

Generalized Strategic Classification and the Case of Aligned Incentives

Sagi Levanon and Nir Rosenfeld

Technion CS

@ICML 2022

Outline

- Strategic classification (SC)
- Generalized strategic classification (GSC)
- Incentive alignment
- Learning in GSC – not what you think!

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:

train:

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

user features



test:

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

*modified
features*



*in response
to learned model*



strategic

classification:

train:

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

user features
↓

test:

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

modified features ↑
in response to learned model ↑

want positive predictions (e.g., loan approved)



strategic

classification:

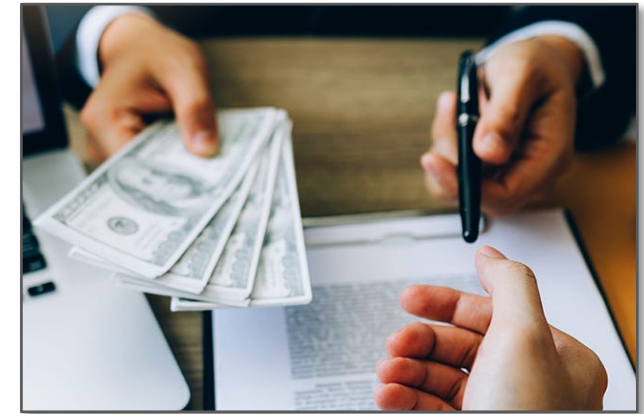
train: $\operatorname{argmin}_h \mathbb{E}[1\{y \neq h(x)\}]$

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

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

utility = prediction
(want $\hat{y} = 1$)

cost of modifying
(assumed known)



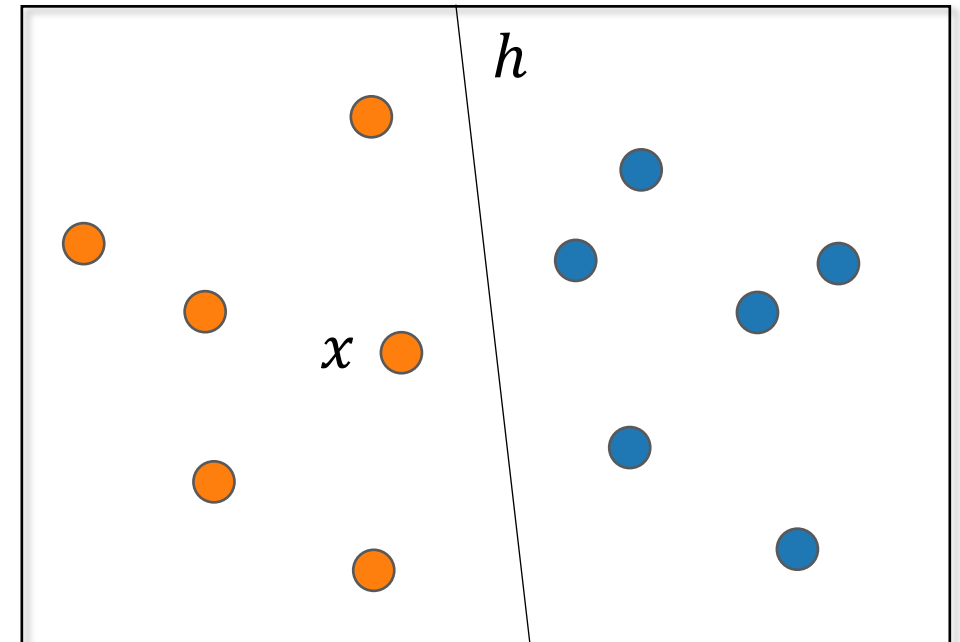
strategic

classification:

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

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

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



optimal non-strategic classifier

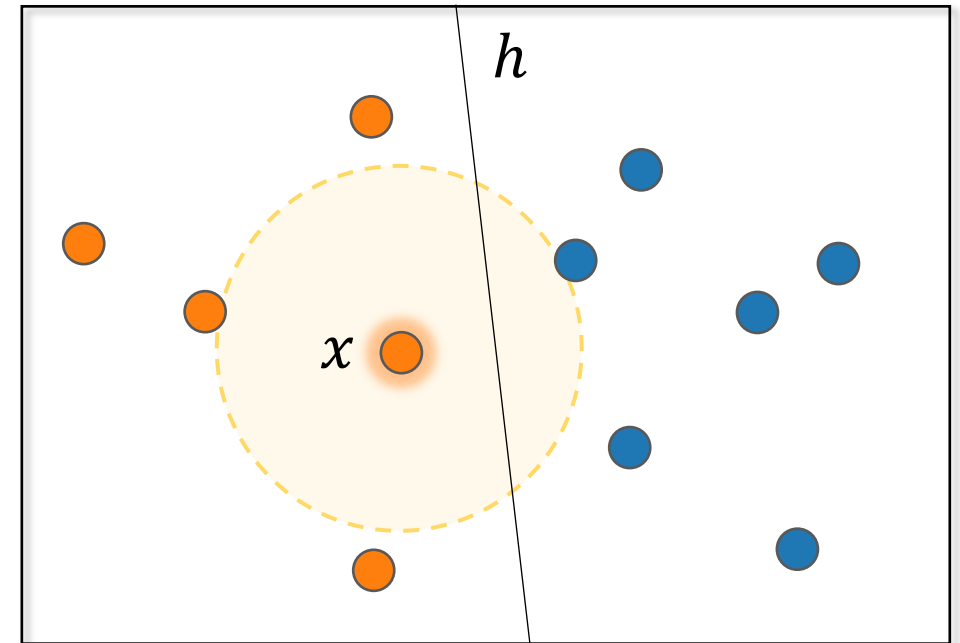
strategic

classification:

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

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

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



optimal non-strategic classifier

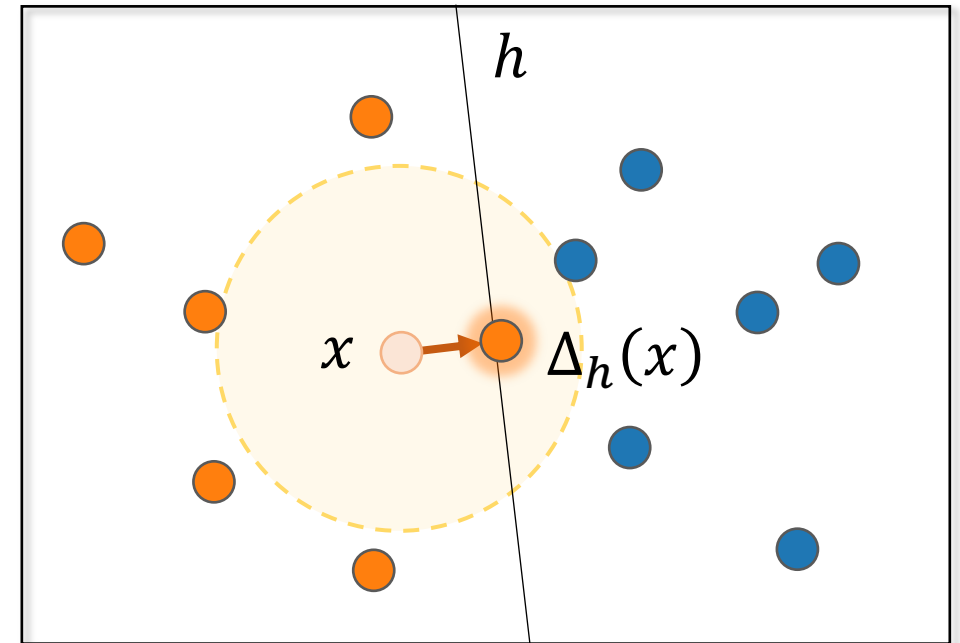
strategic

classification:

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

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

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



optimal non-strategic classifier

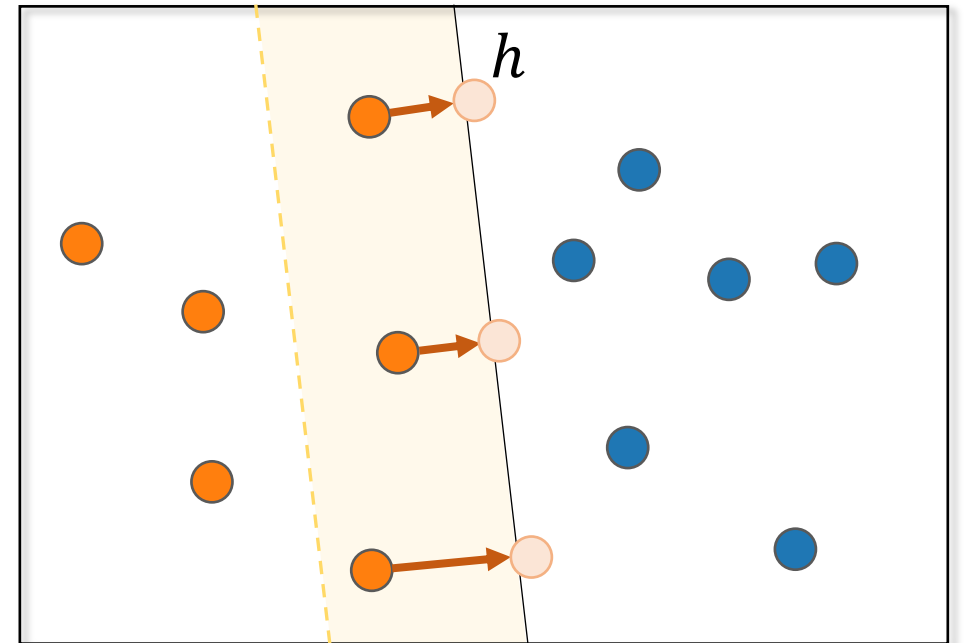
strategic

classification:

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

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

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



optimal non-strategic classifier
not optimal under strategic behavior

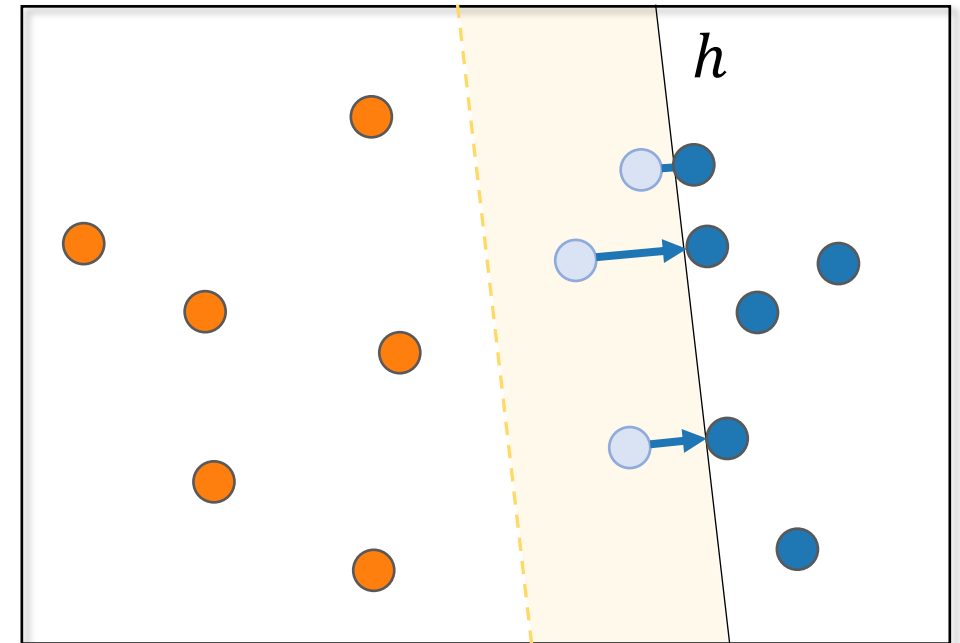
strategic

classification:

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

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

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



optimal classifier under strategic behavior

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

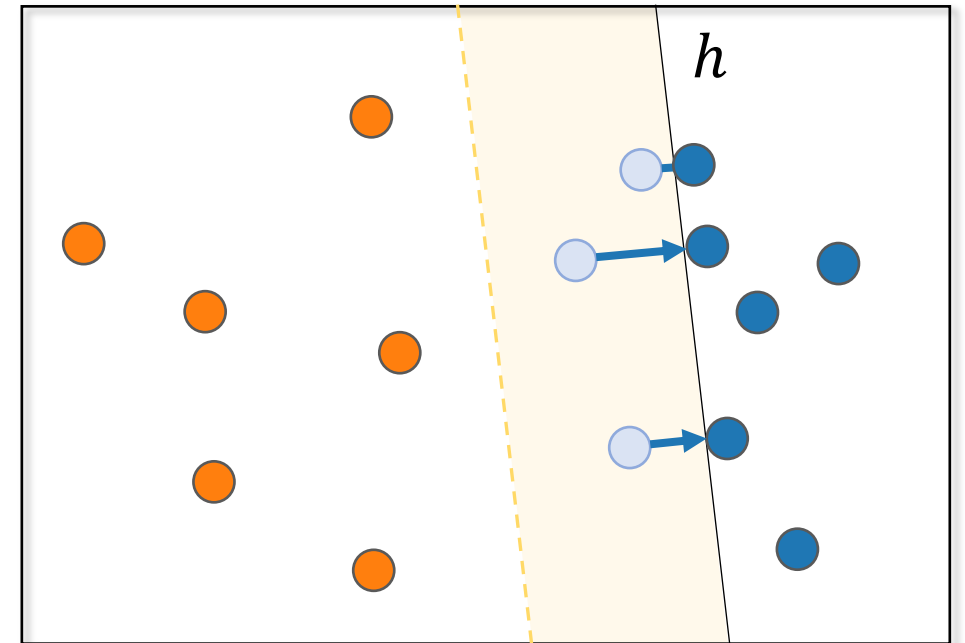
strategic

classification:

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

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

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



optimal classifier under strategic behavior

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

strategic

classification:

train:

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

natural solution:

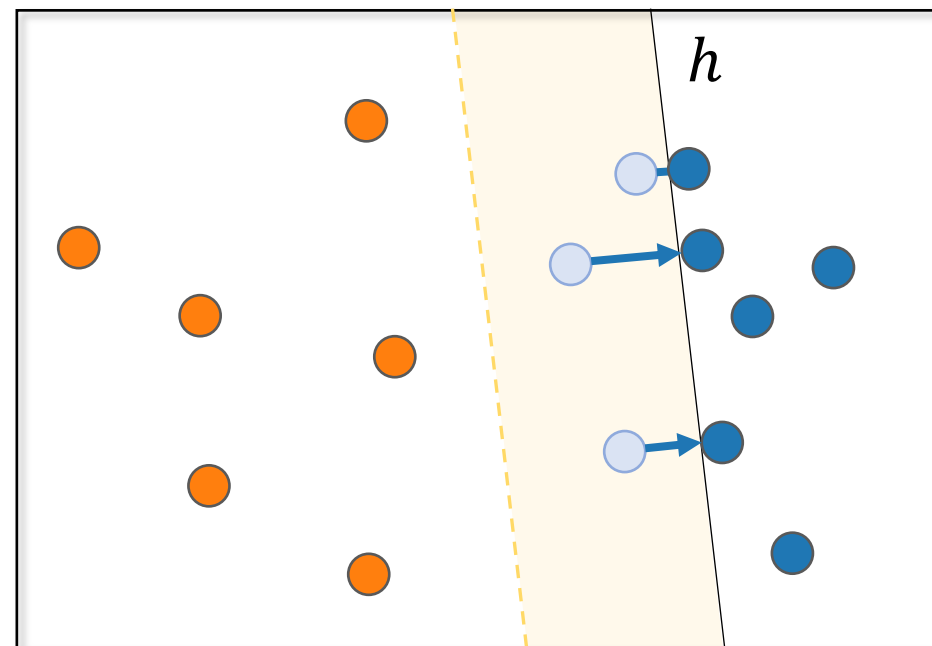
test:

$$h(\Delta_h(x)) = y$$

*main takeaway:
not so simple!*

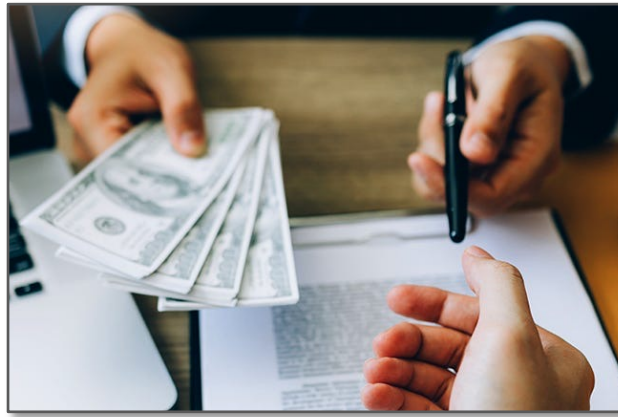
response:

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



optimal classifier under strategic behavior

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

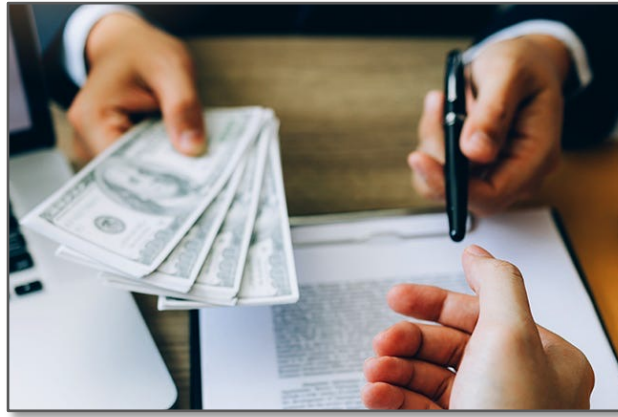


- **pro:** cleanly captures natural, prevalent tension
- **con:** narrow perspective, limited modeling power

?

?

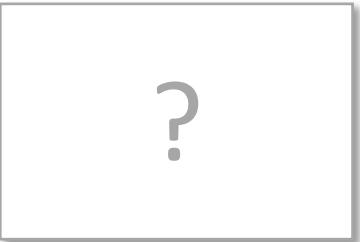
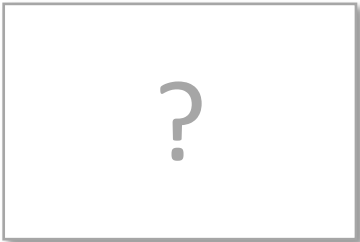
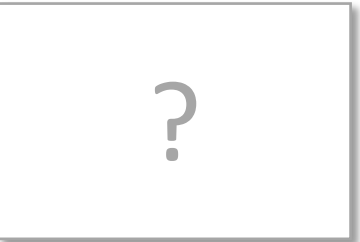
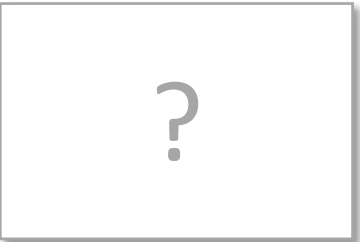
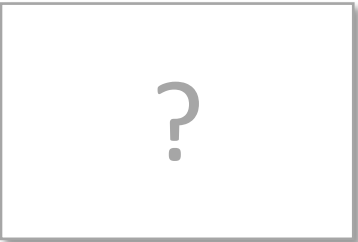
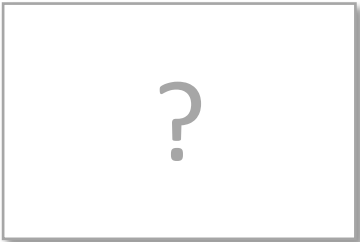
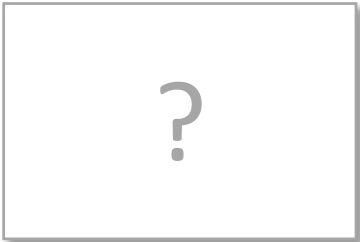
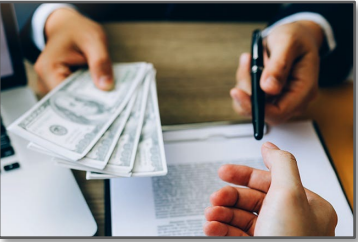
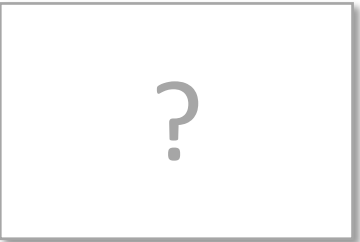
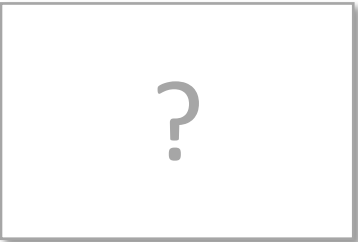
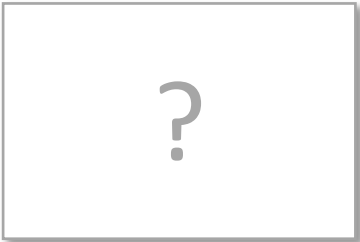
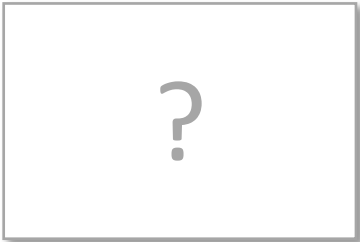
?



?

?

?



generalized
strategic
classification:

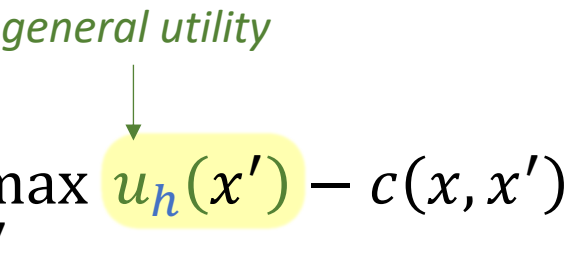
response: $\Delta_h(x) = \operatorname{argmax}_{x'} h(x') - c(x, x')$

users are rational: *maximize* *utility* - *cost*

**generalized
strategic
classification:**

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

general utility



Generalizing SC:

1. allow arbitrary utility functions u (that depend on h)

**generalized
strategic
classification:**

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

perceived utility ↓

$\neq u_h(x') \leftarrow \text{actual utility}$

Generalizing SC:

1. allow arbitrary utility functions u (that depend on h)
2. let users act on perceived utility \tilde{u} (can differ from *true* utility)

**generalized
strategic
classification:**

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

perceived utility (green arrow pointing to $\tilde{u}_h(x', z)$)
private information (orange arrow pointing to z)
 $\neq u_h(x') \leftarrow$ *actual utility* (green arrow pointing to $u_h(x')$)

Generalizing SC:

1. allow arbitrary utility functions u (that depend on h)
2. let users act on perceived utility \tilde{u} (can differ from *true* utility)
3. permit users to hold private information z (on which \tilde{u} relies)

generalizes:
*Hardt et al. (2016),
Sundaram et al. (2021),
Jagadeesan et al. (2021),
Ghalme et al. (2021), ...*

main questions: can we learn? when? how?

strategic = “gaming”

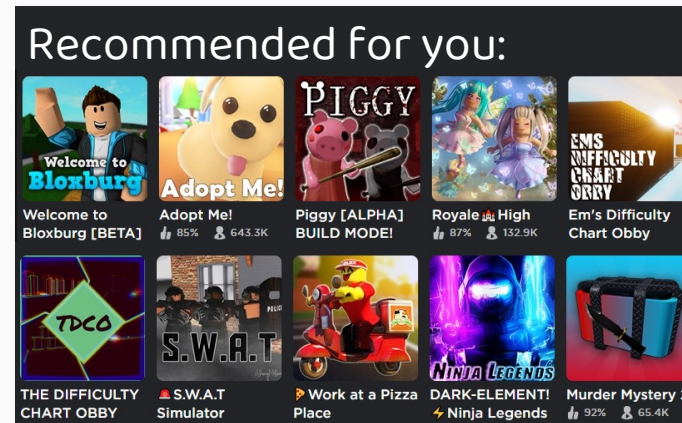


classification *about* humans

system wants: *correct* predictions

users want: *positive* predictions

incentives **align**



classification *for* humans

system wants: *correct* predictions

users want: *correct* predictions

- **Incentive Alignment (IA):** [definition]

$$\exists h \text{ s.t. } \mathbb{E}[\mathbb{1}\{y \neq h(\Delta_h(x, z))\}] < \min_{h'} \mathbb{E}[\mathbb{1}\{y \neq h'(x)\}]$$

- **Begin simple:** users act on **noisy label** beliefs:

$$z = \tilde{y} = y \text{ w.p. } 1 - \epsilon$$

$$\Rightarrow \Delta_h(x, \tilde{y}) = \operatorname{argmax}_{x'} \underbrace{\mathbb{1}\{h(x') = \tilde{y}\}}_{\tilde{u}} - c(x, x')$$

Expect learning to improve if:

- private (noisy) labels informative of true labels
- it can utilize users' willingness to invest effort

- **Incentive Alignment (IA):** [definition]

$$\exists h \text{ s.t. } \mathbb{E}[\mathbb{1}\{y \neq h(\Delta_h(x, z))\}] < \min_{h'} \mathbb{E}[\mathbb{1}\{y \neq h'(x)\}]$$

- **Begin simple:** users act on **noisy label** beliefs:

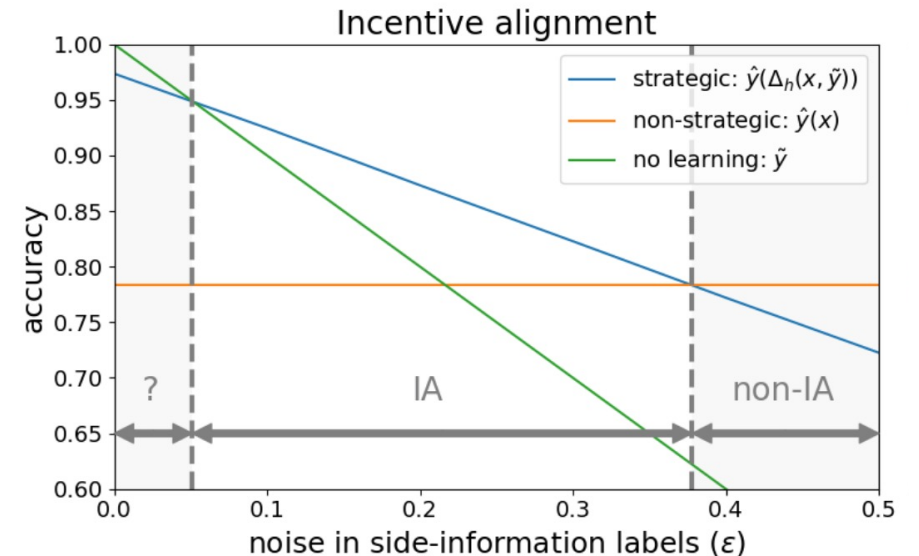
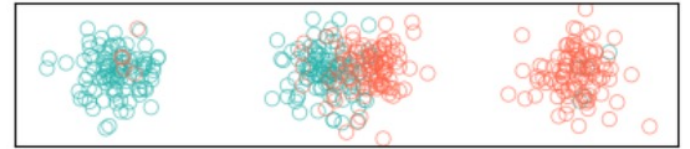
$$z = \tilde{y} = y \text{ w.p. } 1 - \epsilon$$

$$\Rightarrow \Delta_h(x, \tilde{y}) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$$

Expect learning to improve if:

- private (noisy) labels informative of true labels
- it can utilize users' willingness to invest effort
- **Empirically** – learning aligns incentives for many ϵ :
- **Theoretically** – full characterization (see paper)

Data:



strategic behavior can be helpful!

learning:

problem: train-test discrepancy

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

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

response: $h(\Delta_h(x, \tilde{y})) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$

learning:

natural solution: anticipate user response

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

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

response: $h(\Delta_h(x, \tilde{y})) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$

learning:

0/1 loss:

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

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

response: $h(\Delta_h(x, \tilde{y})) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$

learning:

hinge loss:

train: $\operatorname{argmin}_h \mathbb{E}[\max\{0, 1 - yw^\top \Delta_h(x, \tilde{y})\}]$?

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

response: $h(\Delta_h(x, \tilde{y})) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$

learning:

hinge loss:

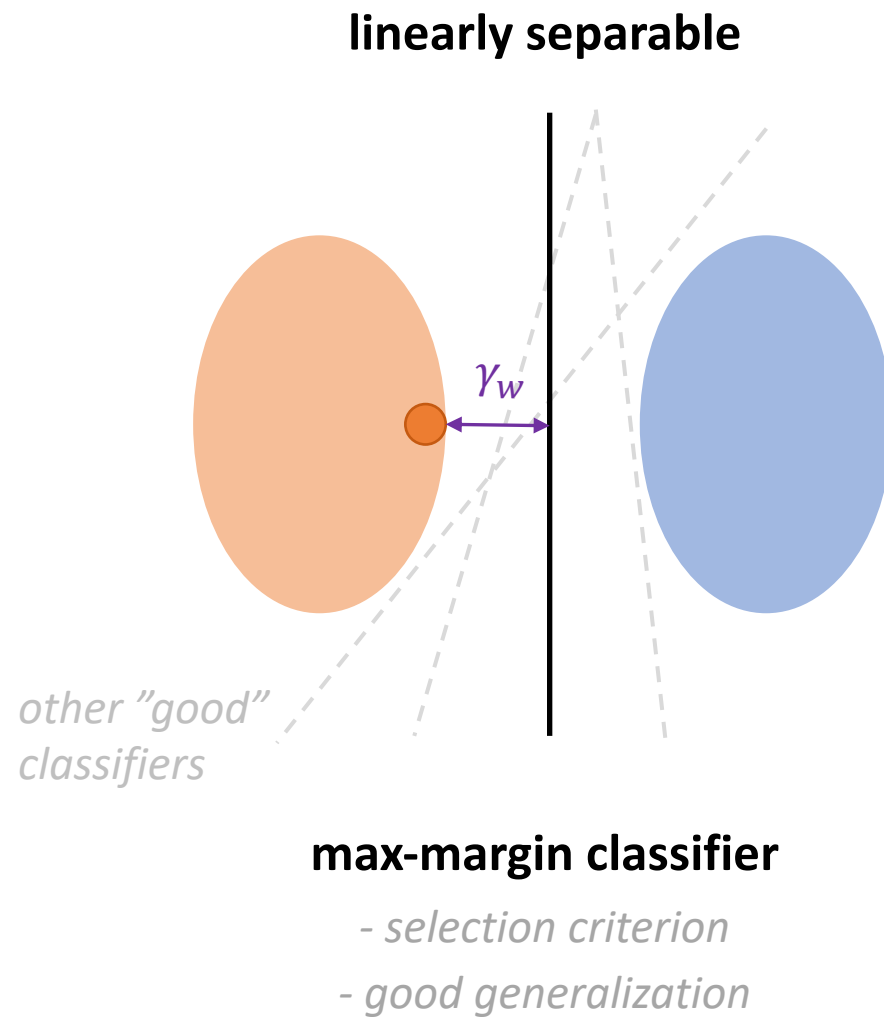
train: $\operatorname{argmin}_h \mathbb{E}[\max\{0, 1 - yw^\top \Delta_h(x, \tilde{y})\}]$ **×**

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

response: $h(\Delta_h(x, \tilde{y})) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$

standard hinge loss:

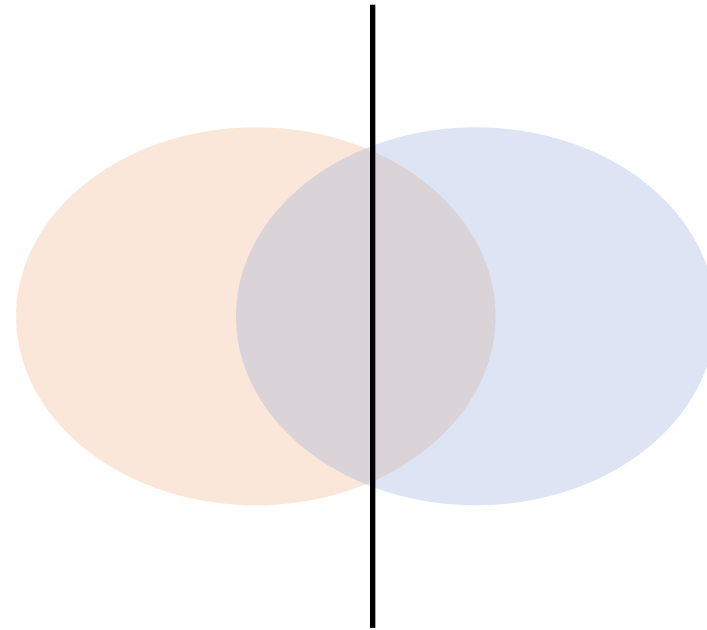
$$\max\{0, 1 - yw^T x\}$$



standard hinge loss:

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

not linearly separable



max-margin classifier

- *selection criterion*
- *good generalization*

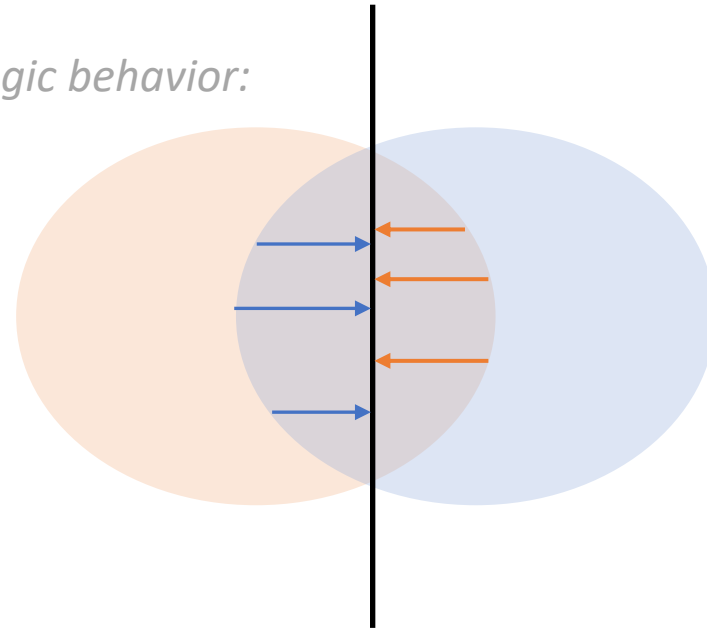
naïve hinge loss:

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

easiest to imagine!

not linearly separable

strategic behavior:



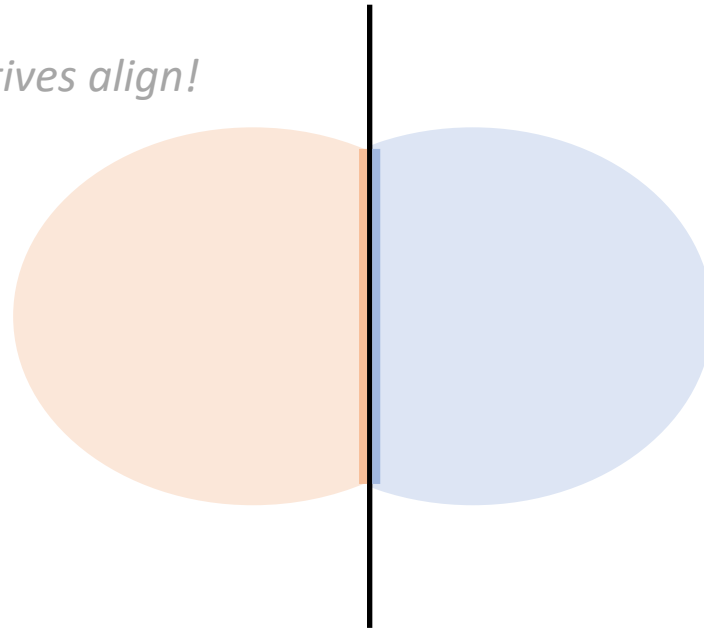
(naïve) max-margin classifier

naïve hinge loss:

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

strategically linearly separable

incentives align!



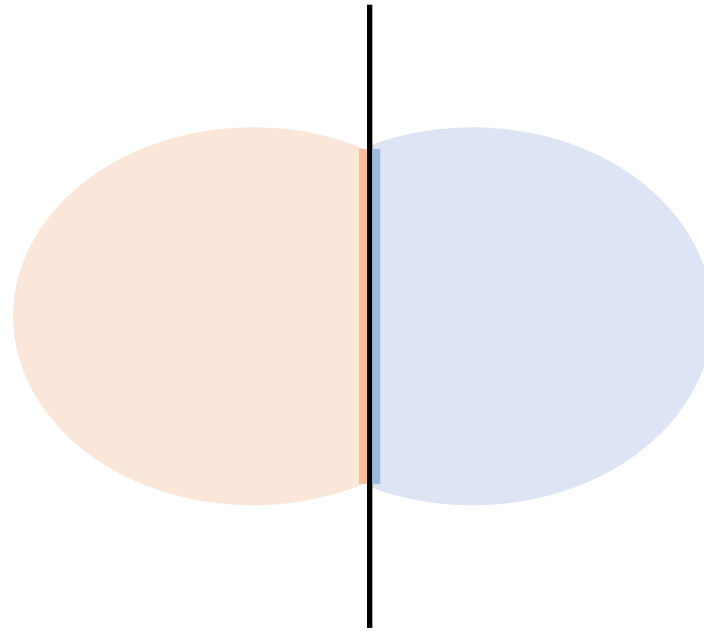
(naïve) max-margin classifier

perfect accuracy!

naïve hinge loss:

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

strategically linearly separable



(naïve) max-margin = 0

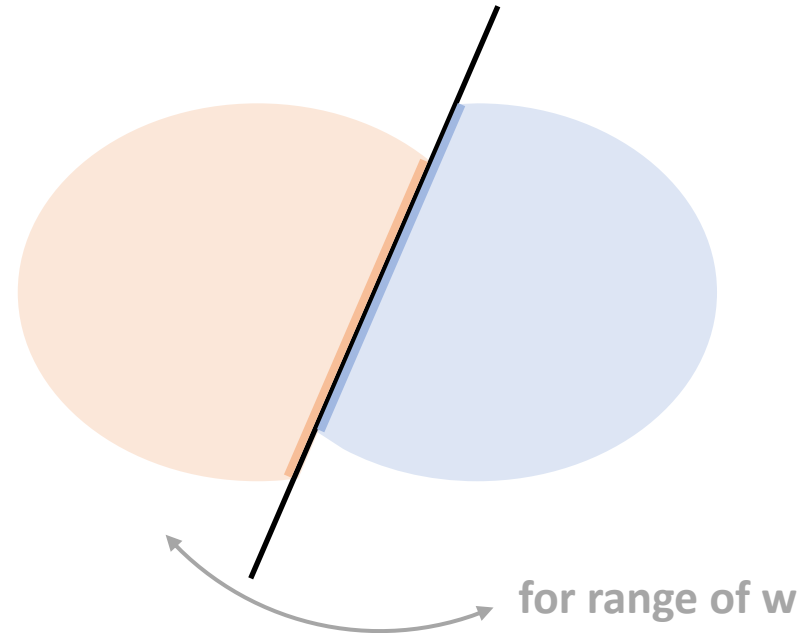
perfect accuracy?

naïve hinge loss:

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

conclusion: adapting hinge to strategic settings requires **rethinking the concept of margin**

strategically linearly separable



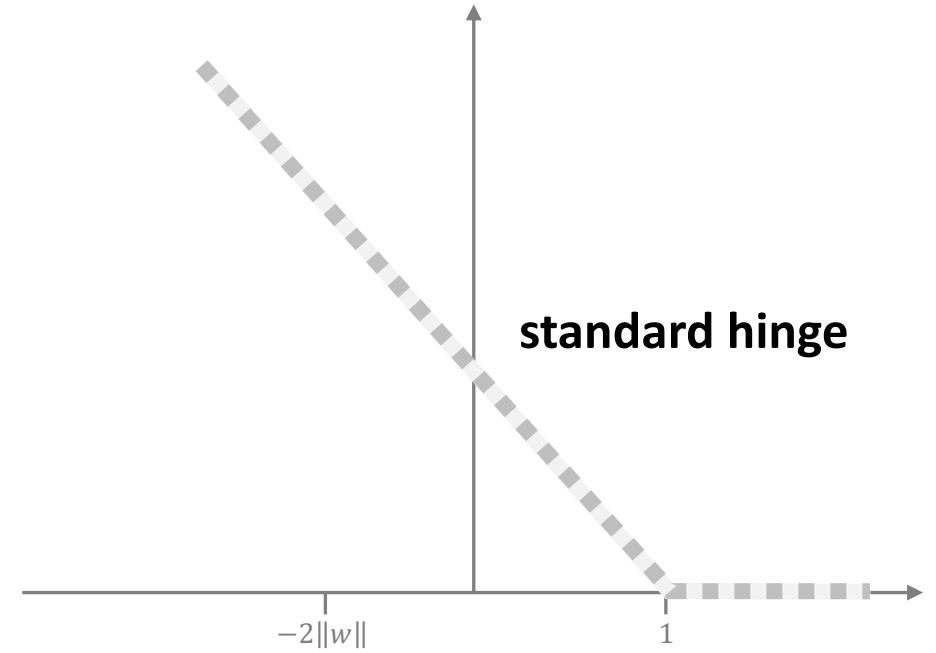
(naïve) max-margin = 0

⇒ vocous selection criterion

⇒ won't help generalization!

- **standard hinge:**

$$\max\{0, 1 - yw^\top x\}$$



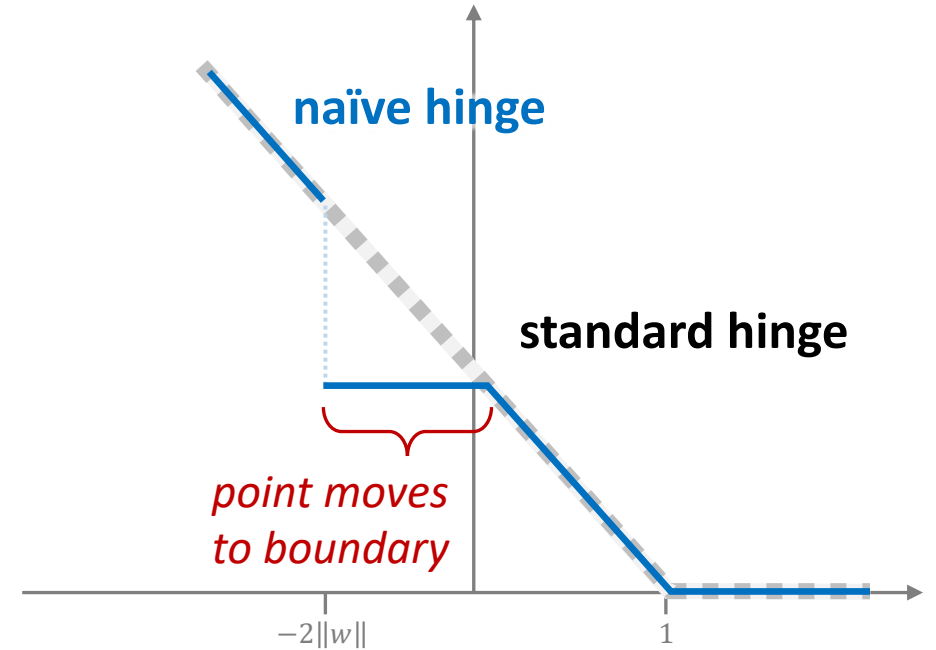
- **standard hinge:**

$$\max\{0, 1 - yw^\top x\}$$

- **naïve hinge:**

$$\max\{0, 1 - yw^\top \Delta_h(x, \tilde{y})\}$$

suspicious: does not depend on x !



- **standard hinge:**

$$\max\{0, 1 - yw^\top x\}$$

- **naïve hinge:**

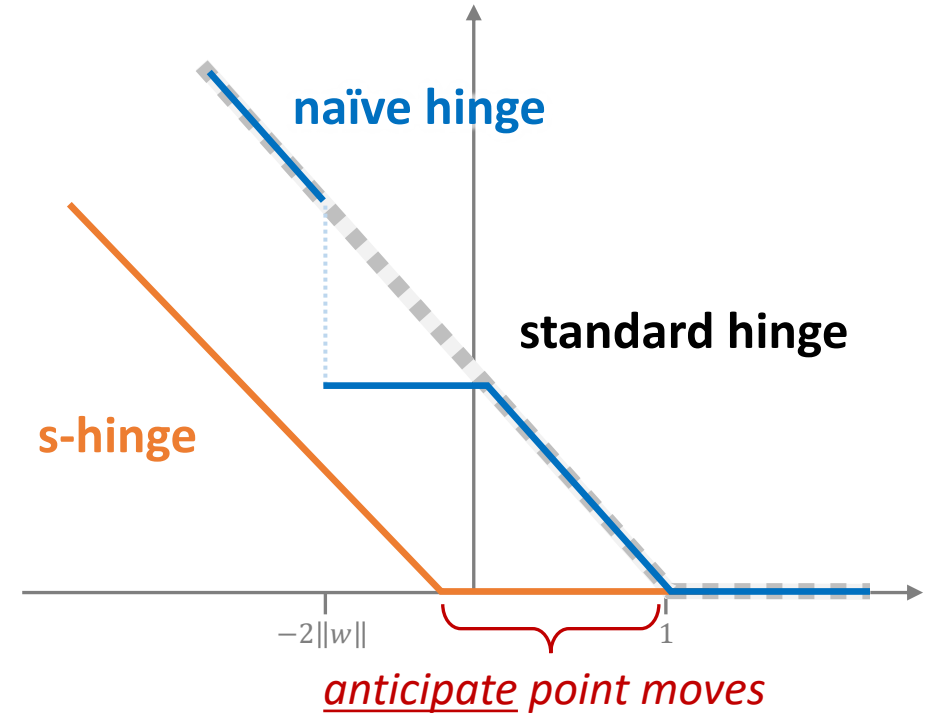
$$\max\{0, 1 - yw^\top \Delta_h(x, \tilde{y})\}$$

- **our strategic hinge:** (s-hinge)

$$\max\{0, 1 - yw^\top x - 2y\tilde{y}\|w\|\}$$

\uparrow
prediction before movement...

\uparrow
*... but relative to "shifted",
 personalized decision boundary*



- **standard hinge:**

$$\max\{0, 1 - yw^\top x\}$$

- **naïve hinge:**

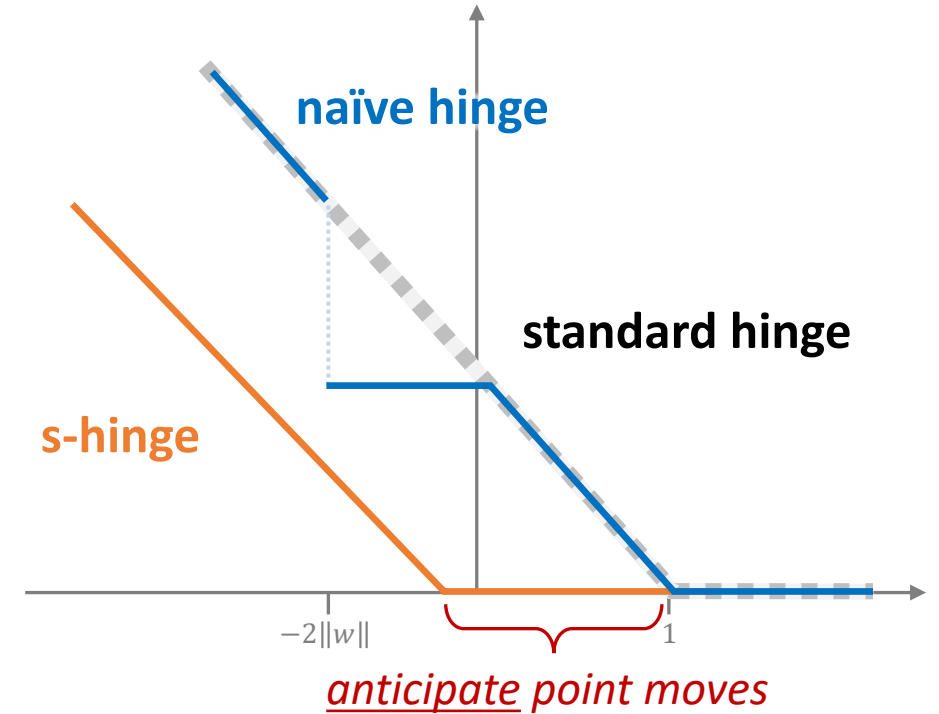
$$\max\{0, 1 - yw^\top \Delta_h(x, \tilde{y})\}$$

- **our strategic hinge:** (s-hinge)

$$\max\{0, 1 - yw^\top x - 2y\tilde{y}\|w\|\}$$

- **Benefits:**

1. Optimization – differentiable; no nasty Δ !
2. Generalization guarantees (up next)
3. Extends to broader strategic settings



general-preference strategic classification:

noisy labels:

$$\Delta_h(x; \tilde{y}) = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = \tilde{y}\} - c(x, x')$$

standard SC: $\mathbb{1}\{h(x') = 1\}$

"adversarial": $\mathbb{1}\{h(x') \neq y\}$

\vdots

general preferences: $\mathbb{1}\{h(x') = y'\}, y' \in \{\pm 1\}$

also Sundaram et al. (2021)

general-preference strategic classification:

general preferences:

$$\Delta_h(x; y') = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = y'\} - c(x, x')$$

1) s-hinge extends to GP:

$$\max\{0, 1 - yw^\top x - 2yy'\|w\|\}$$

differentiable – no Δ !

2) generalization bound:

$$\mathcal{L}_{0/1} \leq \mathcal{L}_{0/1}^{\text{NL}} + \frac{4r\|\hat{w}\|}{\sqrt{m}} + (1 + 2\rho^{\text{GP}}\|w\|) \sqrt{\frac{2 \ln(4\|\hat{w}\|/\delta)}{m}}$$

- extends and closely matches non-strategic case
- relations: $\rho_{\epsilon < 0.5}^{\text{NL}} \leq \rho \leq \rho_{\epsilon > 0.5}^{\text{NL}} \leq \rho^{\text{SC}} = \rho^{\text{ADV}} = \rho^{\text{GP}}$

general-preference strategic classification:

general preferences:

$$\Delta_h(x; y') = \operatorname{argmax}_{x'} \mathbb{1}\{h(x') = y'\} - c(x, x')$$

1) s-hinge extends to GP:

$$\max\{0, 1 - yw^\top x - 2yy'\|w\|\}$$

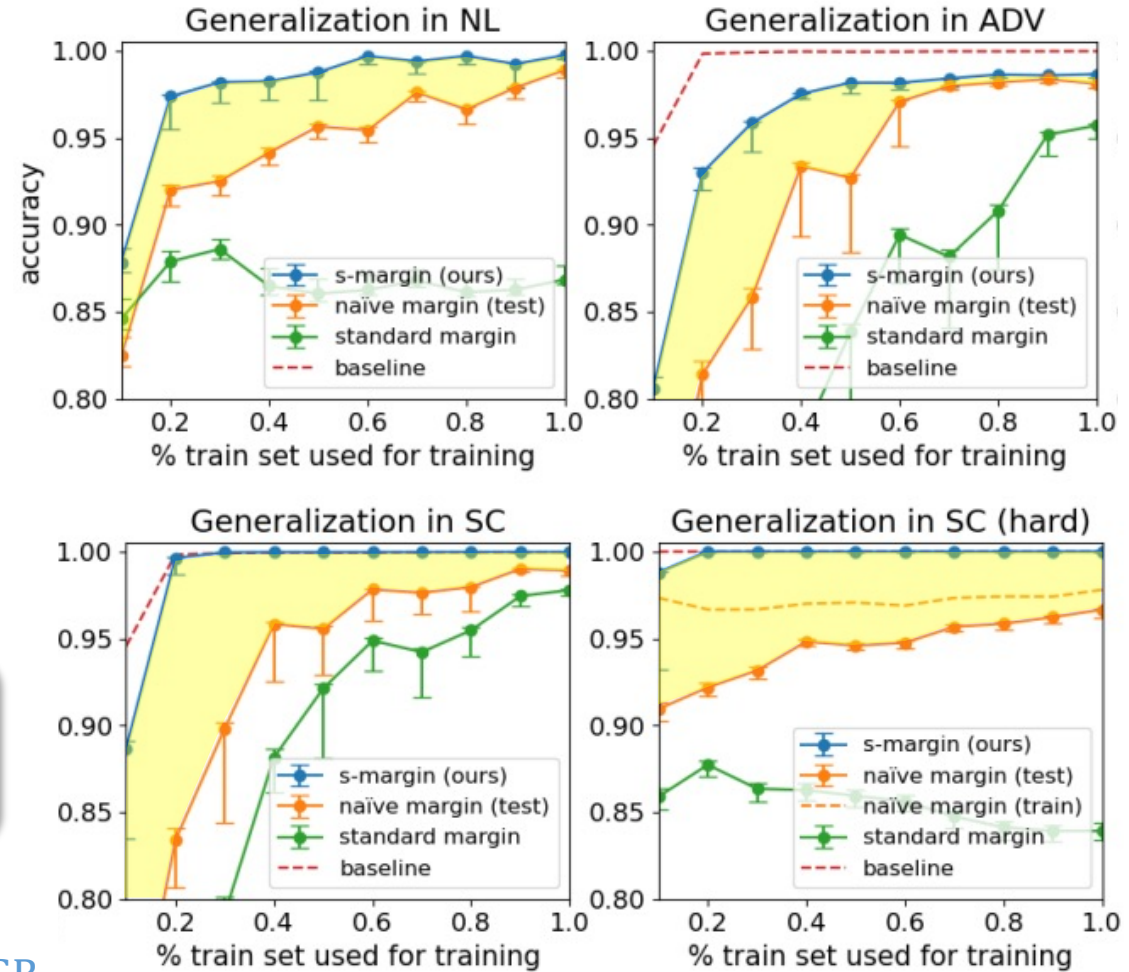
differentiable – no Δ !

2) generalization bound:

$$\mathcal{L}_{0/1} \leq \mathcal{L}_{0/1}^{\text{NL}} + \frac{4r\|\hat{w}\|}{\sqrt{m}} + (1 + 2\rho^{\text{GP}}\|w\|) \sqrt{\frac{2 \ln(4\|\hat{w}\|/\delta)}{m}}$$

- extends and closely matches non-strategic case
- relations: $\rho_{\epsilon < 0.5}^{\text{NL}} \leq \rho \leq \rho_{\epsilon > 0.5}^{\text{NL}} \leq \rho^{\text{SC}} = \rho^{\text{ADV}} = \rho^{\text{GP}}$

3) empirical generalization gaps:



generalized strategic classification:

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

perceived utility
side information

- **standard hinge:**

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

correctness distance

generalized strategic classification:

$$\Delta_h(x; z) = \operatorname{argmax}_{x'} \underbrace{\tilde{u}_h(x', z)}_{\text{perceived utility}} - \underbrace{c(x, x')}_{\text{side information}}$$

- **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:**

$$\max\{0, 1 - \underbrace{\operatorname{sign}(yw^\top \Delta_h(x, z))}_{\text{correctness}} \underbrace{|w^\top \Delta_h(x, z)|}_{\text{distance}}\} \quad \begin{aligned} & - \text{no original } x \\ & - \text{distance after move} \\ & - \text{distance to hyperplane} \end{aligned}$$

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

$$\max\{0, 1 - \underbrace{\operatorname{sign}(yw^\top \Delta_h(x, z))}_{\text{correctness}} \underbrace{d_\Delta(x, z; w)}_{\text{anticipated distance}}\}$$

generalized strategic classification:

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

- standard hinge:

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

- naïve hinge:

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

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

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

reinterpretation of “margin”:

distance to nearest x' who's movement *flips label*:
subsumes non-strategic case

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

flip label s.t. $h(\Delta_h(x, z)) \neq h(\Delta_h(x', z))$
(after movement)

- *subsumes standard hinge for non-strat.*
- *for GP, get s-hinge as special case*
- *for general GSC, may be hard to compute*

extended generalization bound:

$$\mathcal{L}_{0/1} \leq \mathcal{L}_{0/1}^{\text{NL}} + \frac{8r\|\hat{w}\|}{\sqrt{m}} + (1 + 2\rho^{\text{GSC}}\|w\|) \sqrt{\frac{2 \ln(4\|\hat{w}\|/\delta)}{m}}$$

strategic content recommendation:

perceived utility

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

personal history

$$= \{(a_j, y_j)\}_{j=1}^n$$

modify features to have previous experiences classified correctly

$$\frac{1}{n} \sum_{j=1}^n \mathbb{1}\{h(x', a_j) = y_j\}$$

want correct predictions on future items:

like/dislike

$$\leftarrow u_h(x) = \mathbb{1}\{h(x, a) = y\}$$

true utility *items*
(e.g., movies)

act in hope of improving (future) accuracy:

$$\longrightarrow u_h(\Delta_h(x, z)) = \mathbb{1}\{h(\Delta_h(x, z), a) = y\}$$

- **Challenge:** in essence, users *also* aim to minimize 0-1 loss
- Get *coupling* of learning problems (system+users) – **implicit cooperation!**

strategic content recommendation:

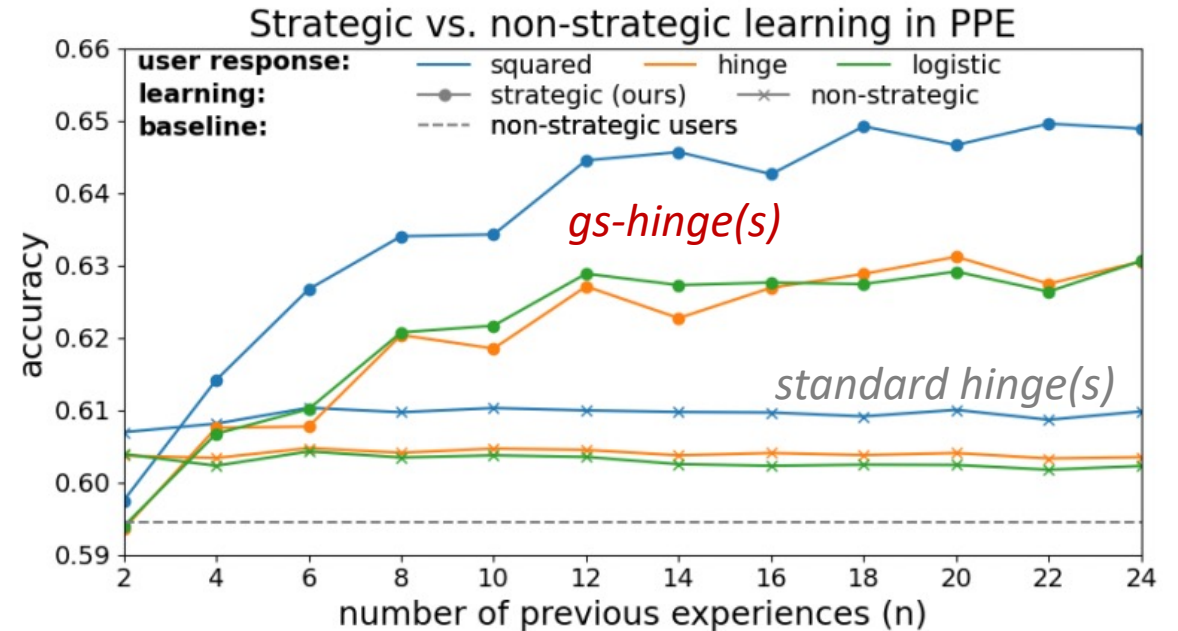
$$\Delta_h(x; z) = \operatorname{argmax}_{x'} \tilde{u}_h(x', z) - c(x, x')$$

personal history

$$= \{(a_j, y_j)\}_{j=1}^n$$

modify features to have previous experiences classified correctly

$$\frac{1}{n} \sum_{j=1}^n \mathbb{1}\{h(x', a_j) = y_j\}$$



- **Challenge:** in essence, users *also* aim to minimize 0-1 loss
- Get *coupling* of learning problems (system+users) – **implicit cooperation!**
- We propose differentiable proxy
- **Results:** highly incentive-aligned task!

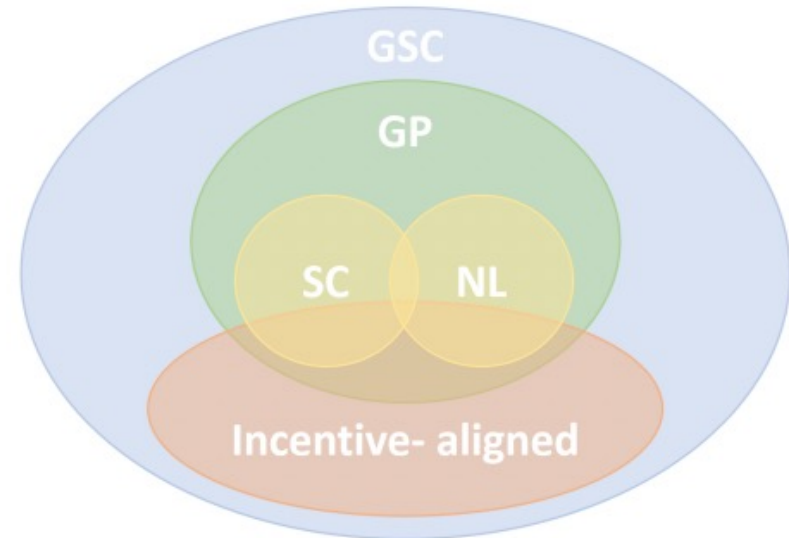
Conclusions

1. Strategic behavior **comes in many flavors**

- Hierarchy of problem classes
- Mild variations in user response
⇒ very different learning problems
- We highlight **incentive alignment**
- Others interesting classes?
(e.g., bounded-rational, behavioral, Bayesian)

2. Strategic behavior **can break proxy losses**

- Choose with care!
- We propose **gs-hinge**:
differentiable (sometimes) + theoretical guarantees + empirical results
- Other losses? (e.g., strategic cross entropy)



Thanks!

(come to our poster)

check out our other paper on
strategic representation

(also @ICML2022)

