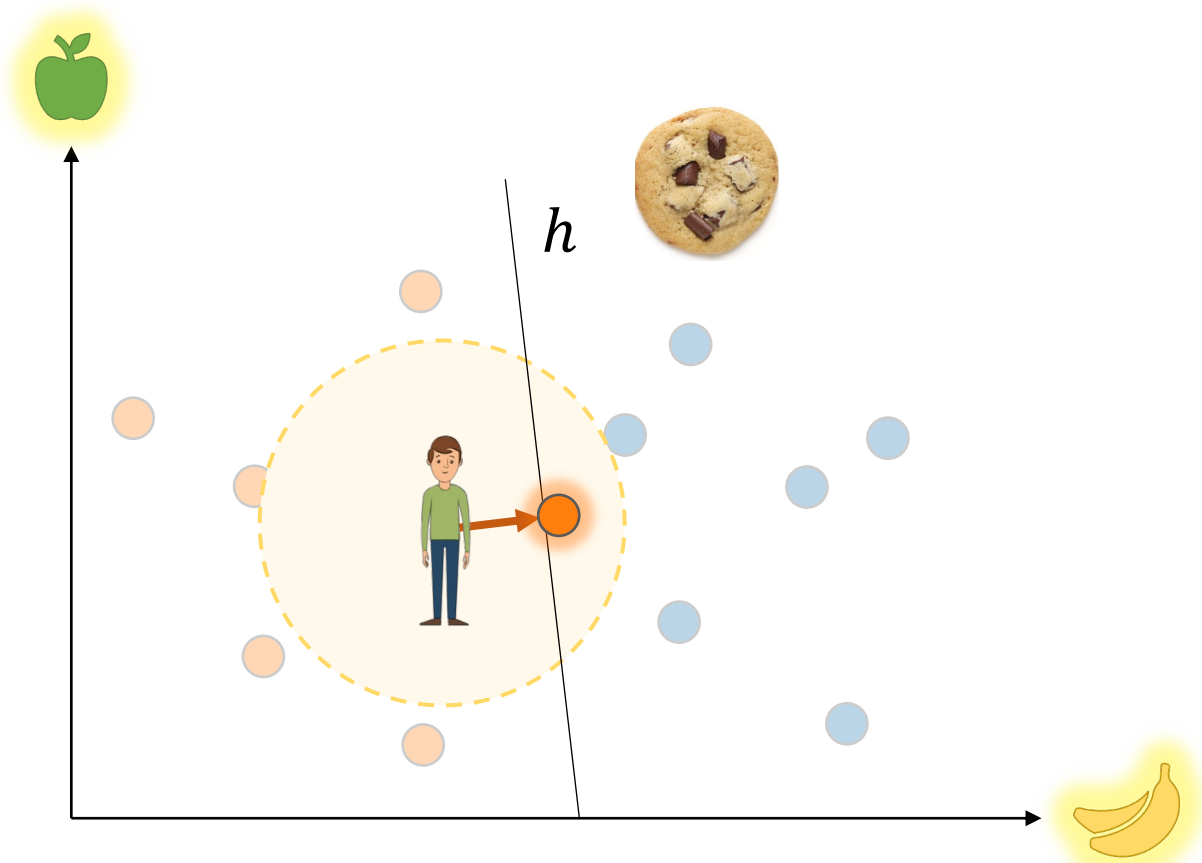
$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

utility *cost*



ask: where do costs come from?

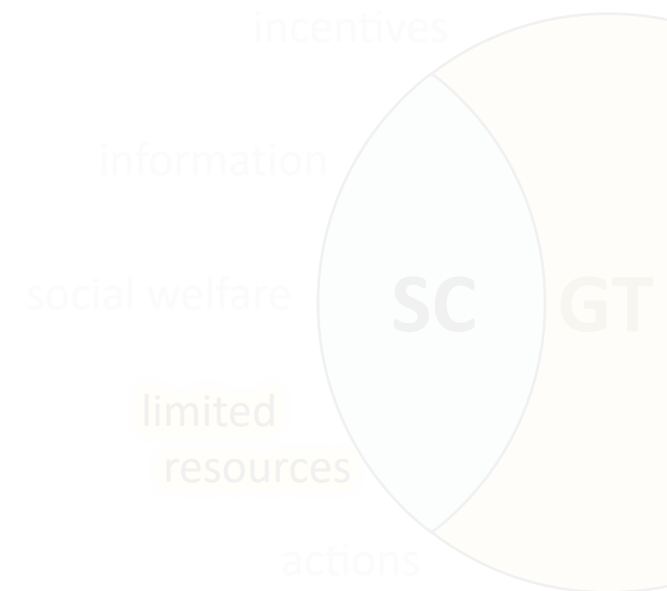
(ask first: what are features?)

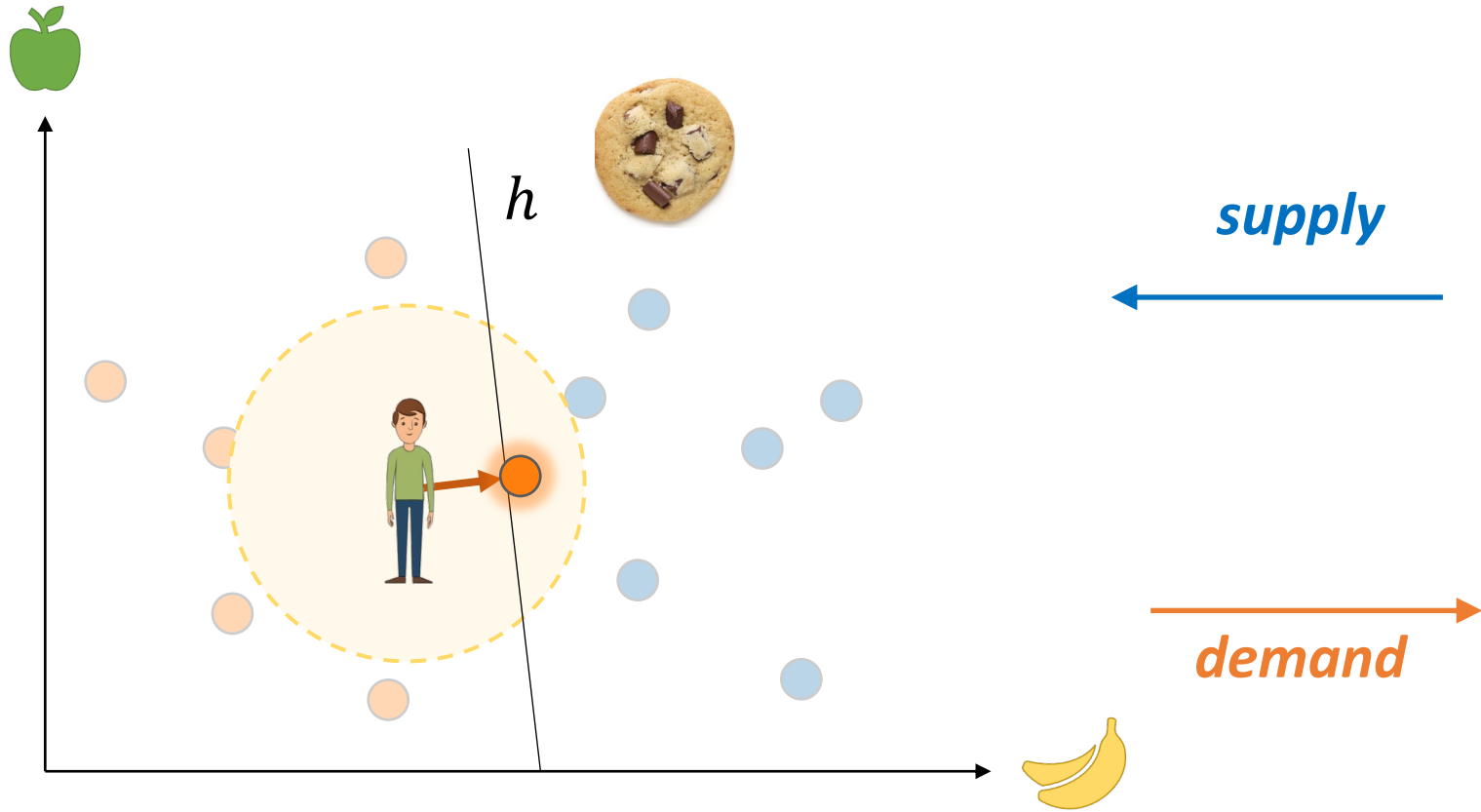


$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

utility
cost of apples and bananas

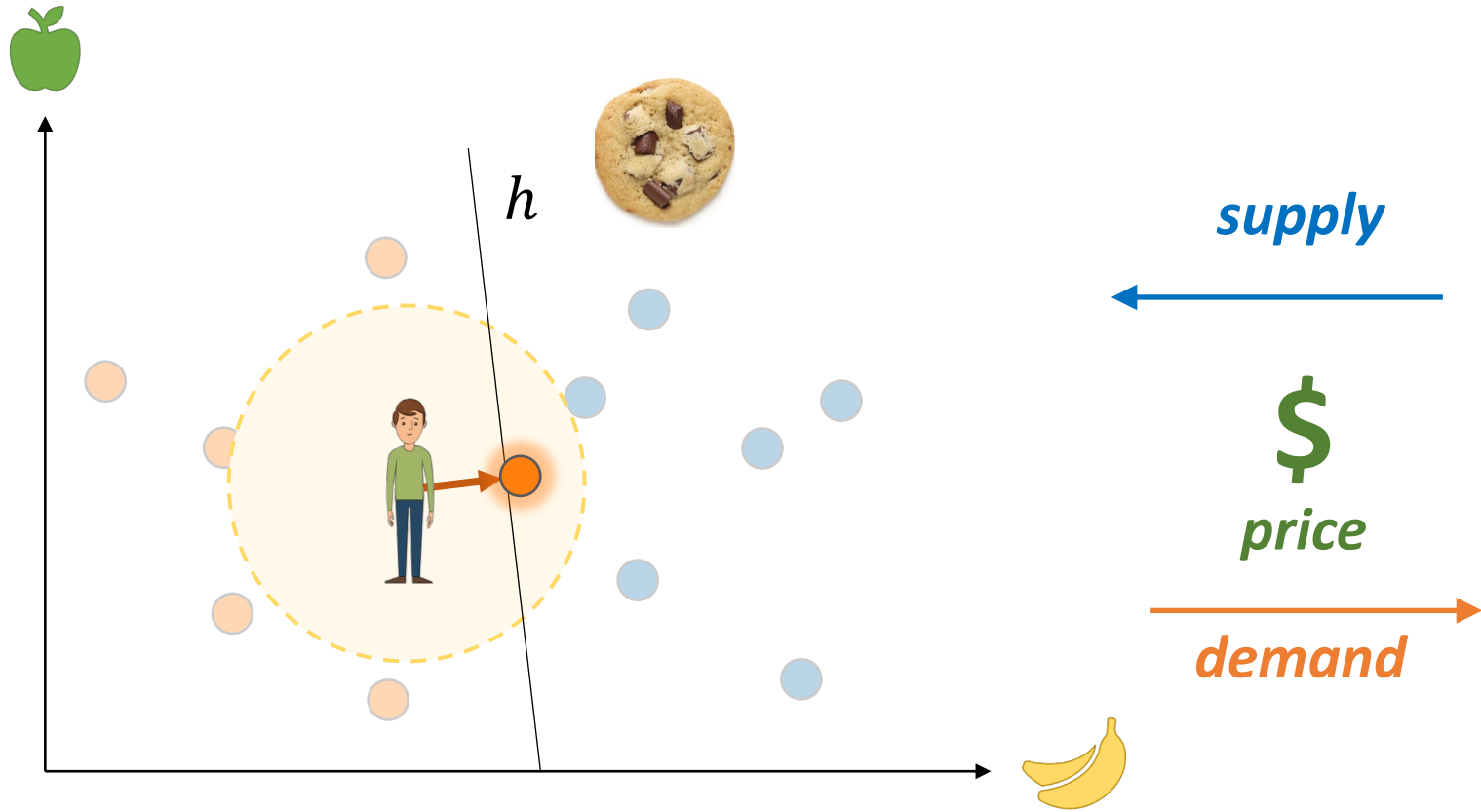
ask: where do fruits come from?





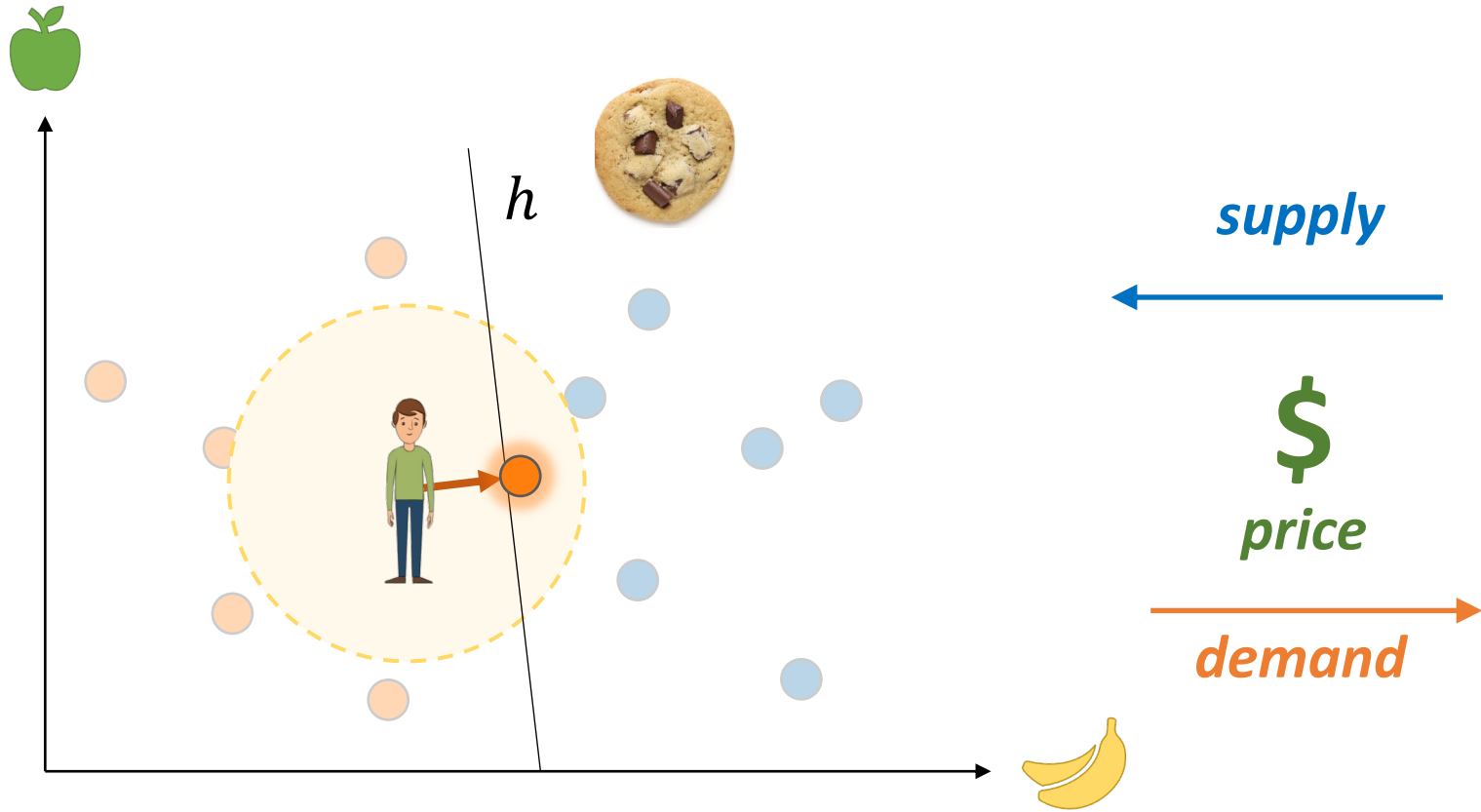
$$\Delta_h(x) = \underset{x'}{\operatorname{argmax}} h(x') - c(x, x')$$

utility
cost of apples and bananas



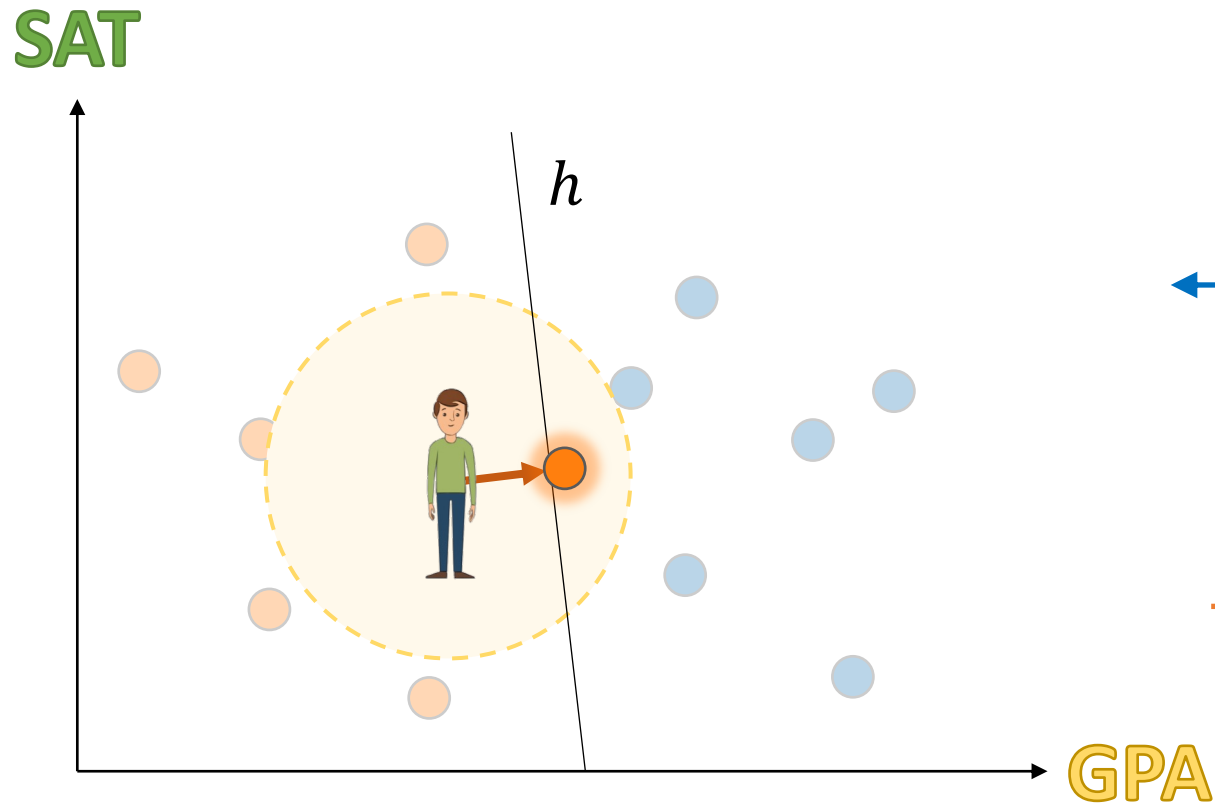
$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

utility
price of apples and bananas



$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

utility
price of apples and bananas

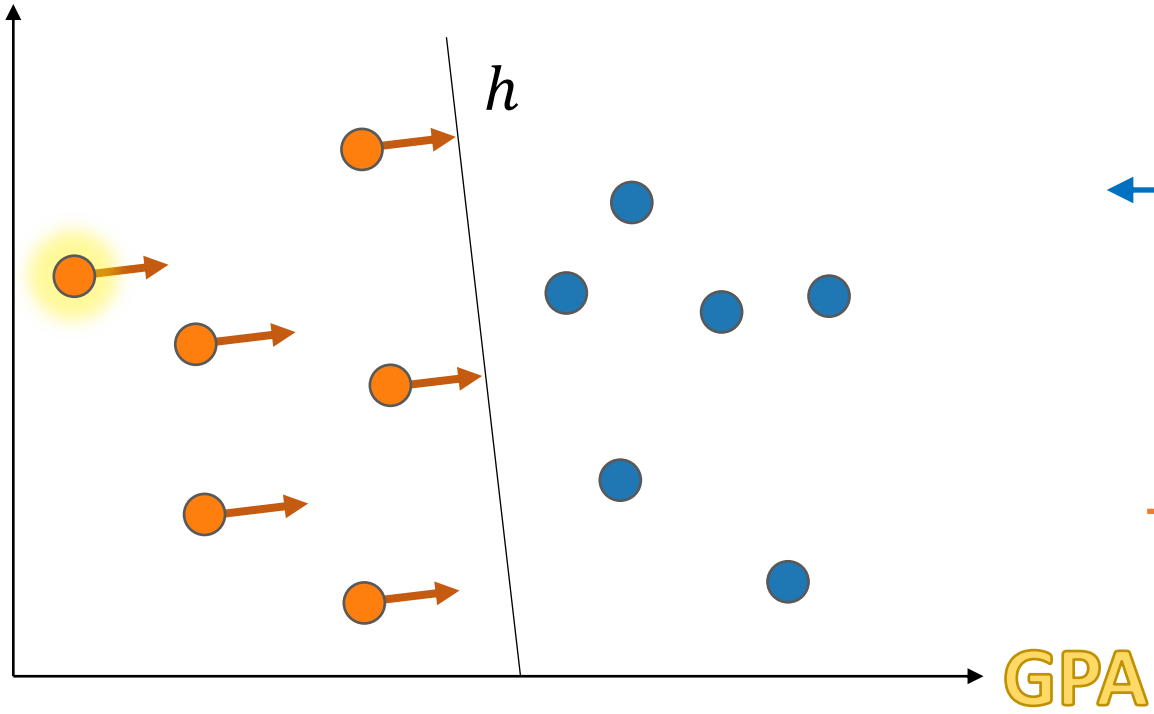


$$\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$

utility
price of prep courses

claim: classifier induce markets!

SAT



supply



\$
price

demand



**SAT
PREP
COURSE**

March 18, 25,
April 1, 15, 22, 29

Fee: \$80 payable by
check only
Please make checks payable
to Laura Lavacca

Students are asked to
purchase the Official SAT
Study Guide 2020 Edition.

Includes:
8 Real Tests
and official answer
explanations

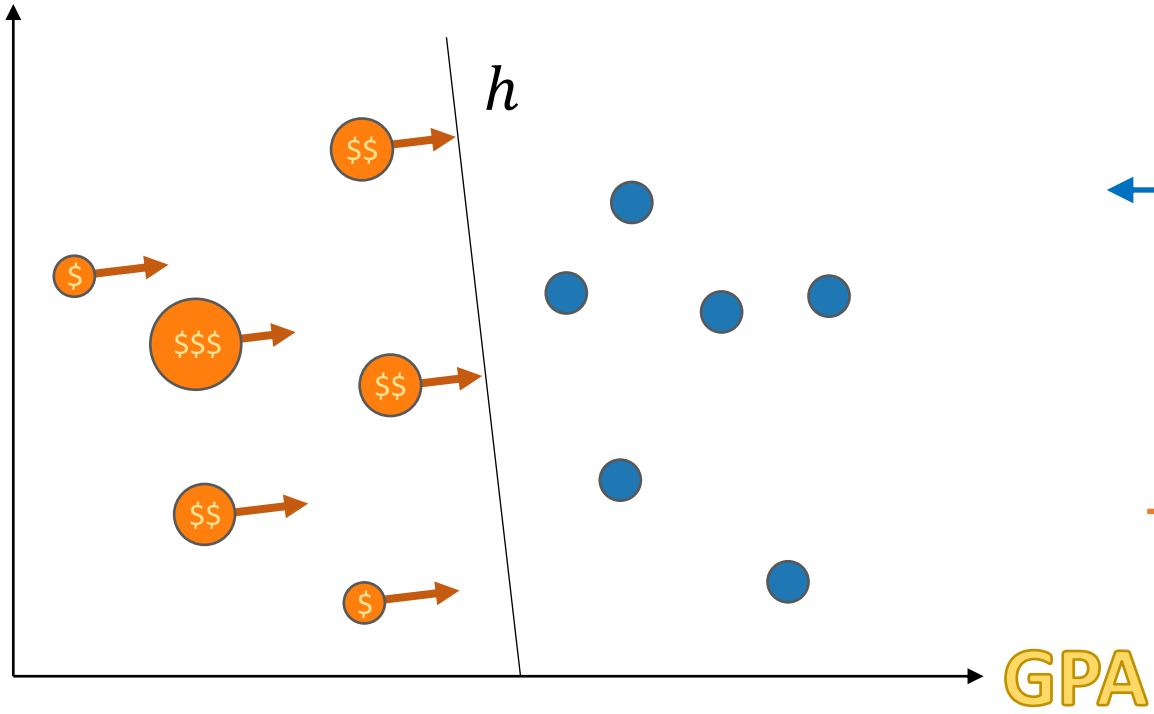
6-WEEK COURSE INCLUDES PRACTICE TEST & FEEDBACK

THE OFFICIAL
SAT
STUDY GUIDE

Prepare for the SAT
with sample questions,
practice tests, and more.

2020 EDITION

SAT



**SAT
PREP
COURSE**

March 18, 25,
April 1, 15, 22, 29

Fee: \$80 payable by
check only
Please make checks payable
to Laura Lavacca

Students are asked to
purchase the Official SAT
Study Guide 2020 Edition.

6-WEEK COURSE INCLUDES PRACTICE TEST & FEEDBACK

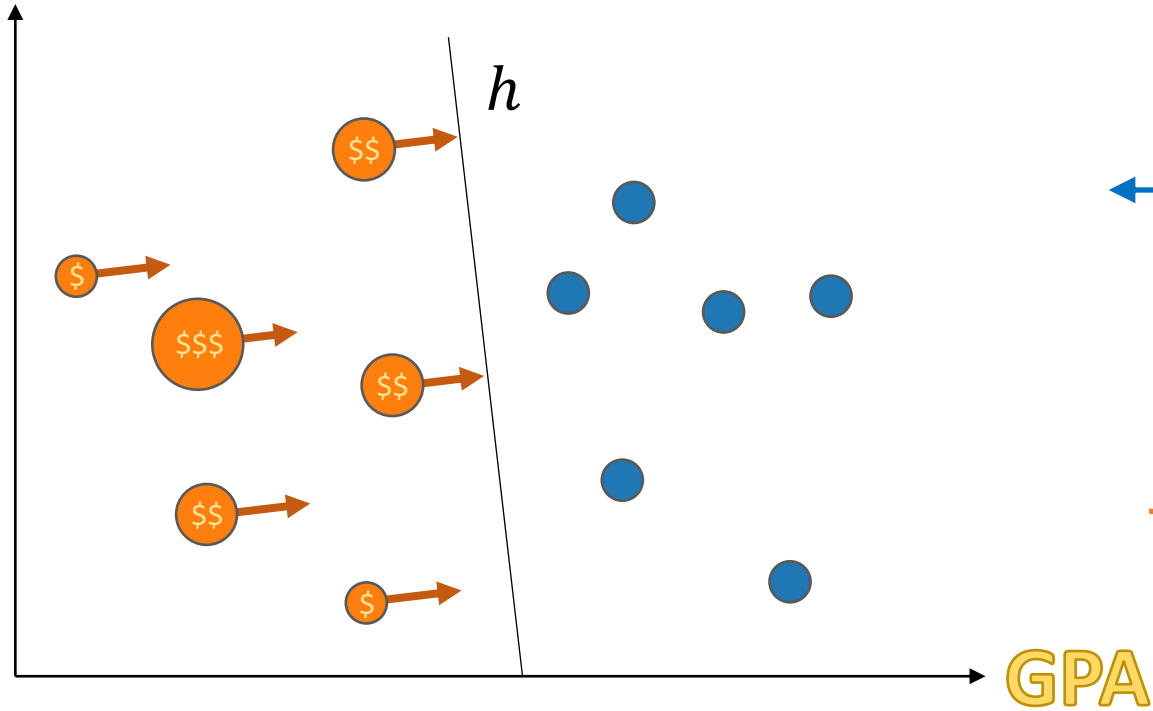
THE OFFICIAL
SAT
STUDY GUIDE

Prepare for the SAT
with sample questions,
practice tests, and more.

Includes:
8 Real Tests
and official answer
explanations

2020 EDITION

SAT



supply



\$
price

demand



**SAT
PREP
COURSE**

March 18, 25,
April 1, 15, 22, 29

Fee: \$80 payable by
check only
Please make checks payable
to Laura Lavacca

Includes:
8 Real Tests
and official answer
explanations

Students are asked to
purchase the Official SAT
Study Guide 2020 Edition.

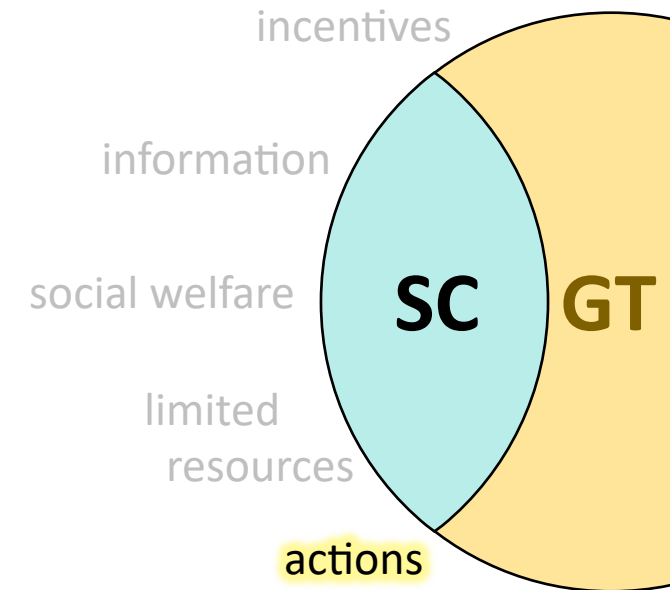
6-WEEK COURSE INCLUDES PRACTICE TEST & FEEDBACK

ask: can learning anticipate and account for the markets it induces?

strategic **modification**:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

$$\text{s.t. } \Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$$



ask: what other actions can users take?

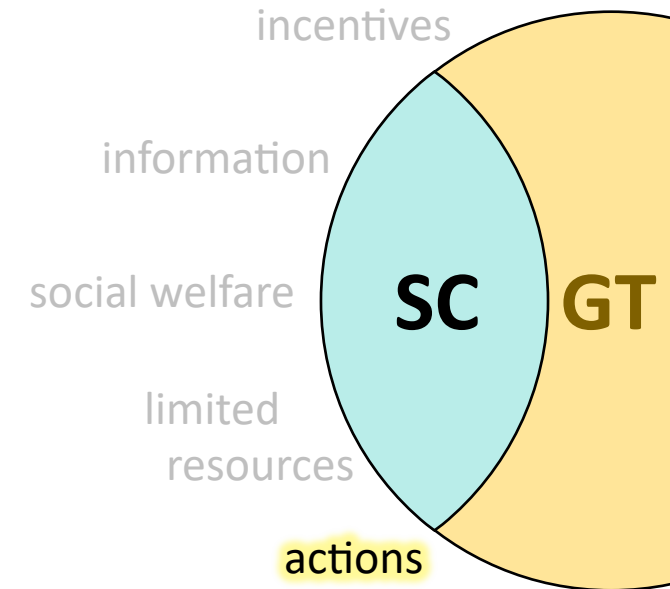
strategic **modification**:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

strategic **participation**:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m a_h(x_i) \ell(y_i, h(x_i))$$

s.t. $a_h(x) = \mathbb{1}\{\text{worthwhile to apply}\}$



ask: what other actions can users take?

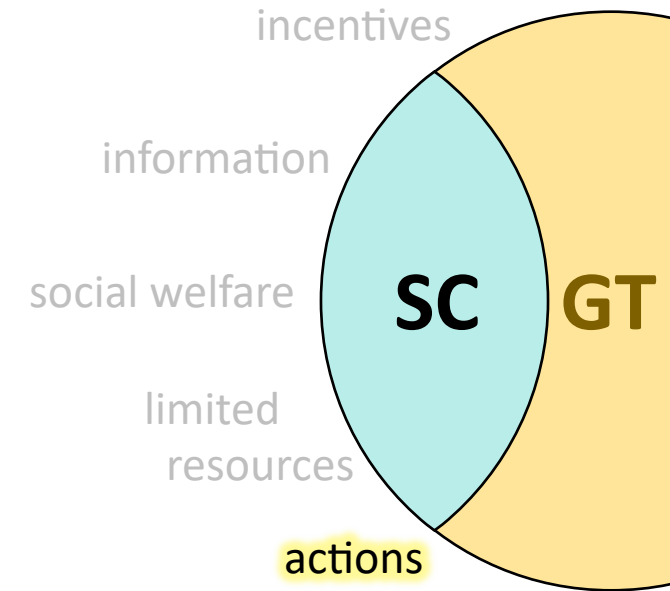
strategic **modification**:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(\Delta_h(x_i)))$$

strategic **participation**:

$$\operatorname{argmin}_h \frac{1}{m} \sum_{i=1}^m a_h(x_i) \ell(y_i, h(x_i))$$

s.t. $a_h(x) = \mathbb{1}\{\text{worthwhile to apply}\}$



ask: what other actions can users take?

strategic participation:

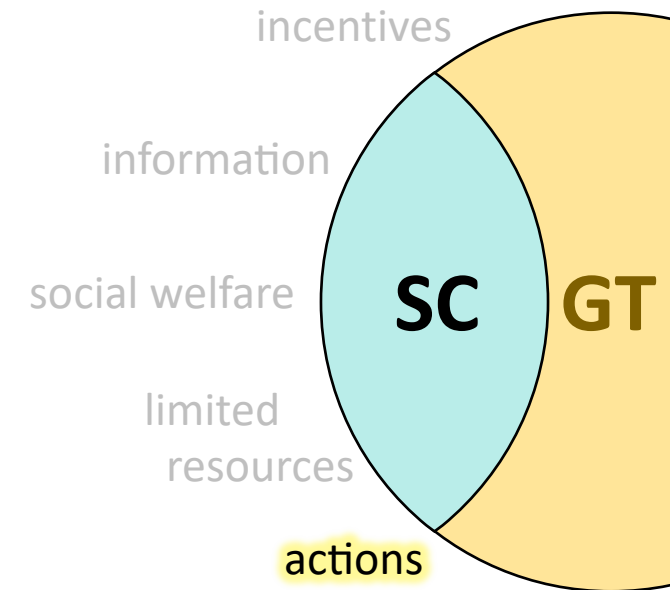


test

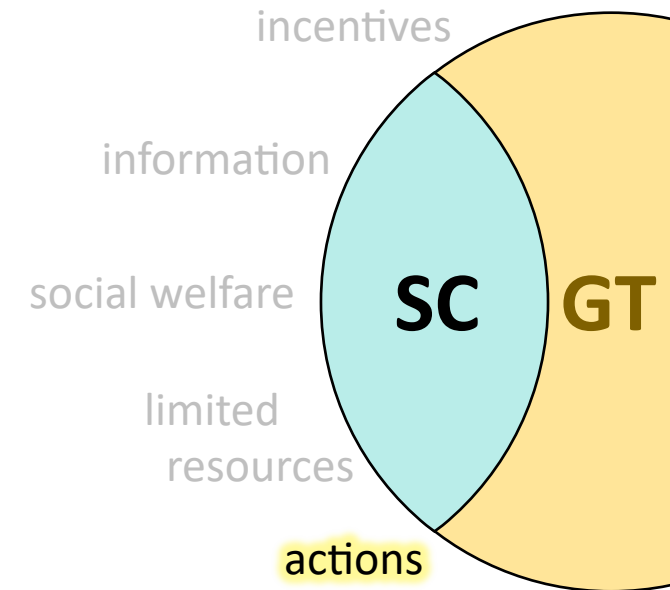
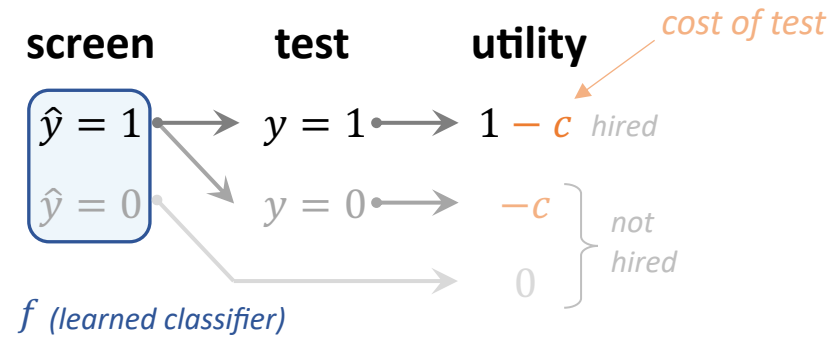
$$y = 1$$

$$y = 0$$

(interview, trial period, ...)

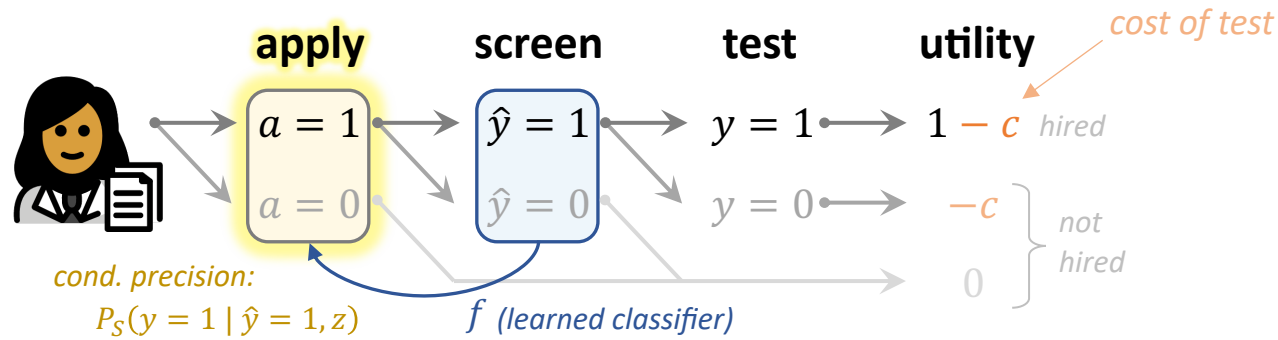


strategic participation:

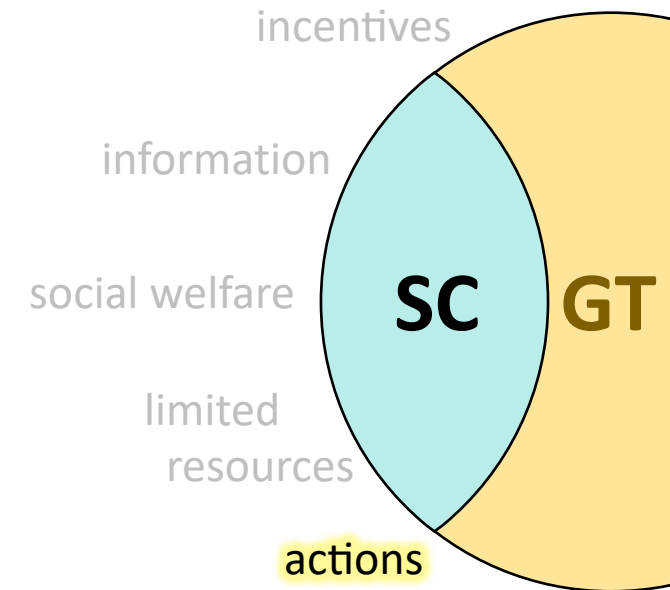
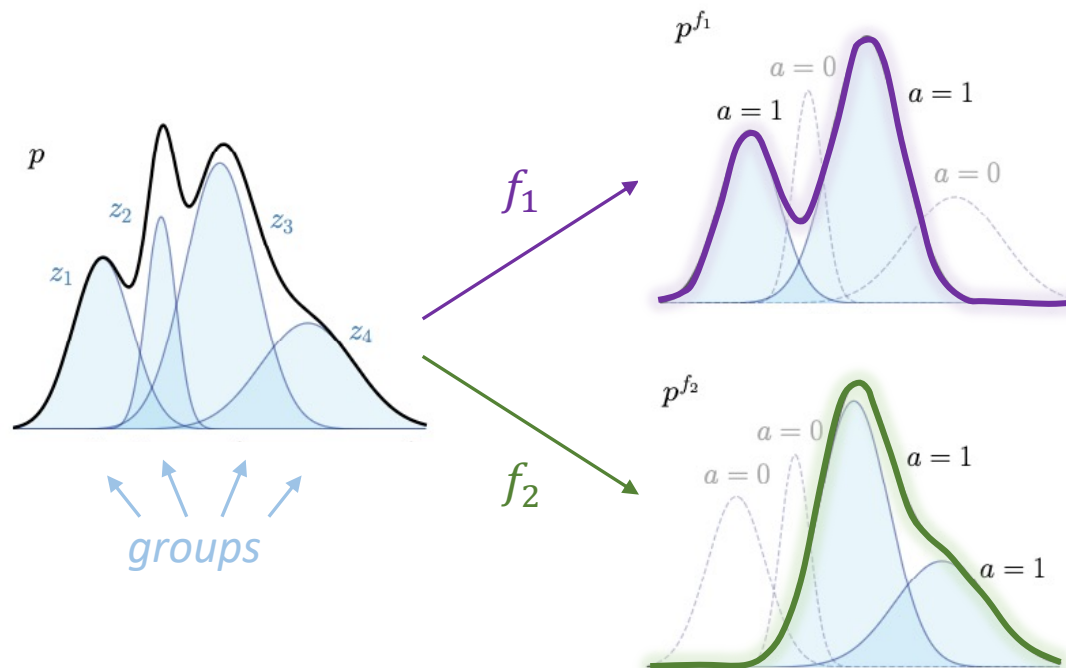


ask: how does learning affect applications?

strategic participation:

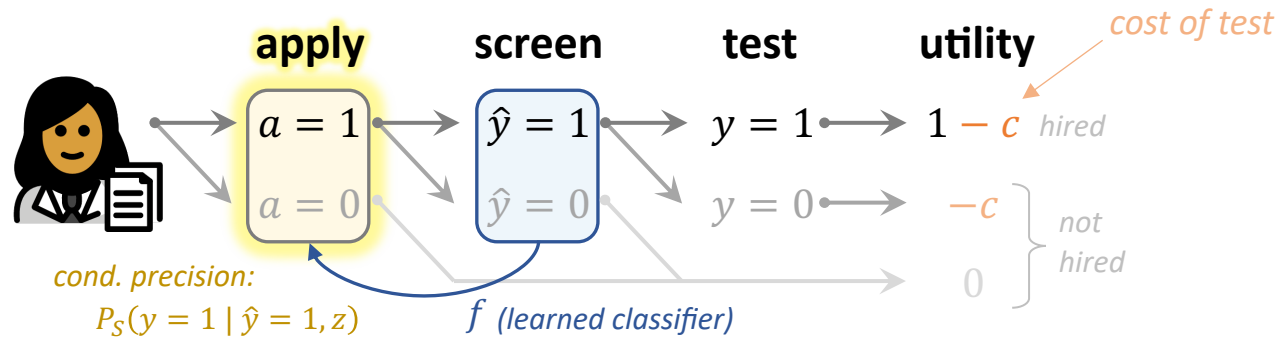


- **Observation:** learning rule determines **self-selection**

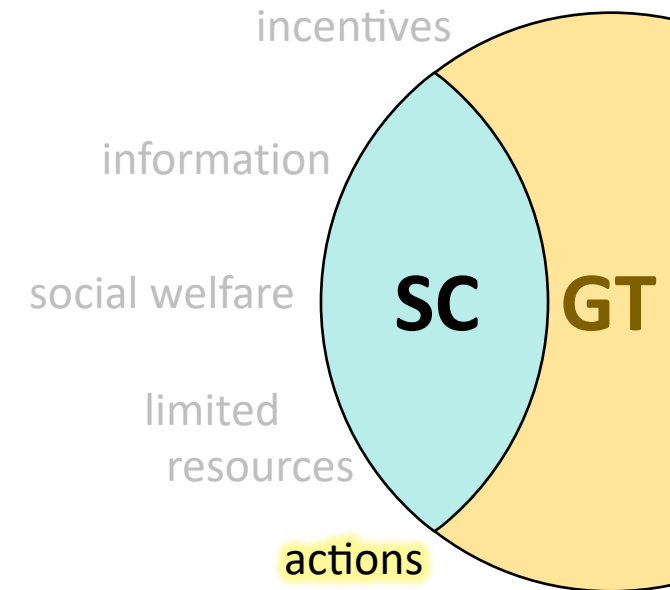
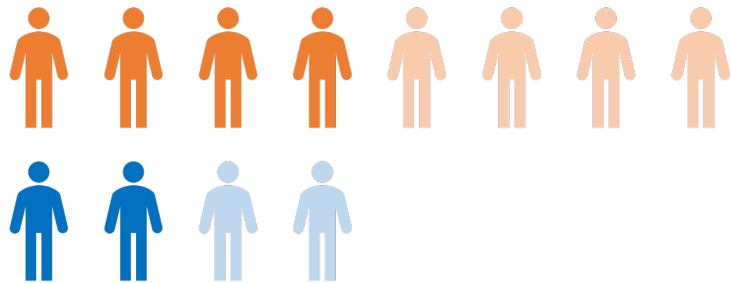


main result: learning has capacity to **fully determine applications!**

strategic participation:

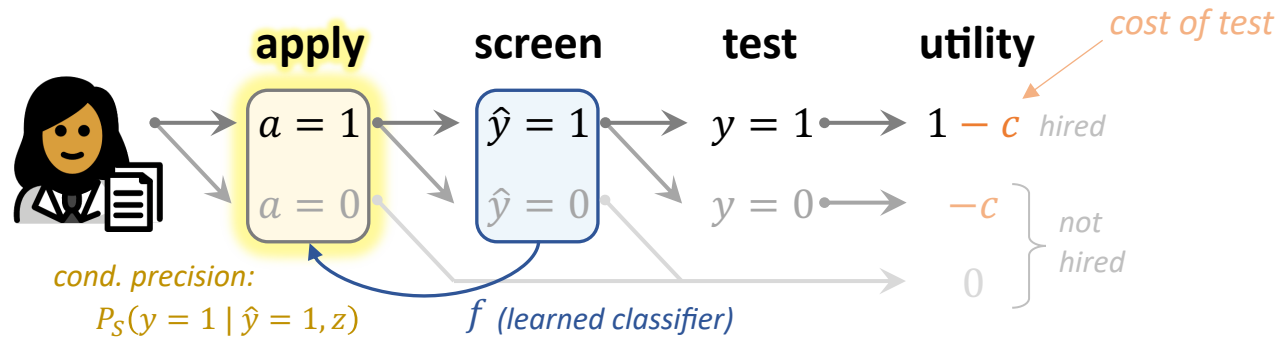


- **Observation:** learning rule determines **self-selection**
- **Implications:** can create *false appearance* of fairness, discriminates by making application too costly/risky

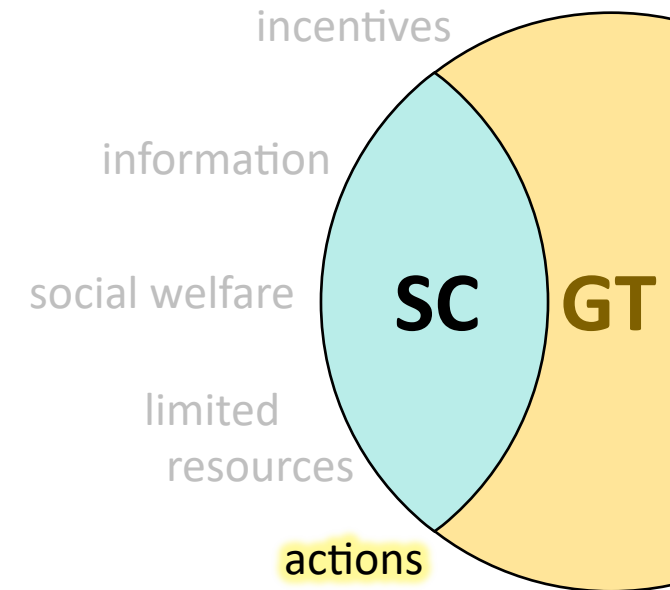
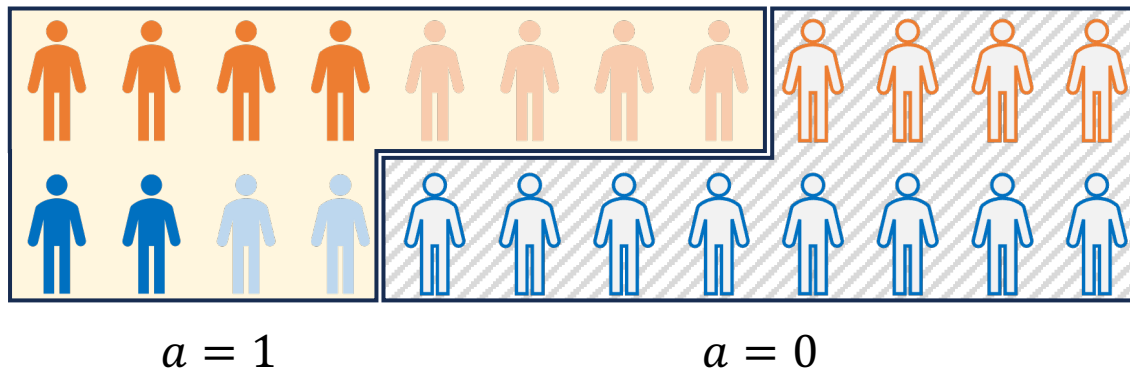


ask: how does learning affect applications?

strategic participation:



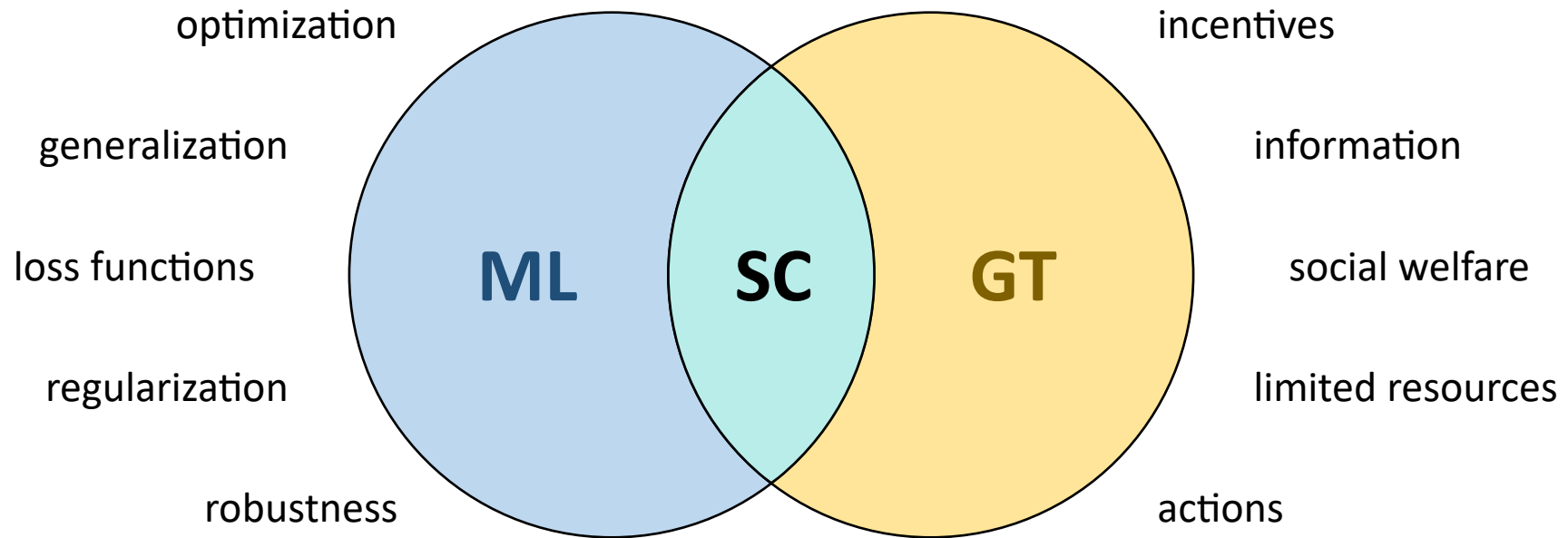
- **Observation:** learning rule determines **self-selection**
- **Implications:** can create *false appearance* of fairness, discriminates by making application too costly/risky

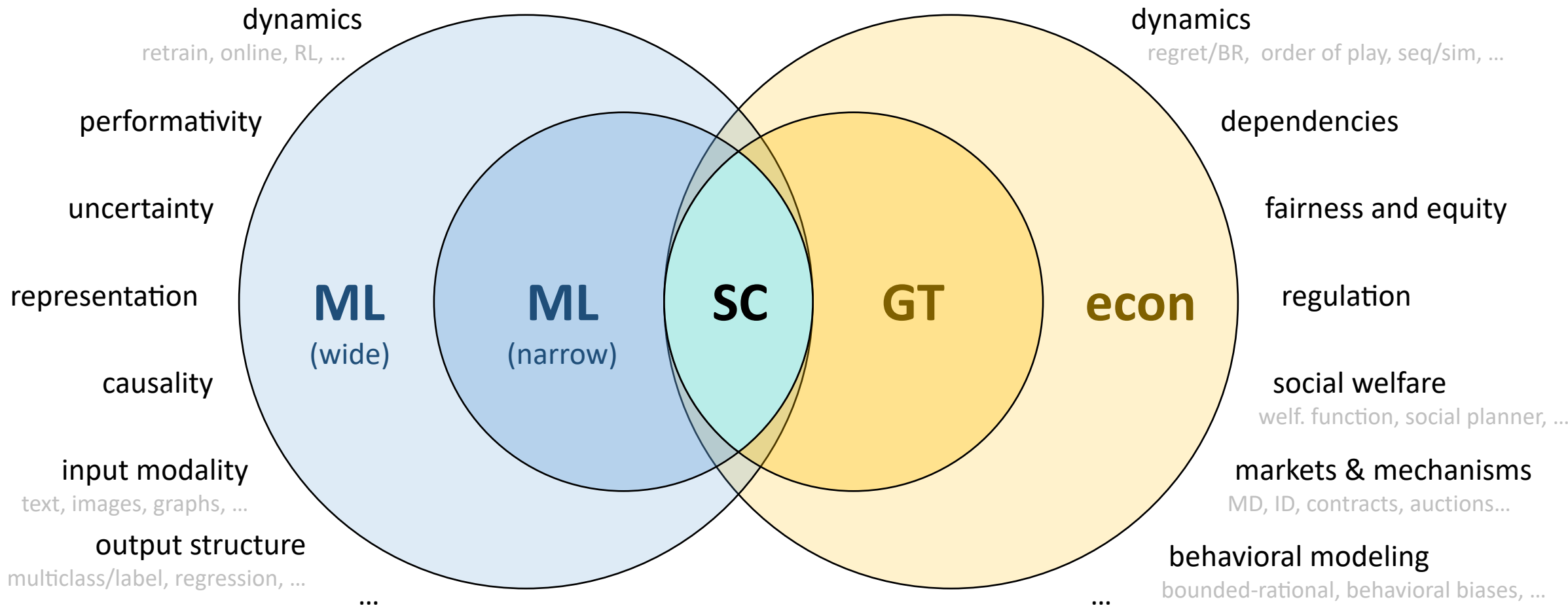


ask: how does learning affect applications?

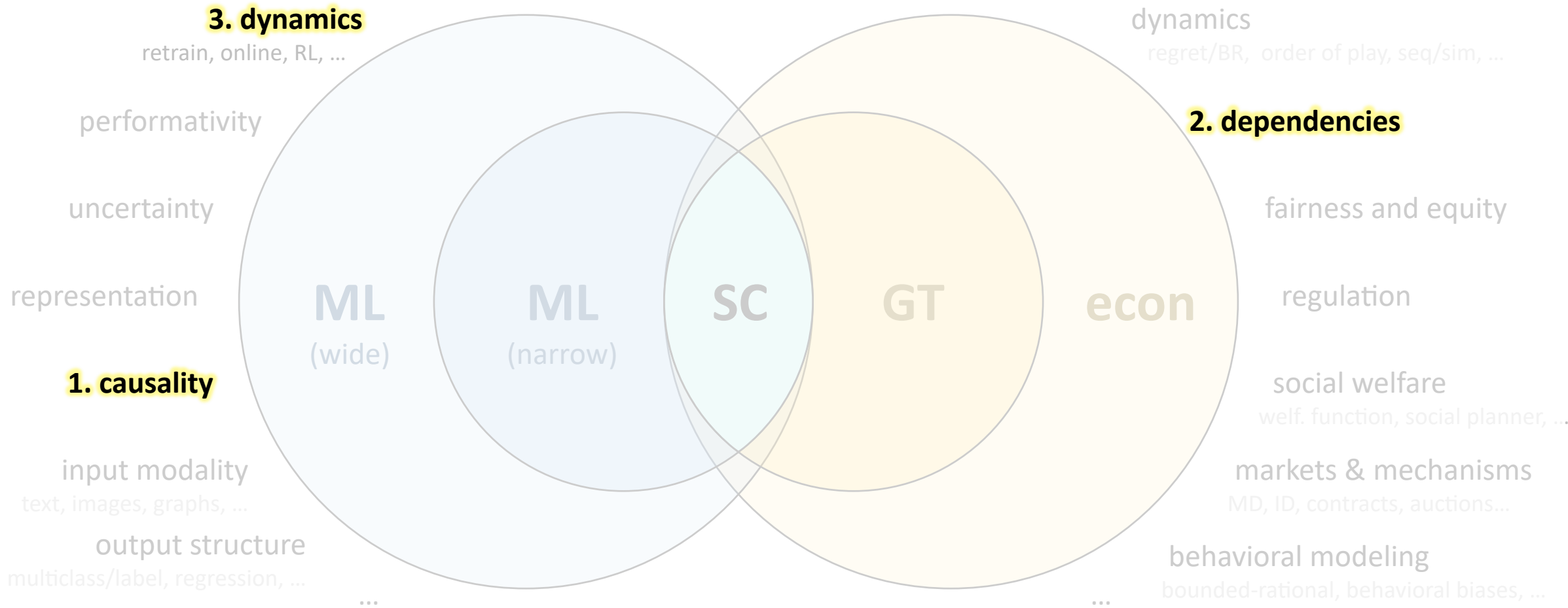
Beyond

the standard setup





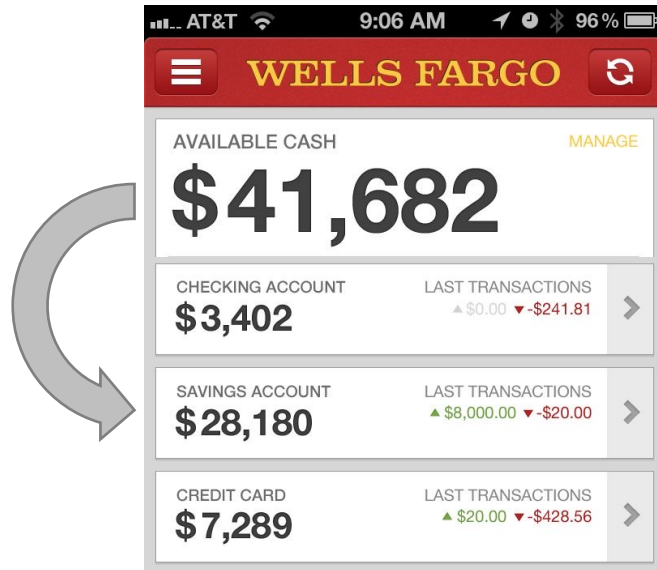
revisit old fronts + tackle new ones!



revisit old **fronts** + tackle new ones!

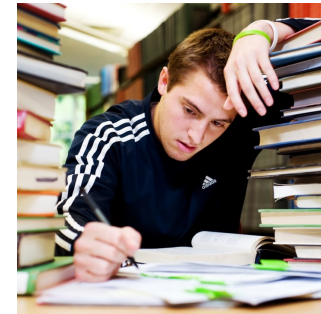
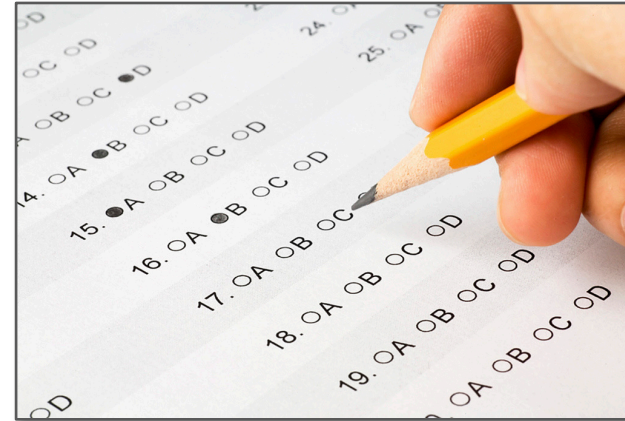
1) Causality

vanilla SC



*superficial changes
⇒ gaming*

causal SC

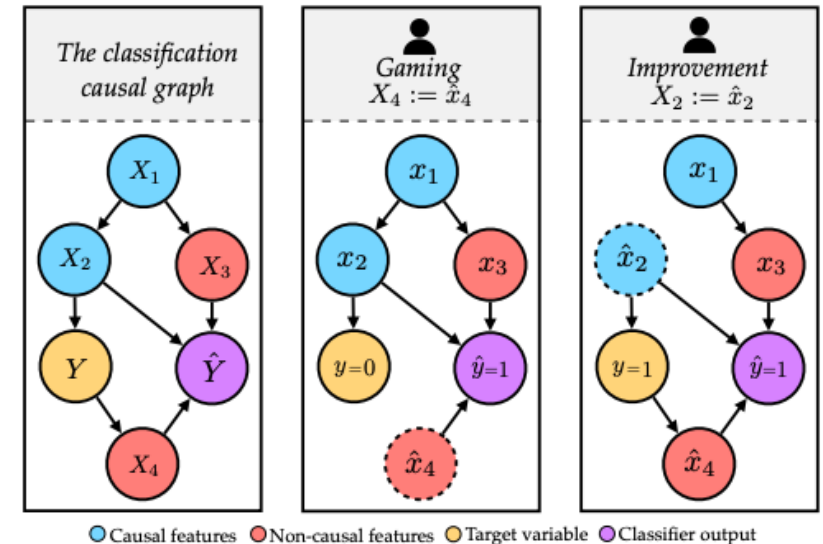


1) Causality

- **Standard SC:** changing x does not affect y (=gaming)
- **More realistic:** changing x can also change y
- Assume exists underlying *causal graph* [Pearl 2009]:

ask: can we learn in causal strategic settings?

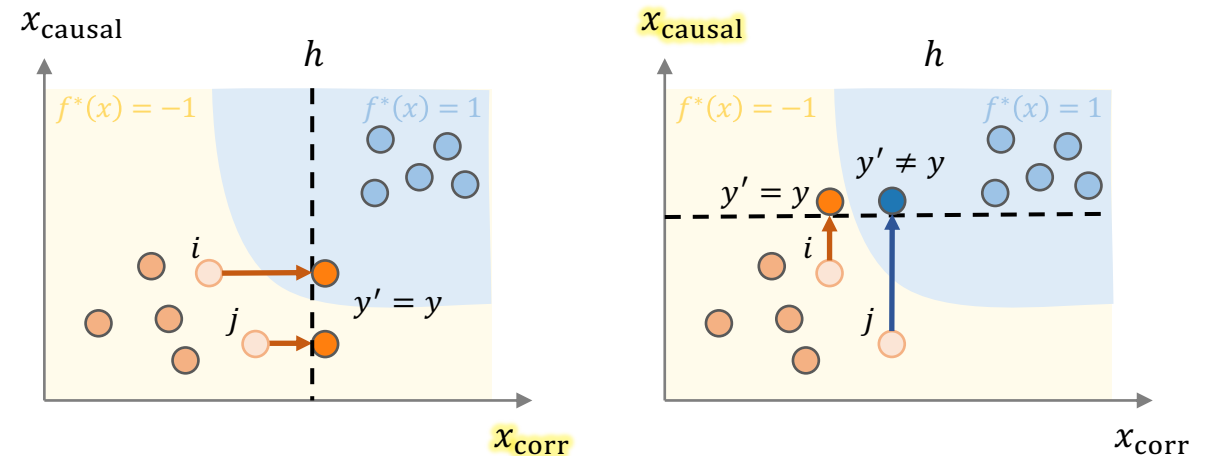
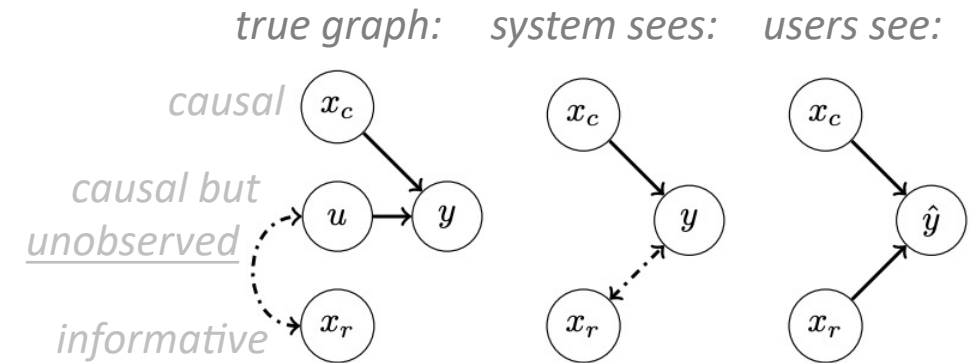
- **Lots of challenges:**
 - graph not necessarily known
 - key variables not necessarily observed (e.g., confounders)
 - structure determines interactions (i.e., what affects what)
- Causal SC is **inherently difficult** – as hard as *causal inference* [MMH ICML20]



(taken from Miller et al. 2020)

Causal SC as distribution shift

- **Q1:** How does **causality** affect **learning**?
- Simplifying assumption: causal vs. correlative features
- **A1:** Entails different **types of distribution shift**:
 - correlative \rightarrow *strategic* shift \rightarrow gaming
 - only causal \rightarrow *covariate* shift \rightarrow missinformation
 - both \rightarrow *mixture* shift \rightarrow interactions
- **Corollary:** choose your battles!

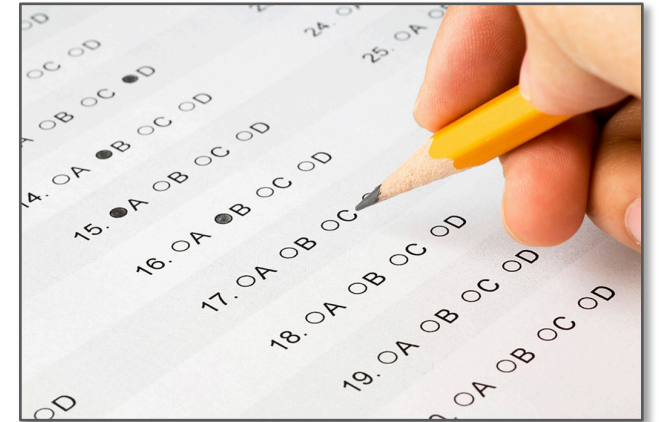


Incentivizing improvement

- **Q1:** How does **causality** affect **learning**?
- **Q2:** How does **causality** affect **social outcomes**?
- **A2:** Causal SC has **potential for *improvement***:

$$\mathbb{E}_x[\mathbb{E}[p(y | \text{do}(\Delta)) - p(y) | x]]$$

- **Goal:** learn h that (also) promotes improvement
- Has long and rich history in economics (e.g., see [KR 19])
- Also considered in (online) SC
(e.g., [SEA'20, BLWZ'21, CWL'21, HNSHW'22, MDW'22])

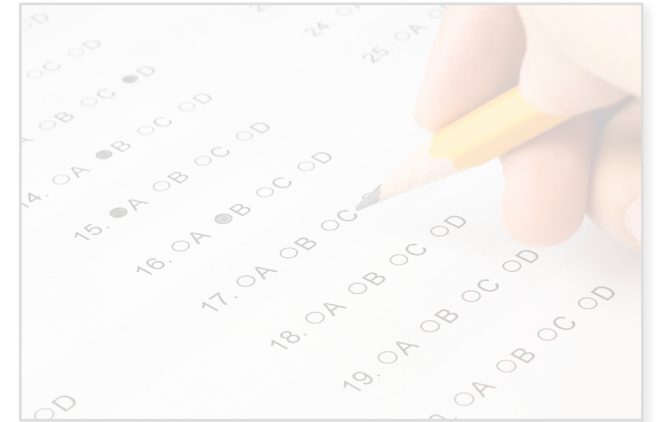


Incentivizing improvement

- **Q1:** How does **causality** affect **learning**?
- **Q2:** How does **causality** affect **social outcomes**?
- **A2:** Causal SC has **potential for *improvement***:

$$\mathbb{E}_x[\mathbb{E}[p(y | \text{do}(\Delta)) - p(y) | x]]$$

- **Goal:** learn h that (also) promotes improvement
- Has long and rich history in economics (e.g., see [KR 19])
- Also considered in (online) SC
(e.g., [SEA'20, BLWZ'21, CWL'21, HNSHW'22, MDW'22])
- But changing x can also **impair** outcomes!



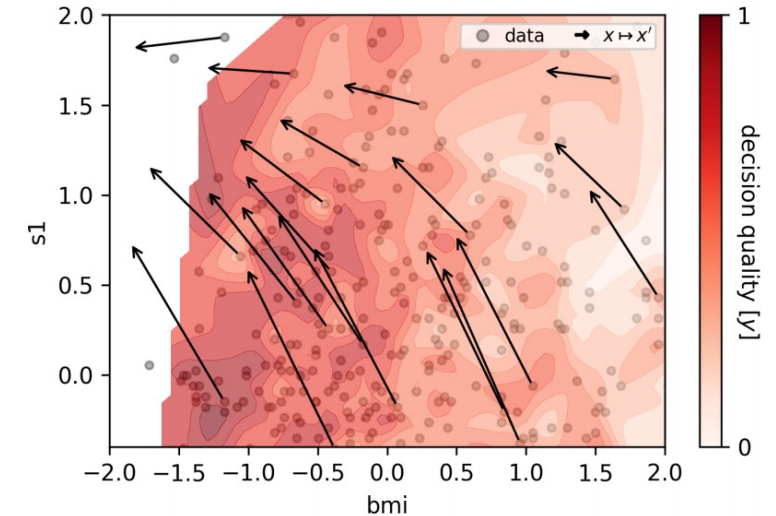
Incentivizing improvement

- **Q1:** How does **causality** affect **learning**?
- **Q2:** How does **causality** affect **social outcomes**?
- **A2:** Causal SC has **potential for improvement**:

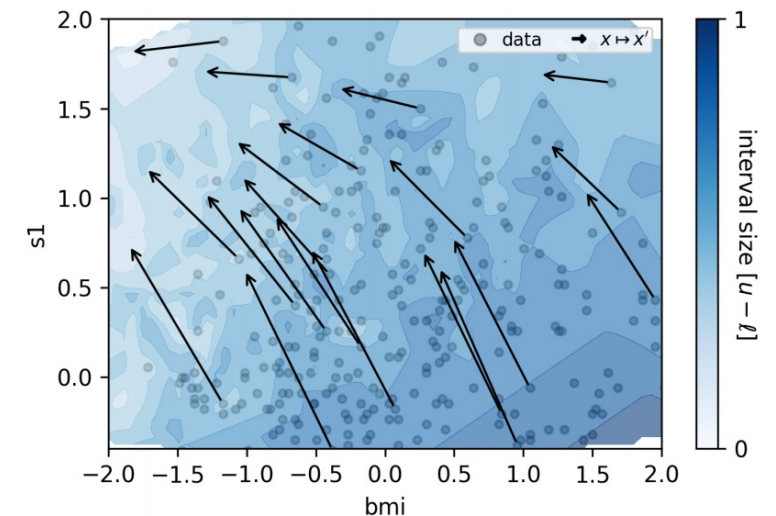
$$\mathbb{E}_x[\mathbb{E}[p(y | \text{do}(\Delta)) - p(y) | x]]$$

- **Goal:** learn h that (also) promotes improvement
- Has long and rich history in economics (e.g., see [KR 19])
- Also considered in (online) SC (e.g., [SEA'20, BLWZ'21, CWL'21, HNSHW'22, MDW'22])
- But changing x can also **impair** outcomes!
- **Solution:** learn safe models by “looking ahead”

causal effect:



uncertainty:



2) Dependencies: How users relate

- **Standard SC:** responses are independent ($\Delta_f(x)$ depends only on x)
- **More realistic:** responses are interdependent

- **Reason #1:** limited resources
 - Actually, all common examples have limit on # of $\hat{y} = 1$
 - This means that users **compete**

*limited
financial
resources*



*limited
teaching
capacity*

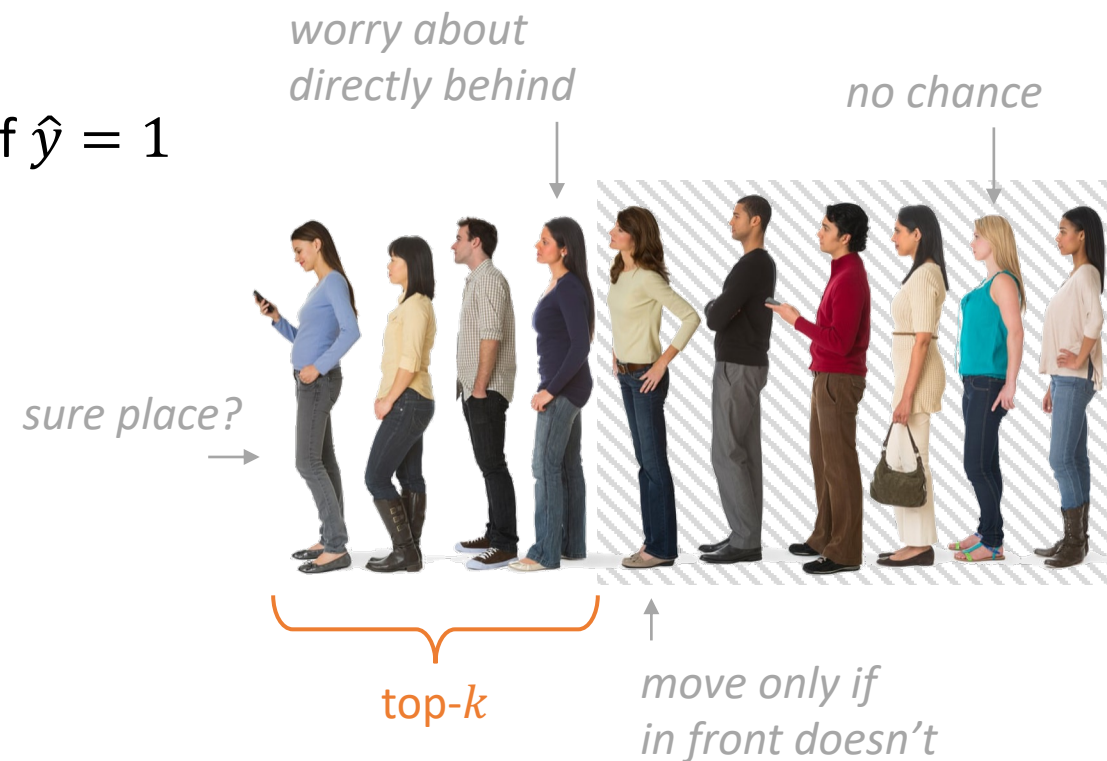


*limited
qualified
personell*



2) Dependencies: How users relate

- **Standard SC:** responses are independent ($\Delta_f(x)$ depends only on x)
- **More realistic:** responses are interdependent
- **Reason #1:** limited resources
 - Actually, all common examples have limit on # of $\hat{y} = 1$
 - This means that users **compete**
 - **Reasonable approach:**
learn to **rank**, then set $\hat{y} = 1$ only for **top- k**
 - Turns out to be *exceedingly hard* [LGB ICML22]
 - Still – major goal!



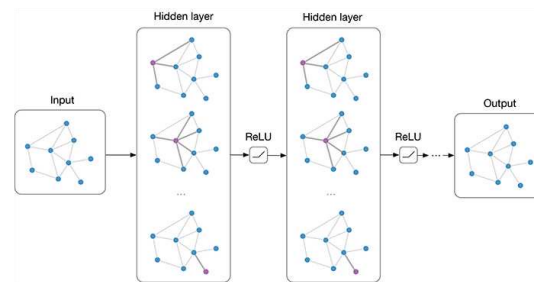
2) Dependencies: How users relate

- **Standard SC:** responses are independent
- **More realistic:** responses are interdependent
- Reason #1: limited resources
- Reason #2: model-induced dependencies

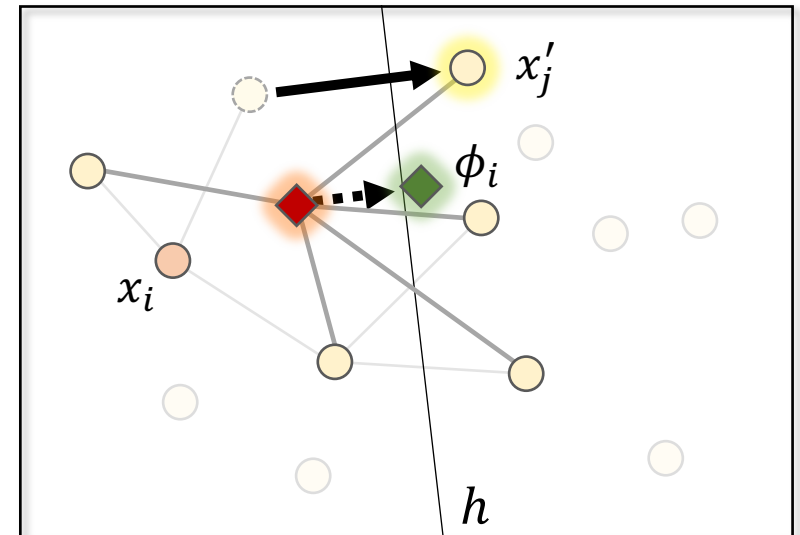
using graph in learning
creates dependencies



social network



GNN

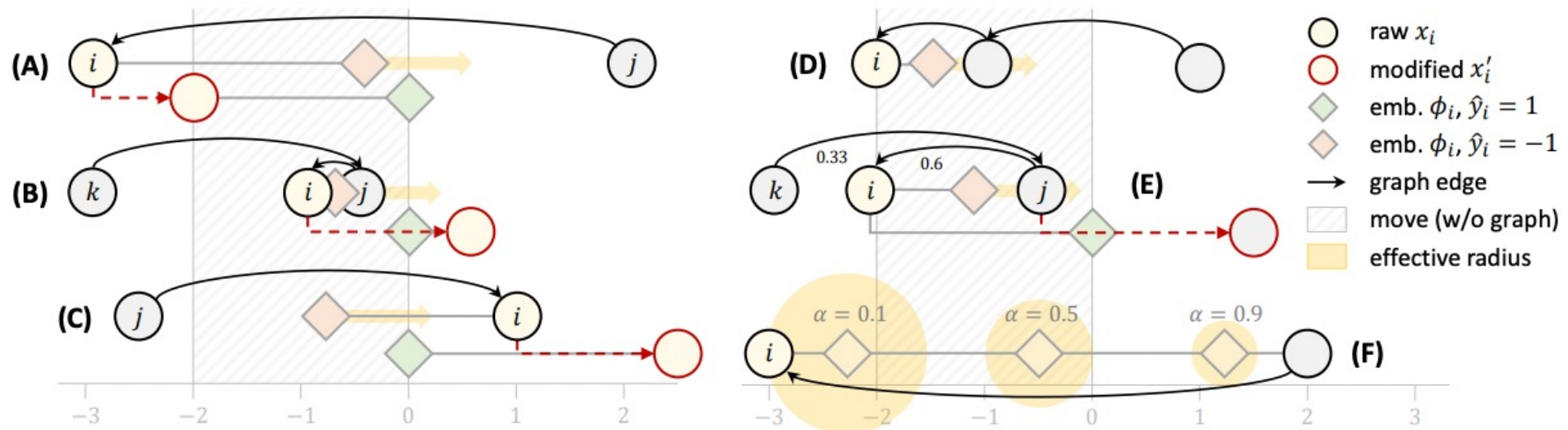


graph-dep. embedding

2) Dependencies: How users relate

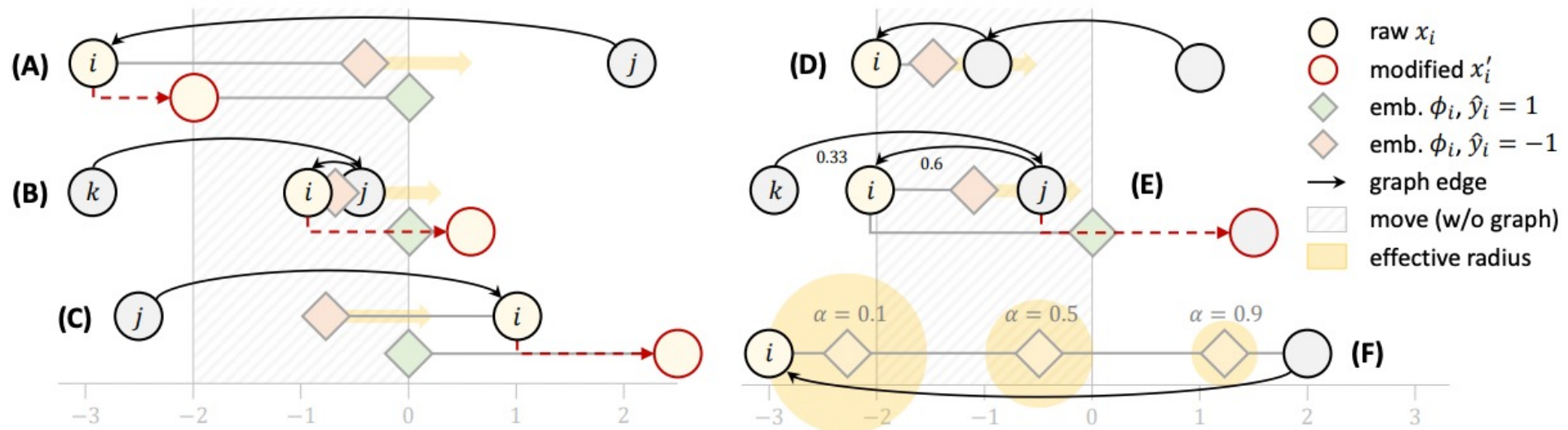
- **Standard SC:** responses are independent
- **More realistic:** responses are interdependent
- Reason #1: limited resources
- Reason #2: model-induced dependencies

using graph in learning
creates **surprising** dependencies



2) Dependencies: How users relate

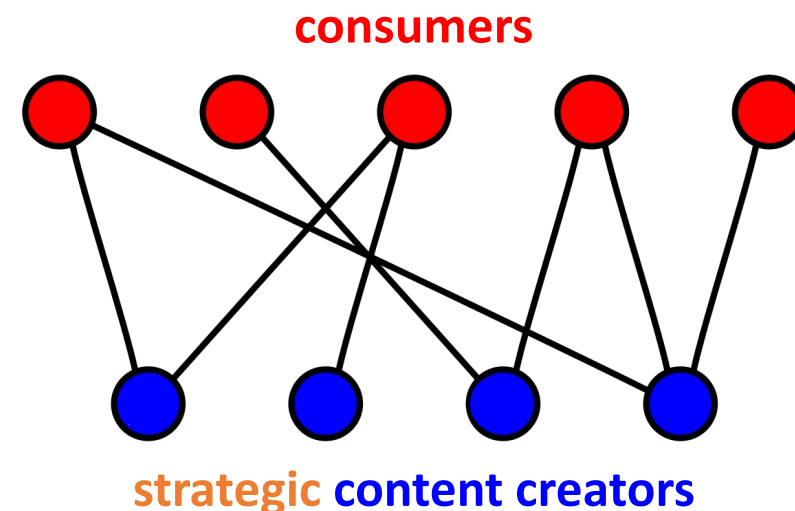
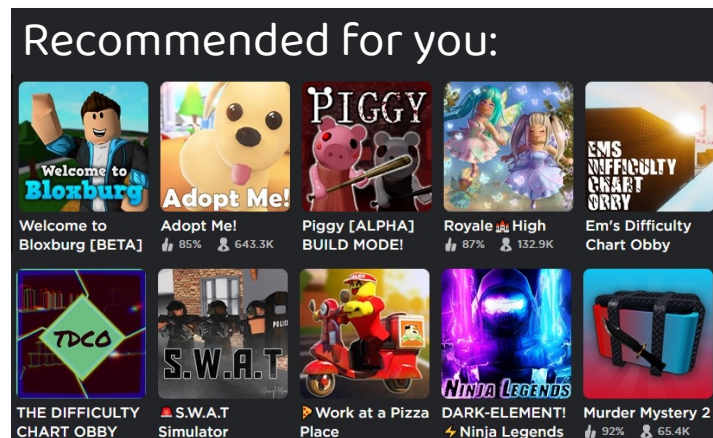
- **Standard SC:** responses are independent
- **More realistic:** responses are interdependent
- Reason #1: limited resources
- Reason #2: model-induced dependencies



2) Dependencies: How users relate

- **Standard SC:** responses are independent
- **More realistic:** responses are interdependent
- Reason #1: limited resources
- Reason #2: model-induced dependencies
- **Reason #3:** economic graph structure

main result: can use graph to incentiveize **diversity**



3) Learning over time

- **Standard SC:** batch setting: **train** → **deploy** → **test**
- Assumes access to **clean data** (otherwise, chicken & egg!)
- **More realistic:** data is dirty (i.e., result of some behavior)

- **One solution:** iterated deployments over time: **train** → **deploy** → **train** → **deploy** → **train** → ...
- **Three main approaches:** → *lots of research; will present here only in brief*
 1. online learning (e.g., bandits)
(e.g., [DRSWW NeurIPS17, CSSVZ ICML23, HPW NeurIPS23, SBM NeurIPS23, ABBN EC21, ...])
 2. performative prediction (retraining revisited) [PZMH ICML20]
 3. dynamical systems

3) Learning over time

- **Standard SC:** batch setting: **train** → **deploy** → **test**
- Assumes access to **clean data** (otherwise, chicken & egg!)
- **More realistic:** data is dirty (i.e., result of some behavior)

- **One solution:** iterated deployments over time: **train** → **deploy** → **train** → **deploy** → **train** → ...
- **Pros:** less restrictive
 - (1) does not require clean data
 - (2) does not assume known Δ_h (or even best-response)
 - (3) permits causal Δ_h (under additional assumptions)
- **Cons:** each deployment is social “experiment”
 - in some cases, exploration is reasonable
 - in other cases – it is very much not

Opportunities & challenges

open questions

Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects:



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects:

- labels beyond binary
 - regression
 - multiclass
 - multilabel
 - sequences
 - structured (e.g., graphs)
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects:

- labels beyond binary
- inputs beyond vectors
 - images
 - text
 - graphs
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects:

- labels beyond binary
- inputs beyond vectors
- models beyond linear
 - neural nets (behavior in latent space)
 - tree-based
 - text-based (prompts)
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects:

- labels beyond binary
- inputs beyond vectors
- models beyond linear
- settings beyond classification
 - unsupervised and semi-supervised
 - generative
 - RL and MARL
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects

2. Econ/GT aspects:

- information
 - power
 - control
 - selective release/withhold
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects

2. Econ/GT aspects:

- information
- other economic settings
 - markets, auctions, contracts, ...
 - competition (between classifiers)
 - cooperation (between users)
 - monopolistic behavior
 - ...



Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects

2. Econ/GT aspects:

- information
- other economic settings
- behavior
 - Bayesian
 - non-rational “behavioral” (=biases)
 - ...

loss aversion *confirmation bias*
future discounting *causal fallacy*
decoy effect *primacy/recency*
bounded rationality *anchoring* *endowment effect*
risk aversion/seeking *choice overload*
k-level reasoning *quantal response* *availability bias*
base rate fallacy *bandwagon effect*
framing/priming ...

Open questions

- Strategic learning is exiting new field with much potential for growth
- But it is also young – so that many challenges still lie ahead:

1. Learning aspects

2. Econ/GT aspects

3. “In the wild”:

- evaluation [BBK 20, HHP 23, CIALRM 23]
- measuring utility/welfare
- estimating costs
- monitoring and regulation

Why supervised learning?

- Most human-centric tasks are **policy problems** (vs. prediction problems)
- So supervised learning is clearly the wrong tool to use
- But it is also *by far* the most prevalent, accessible, and easy to use
- **Vision for the future:**

```
strat_clf = LogisticRegression(penalty='l2', C=0.01, max_iter=500,  
                              want='yhat=1', know='noisy_h', do='game')  
strat_clf.fit(X, Y, cost=c)  
pred = strat_clf.predict(X_test)
```

- **Goal:** make integrating human agency as seamless as possible
- **Not so easy!** And requires much caution and deliberation (c.f. fairness)

Summary

Summary

- SC captures natural tension between learning systems and their users
- Appealing interface between ML and GT – many open question!
- Original setup is clean and simple, but likely to narrow
- Nonetheless, flexible and modular: easy to extend, relax, and generalize

Summary

- SC captures natural tension between learning systems and their users
 - Appealing interface between ML and GT – many open question!
 - Original setup is clean and simple, but likely to narrow
 - Nonetheless, flexible and modular: easy to extend, relax, and generalize
-
- A call to rethink the design of learning algorithms for social settings
 - An opportunity to revise foundations using economic and behavioral modeling
 - High potential for real impact – much more work needed!

Narrative(s)



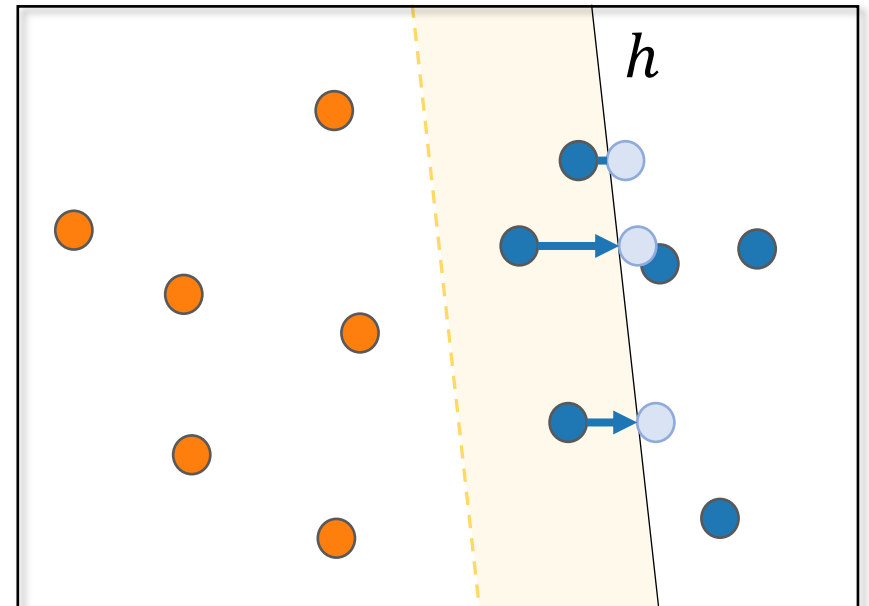
- “users game system”



Narrative(s)



- “users game system”
- “system exploits users”

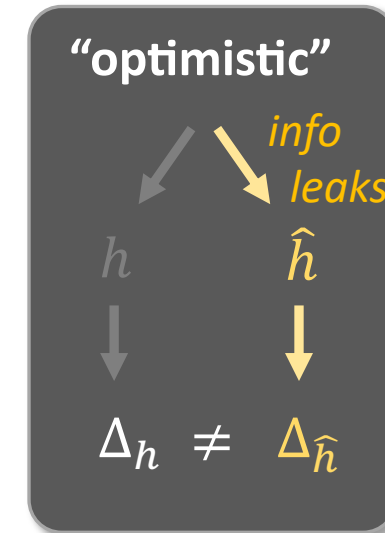


Narrative(s)



- “users game system”
- “system exploits users”
- “system exploits users *unintentionally*”
- “... as long as there is transparency”

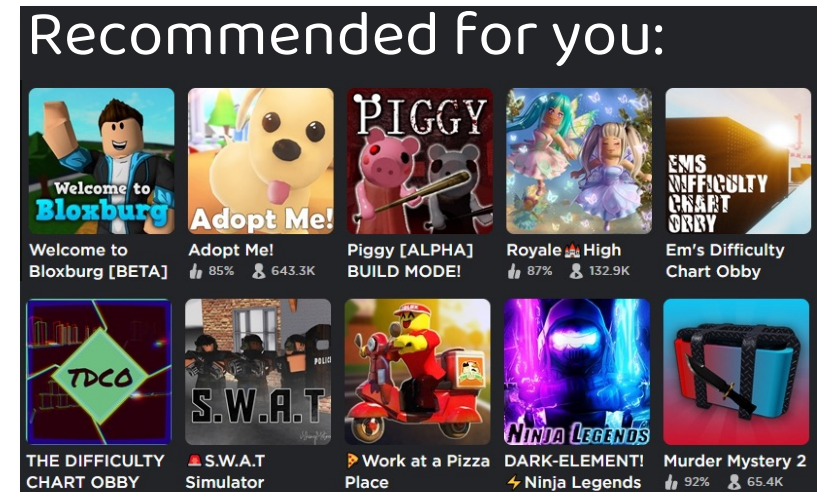
SC in the dark



Narrative(s)



- “users game system”
- “system exploits users”
- “system exploits users *unintentionally*”
- “... as long as there is transparency”
- “potential for cooperation...”



Narrative(s)



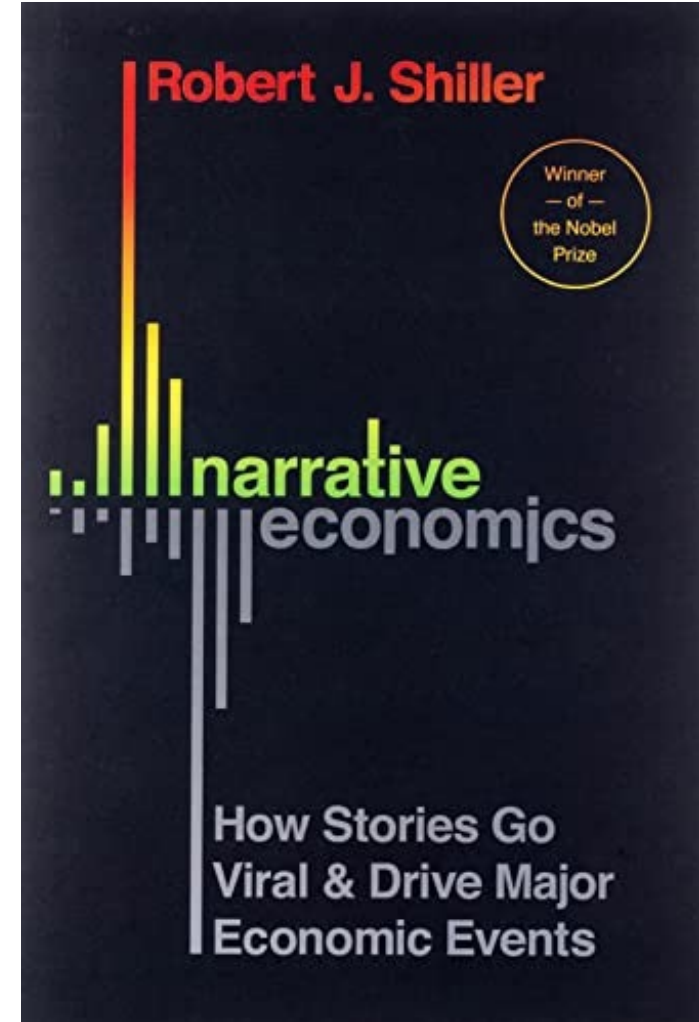
- “users game system”
- “system exploits users”
- “system exploits users *unintentionally*”
- “... as long as there is transparency”
- “potential for cooperation...”
- “its just a market”
- ...

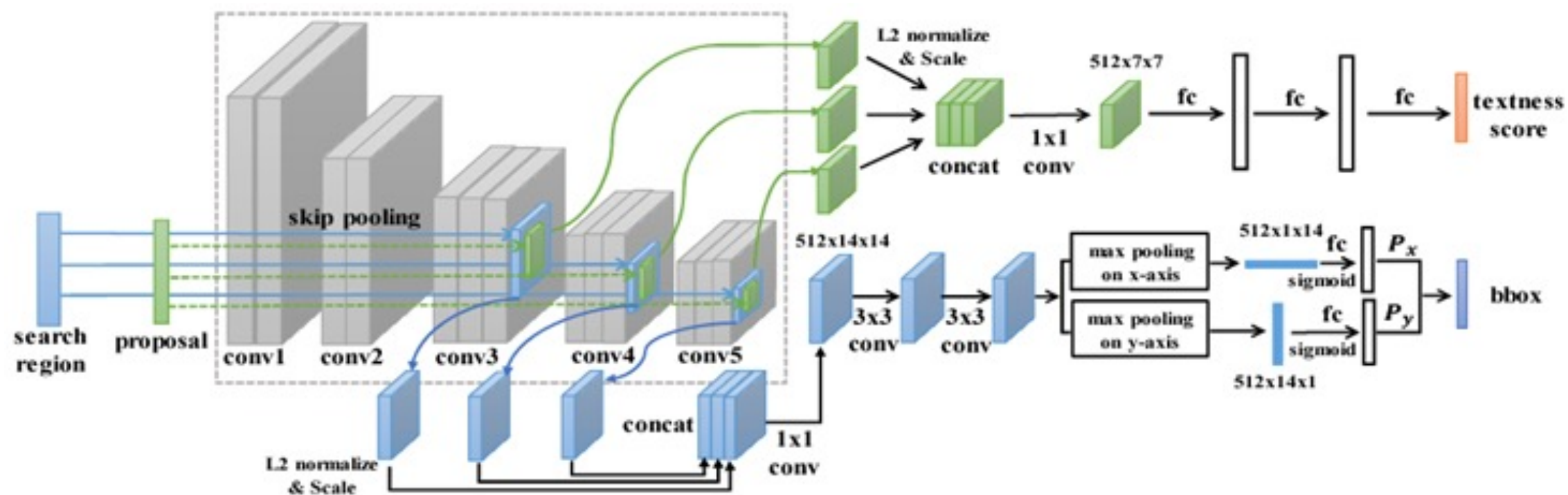


Narrative(s)



- “users game system”
- “system exploits users”
- “system exploits users *unintentionally*”
- “... as long as there is transparency”
- “potential for cooperation...”
- “its just a market”
- ...
-





LES ARTISTES ASSOCIÉS S.A.B. *Mérentent* PARLANT FRANÇAIS 105 ENFANTS ADMIS

**CHARLIE
CHAPLIN**

DANS



les **TEMPS MODERNES**
ECRIT, REALISÉ et PRODUIT par CHARLES CHAPLIN

MODERNE TIJDEN 

IMMAGINE BELGIQUE ETUDES PUBLICITAIRES IMP. J. LUCIERT & FILS - BRUXELLES

LES ARTISTES ASSOCIÉS S.A.B. *Présentent* PARLANT FRANÇAIS 105 ENFANTS ADMIS

**CHARLIE
CHAPLIN**

DANS

THANKS!

les **TEMPS MODERNES**
ECRIT, REALISÉ ET PRODUIT par CHARLES CHAPLIN

MODERNE TIJDEN 