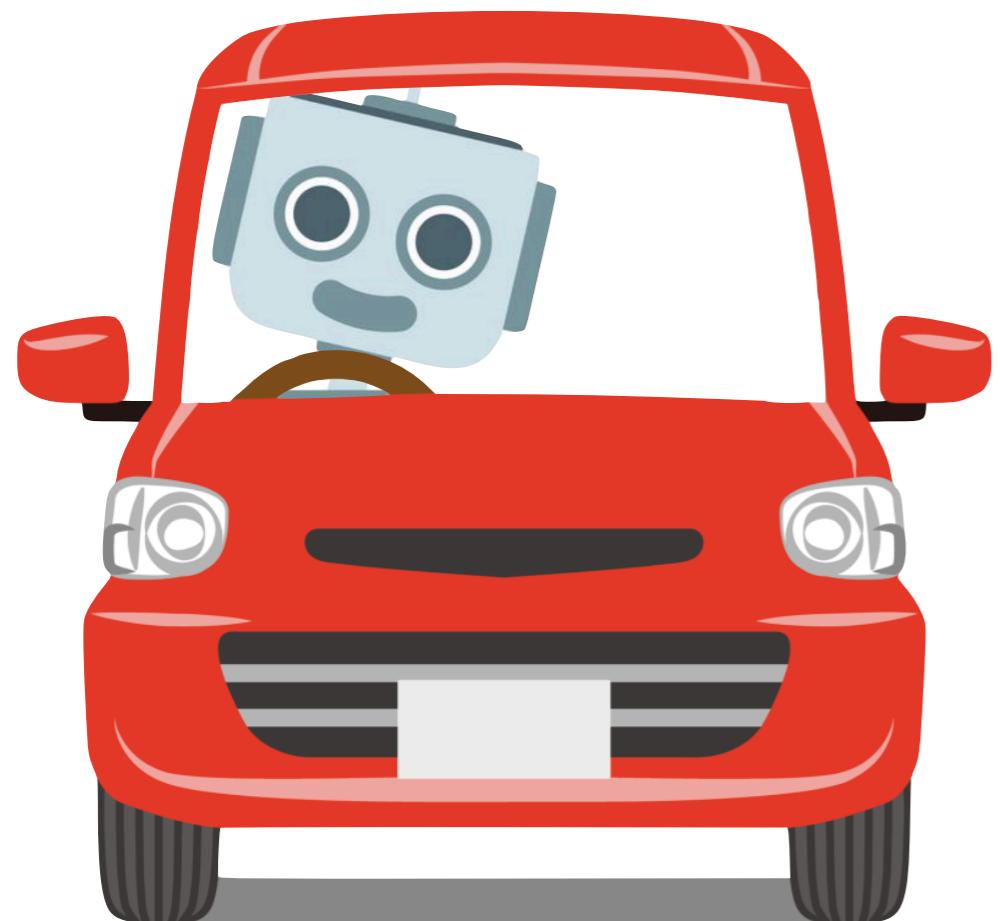


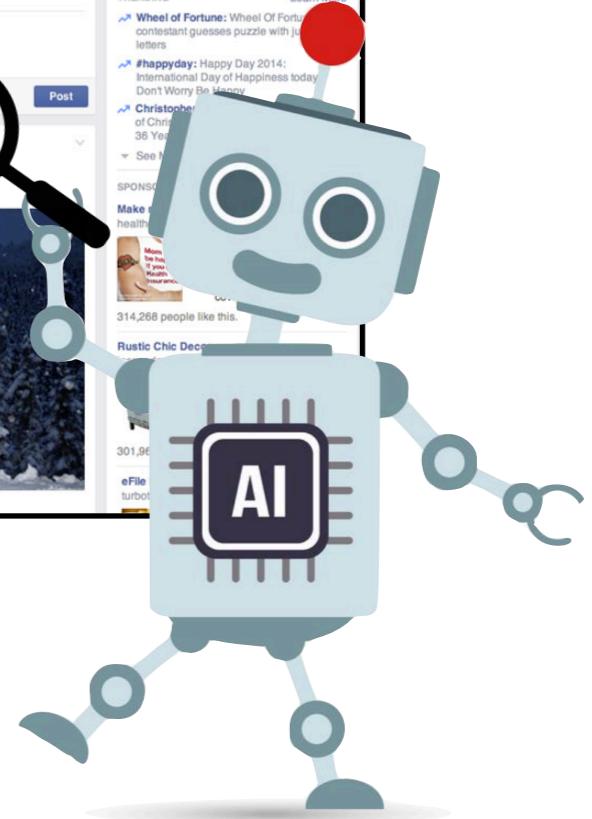
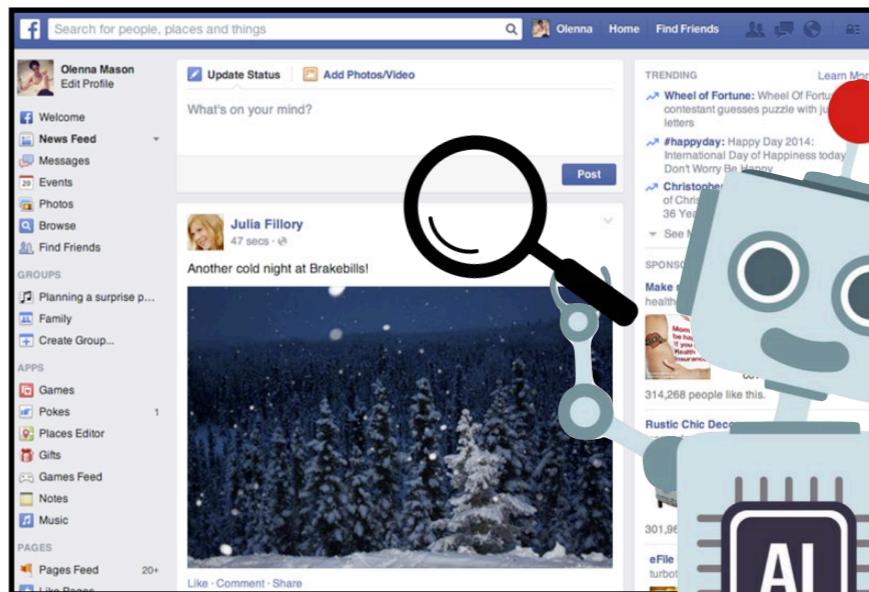
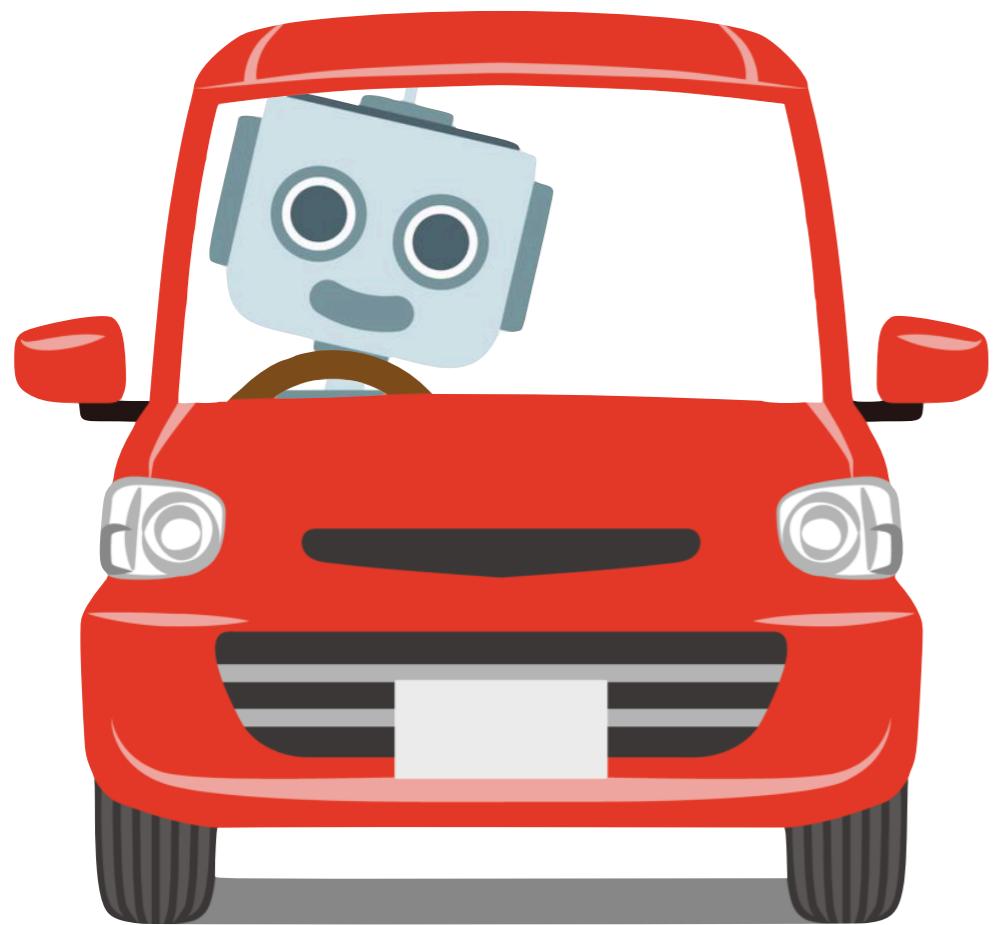
Learning to Reject and Learning to Defer to an Expert

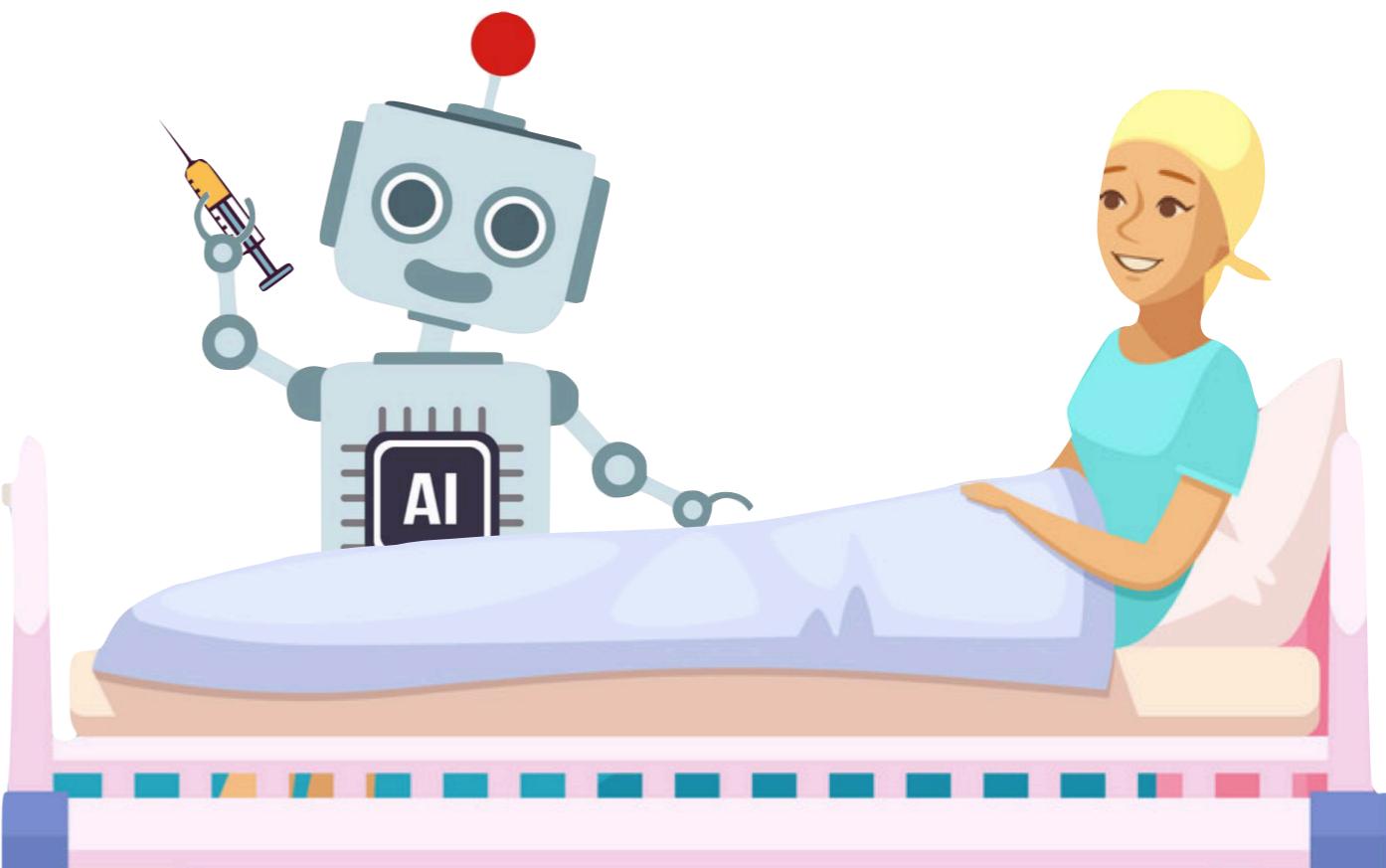
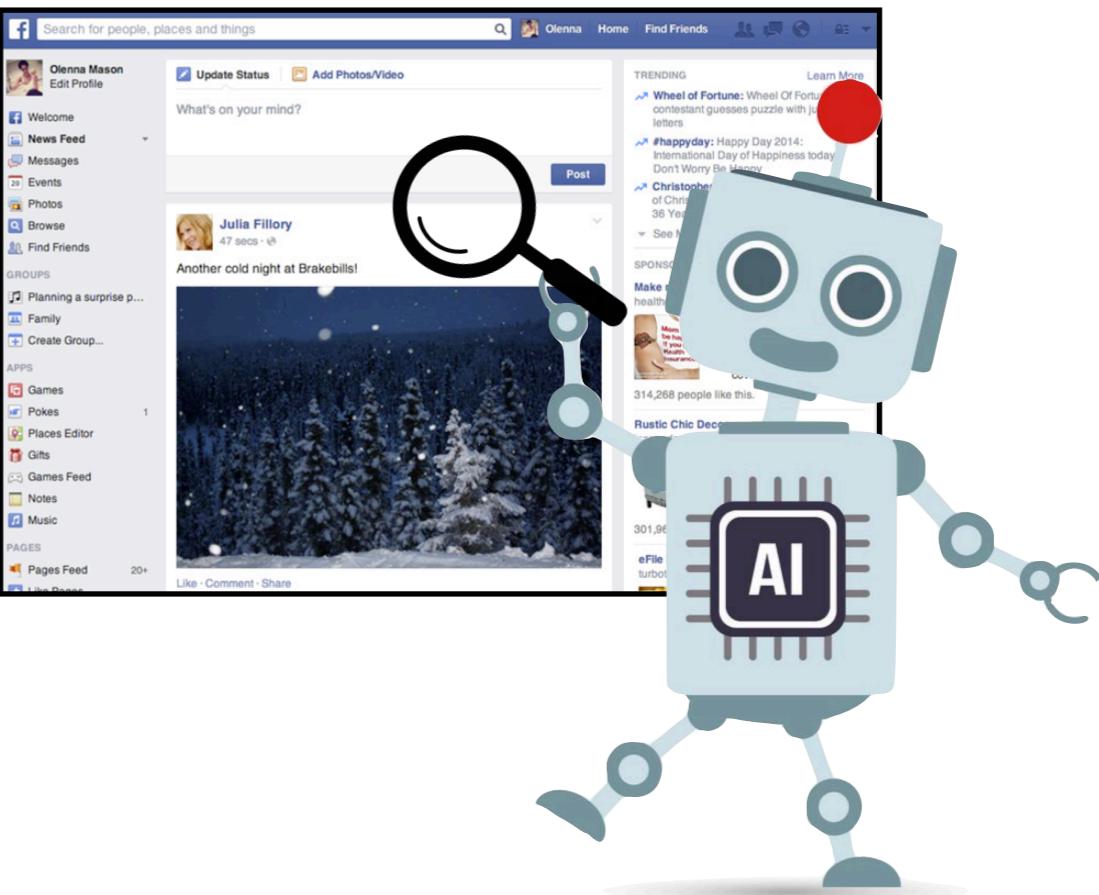
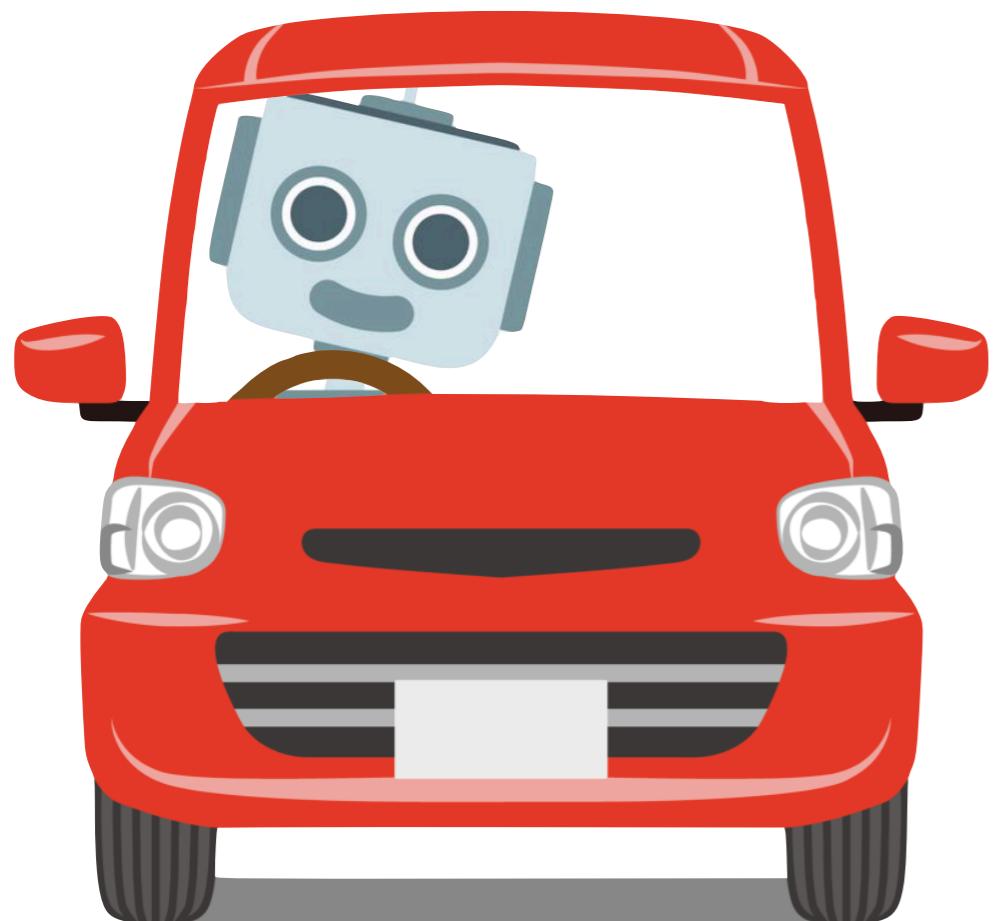
Eric Nalisnick

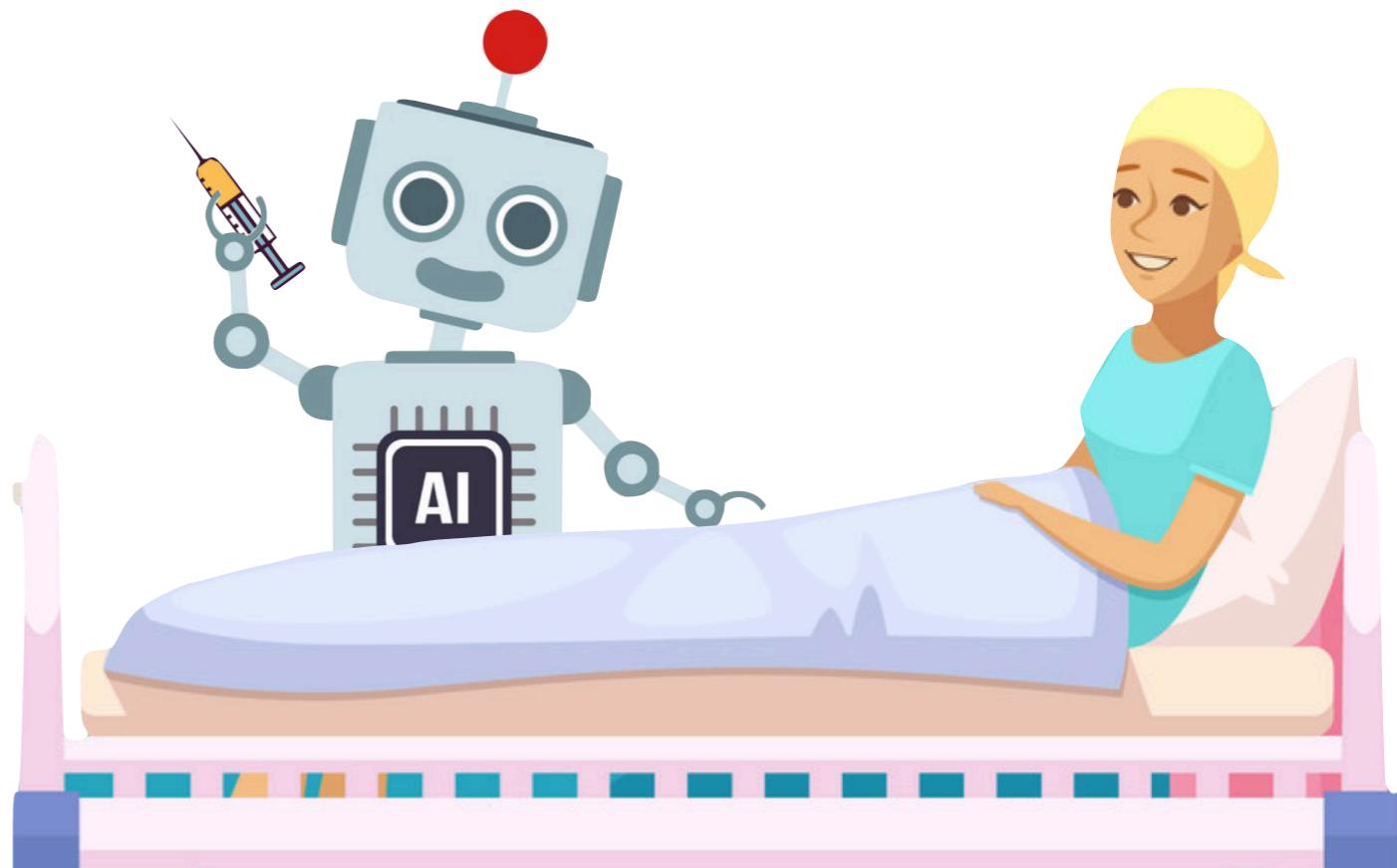
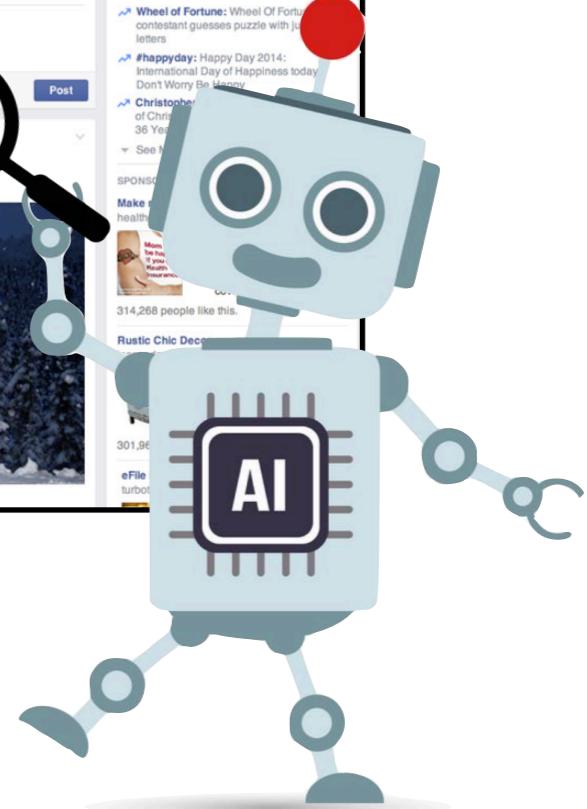
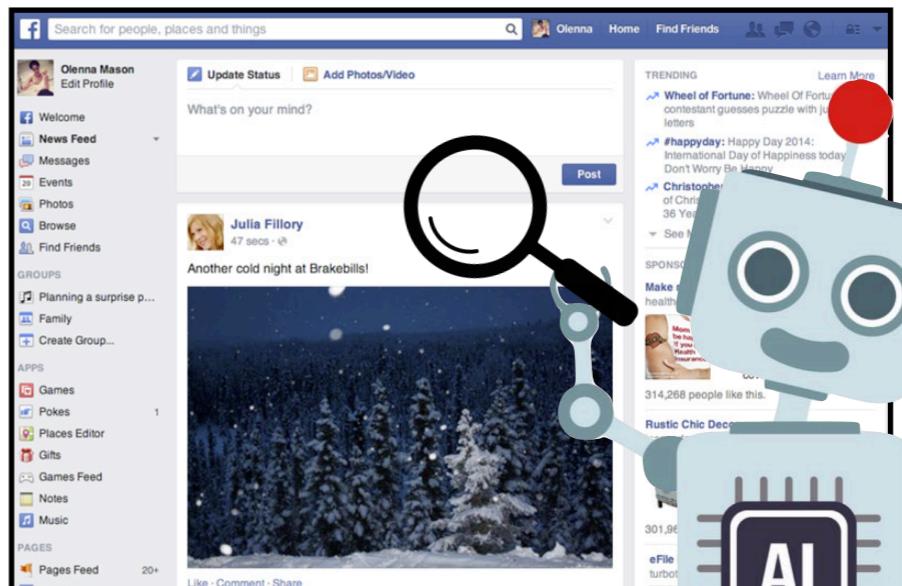
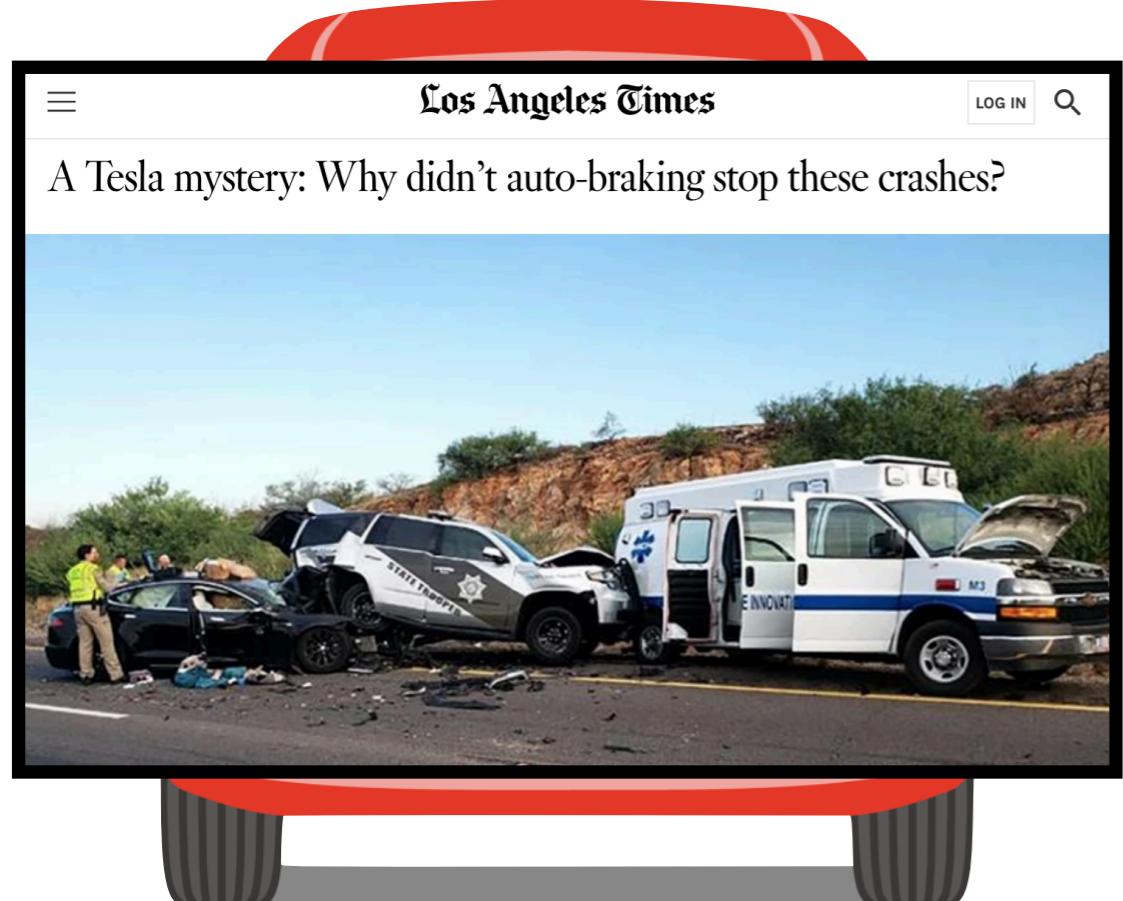
Johns Hopkins University

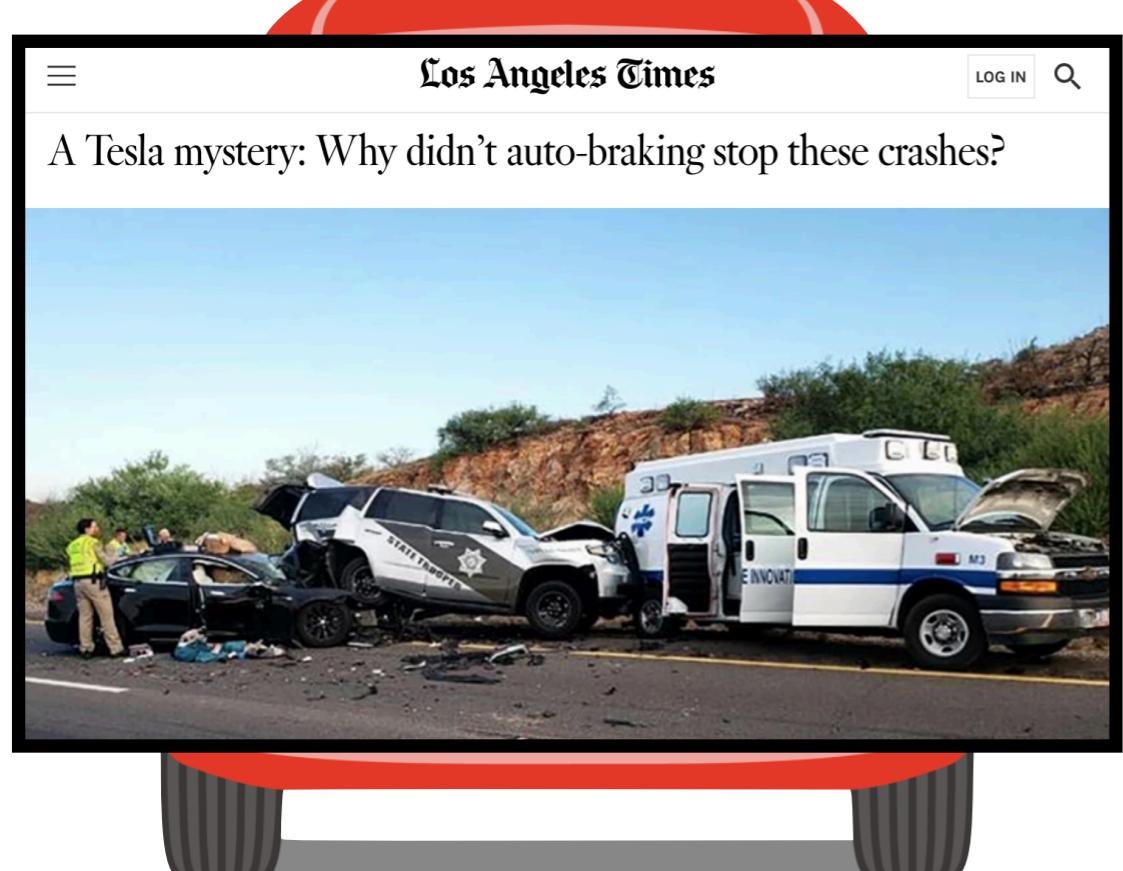












from Net Politics and Digital and Cyberspace Policy Program

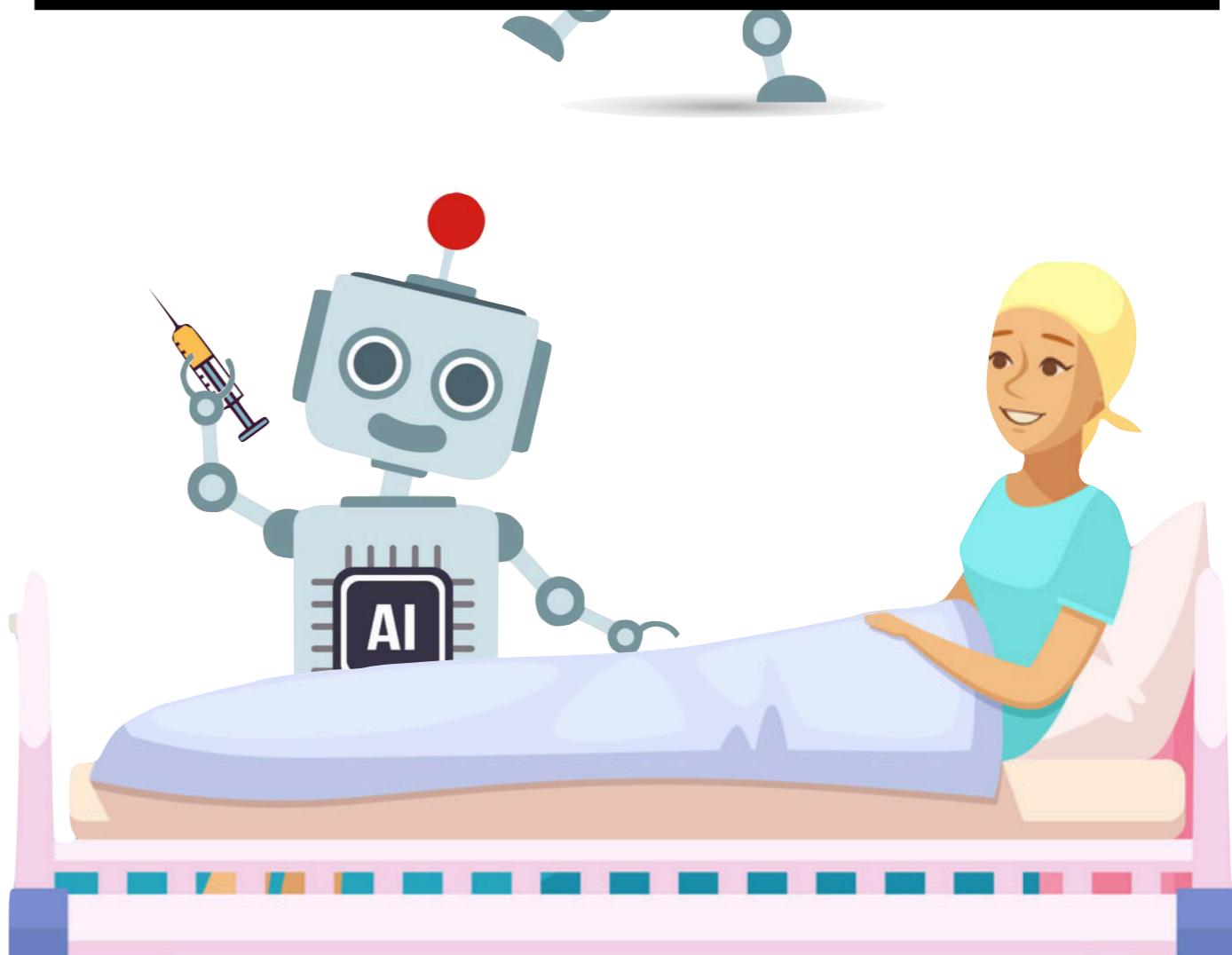
Facebook's Content Moderation Failures in Ethiopia

Facebook has failed to moderate content in underserved countries. Facebook and other social media companies must invest more in local content moderation, instead of relying on global AI systems.

Blog Post by Caroline Allen, Guest Contributor
April 19, 2022 2:36 pm (EST)

[Facebook](#) [Twitter](#) [LinkedIn](#) [Print](#) [Email](#)

Amhara militia members ride in the back of a truck towards a fight with the Tigray People's Liberation Front. Reuters/Tiksa Negeri



Los Angeles Times

A Tesla mystery: Why didn't auto-braking stop these crashes?

The image shows a scene of multiple vehicle collisions on a highway. In the foreground, a dark-colored Tesla sedan has crashed into the side of a white and blue emergency response vehicle, which appears to be a California State Trooper's car. Other vehicles are visible in the background, and several people are standing near the accident site.

from Net Politics and Digital and Cyberspace Policy Program

Facebook's Content Moderation Failures in Ethiopia

Facebook has failed to moderate content in underserved countries. Facebook and other social media companies must invest more in local content moderation, instead of relying on global AI systems.

Blog Post by Caroline Allen, Guest Contributor
April 19, 2022 2:36 pm (EST)

Amhara militia members ride in the back of a truck towards a fight with the Tigray People's Liberation Front. Reuters/Tiksa Negeri

ARTIFICIAL INTELLIGENCE

Google's medical AI was super accurate in a lab. Real life was a different story.

If AI is really going to make a difference, it needs to work when real humans are involved.

By Will Douglas Heaven

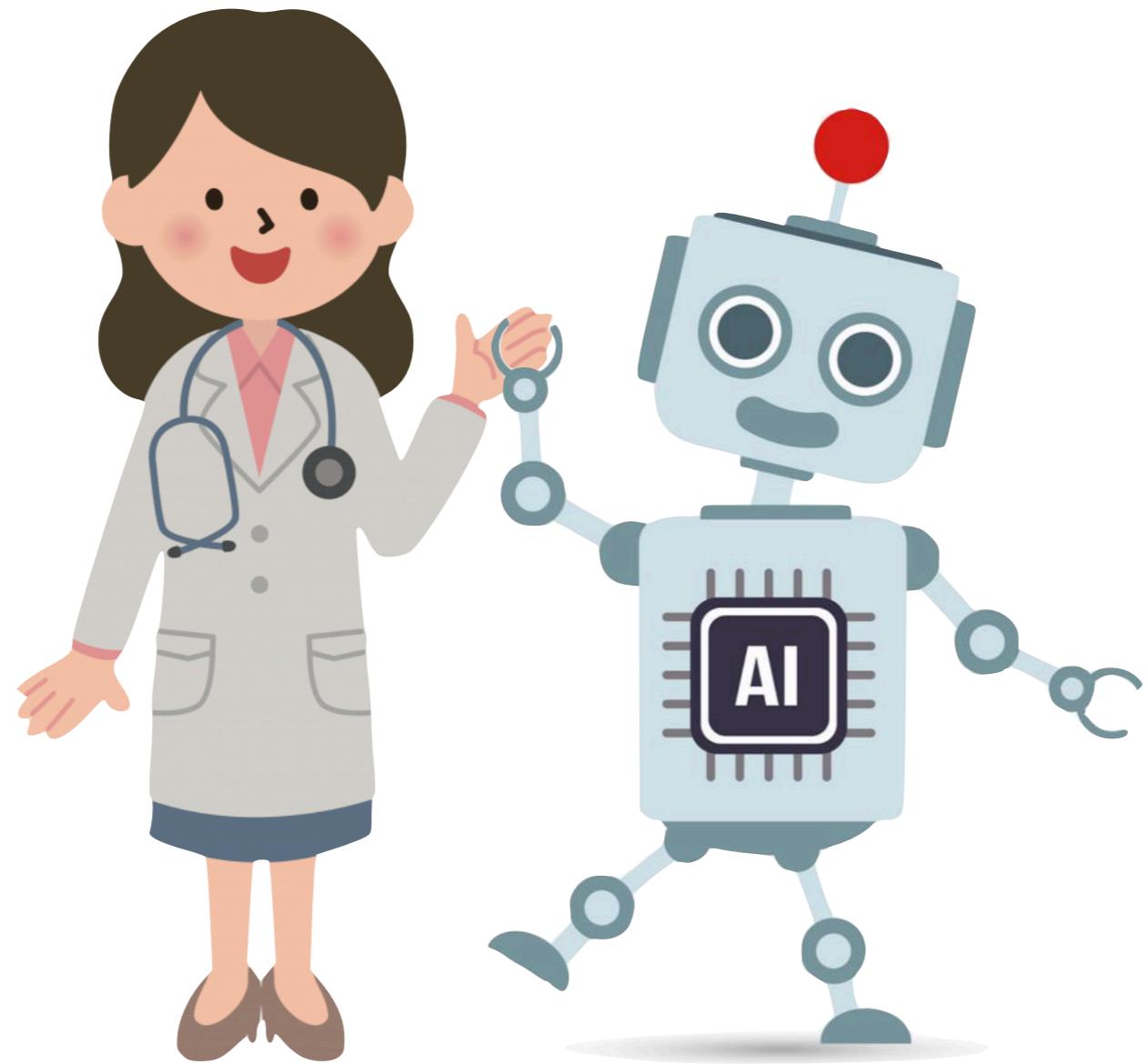
Medscape Tuesday, December 13, 2022

NEWS & PERSPECTIVE DRUGS & DISEASES CME & EDUCATION ACADEMY VIDEO DECISION POINT

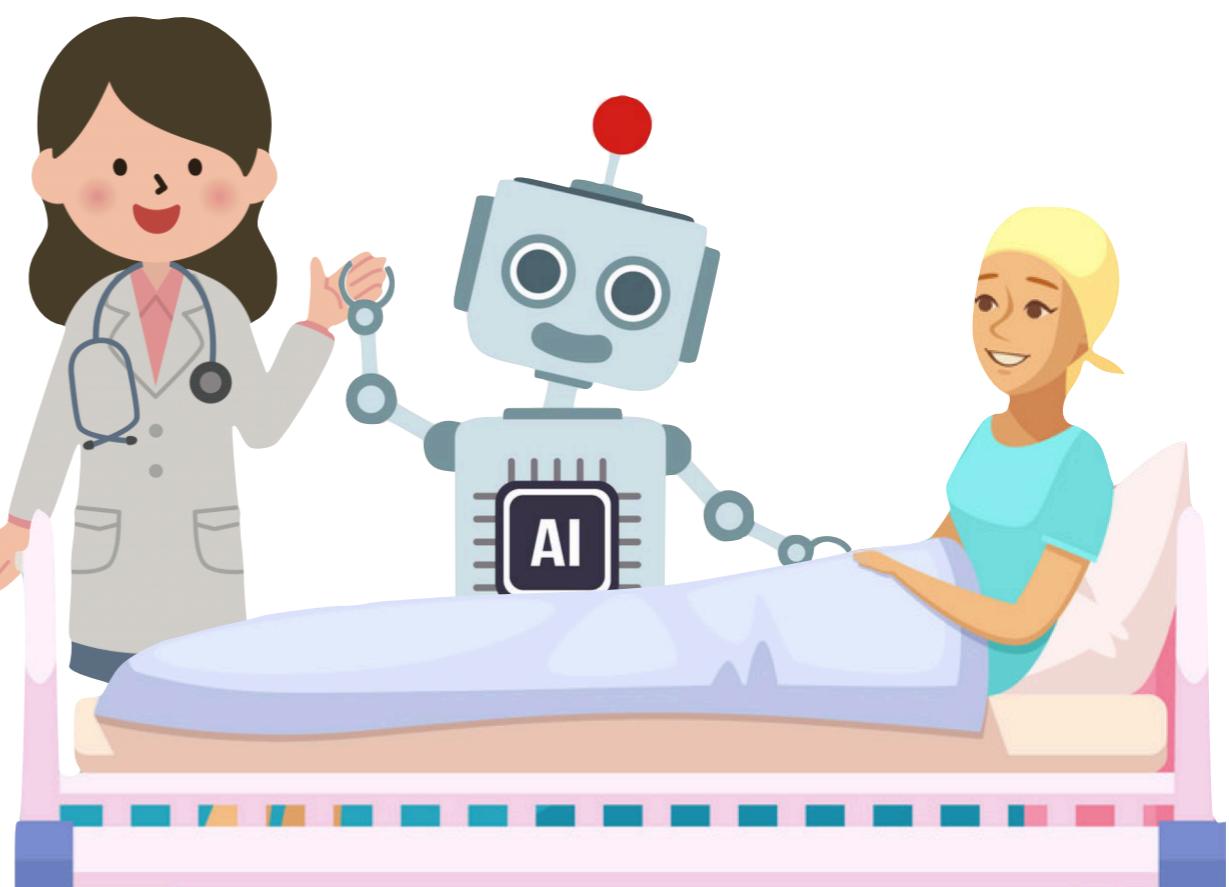
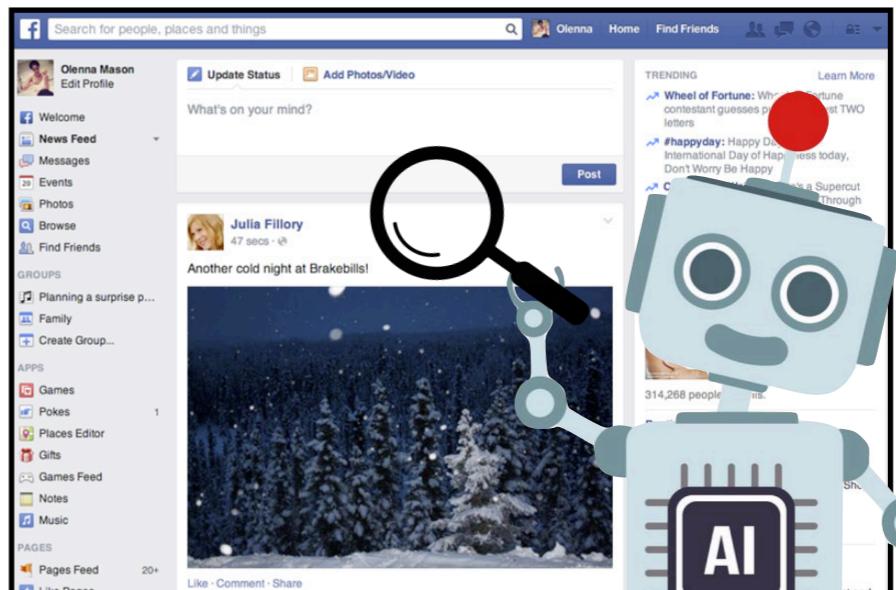
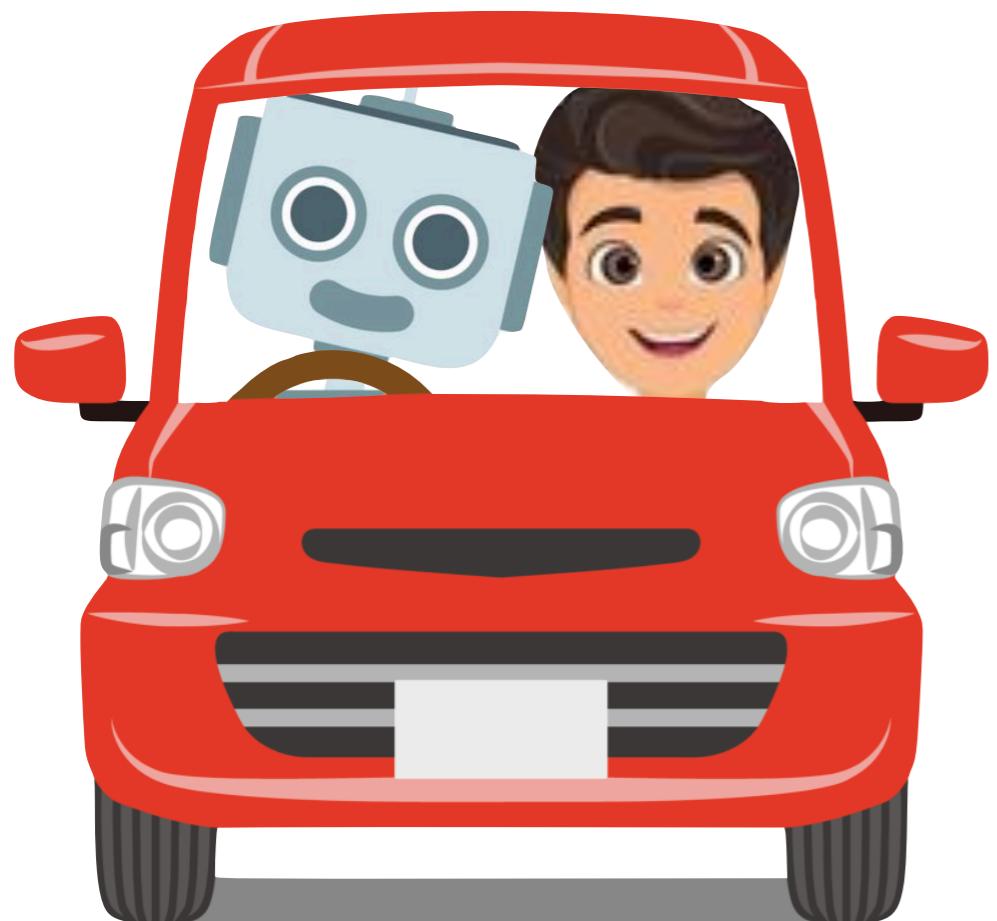
News > Medscape Medical News > Conference News > CHEST 2022

Sepsis Predictor Tool Falls Short in Emergency Setting

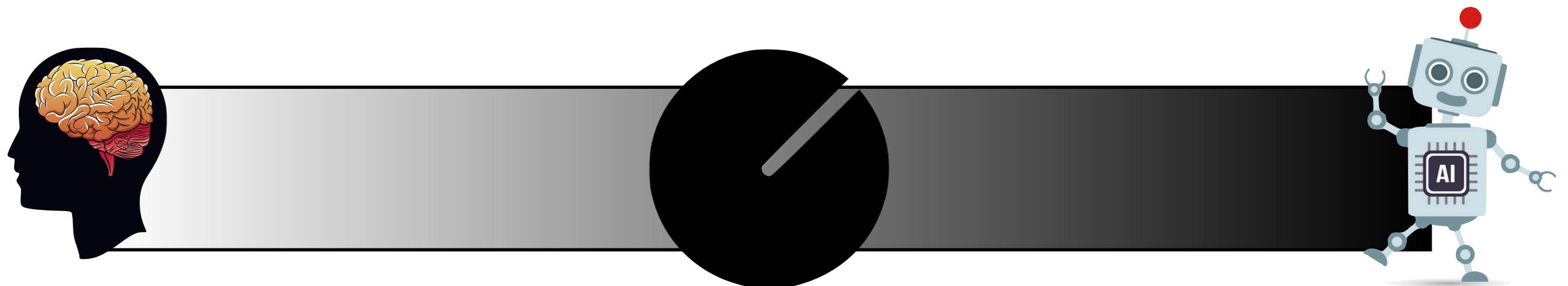
Heidi Splete
October 17, 2022



human-AI collaboration

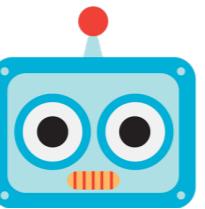


safe, gradual automation





human



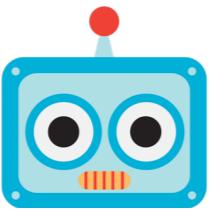
model

\hat{y}

prediction



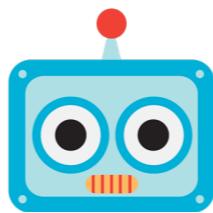
human



model

\hat{y}

prediction



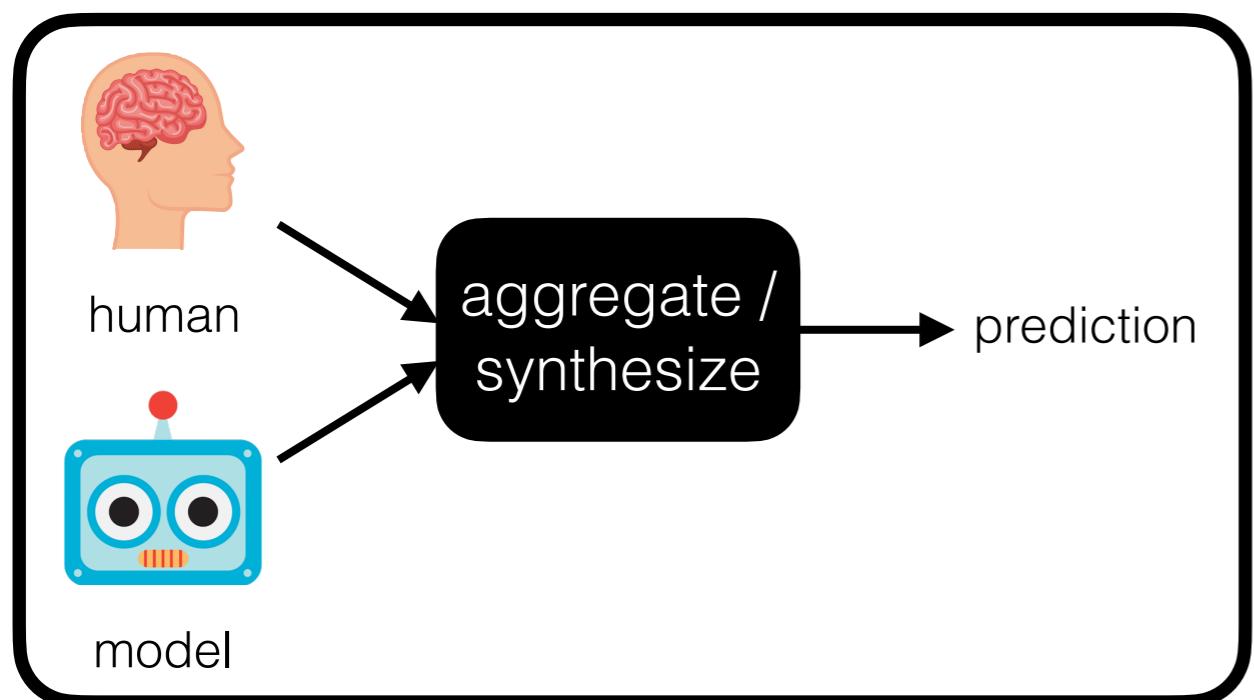
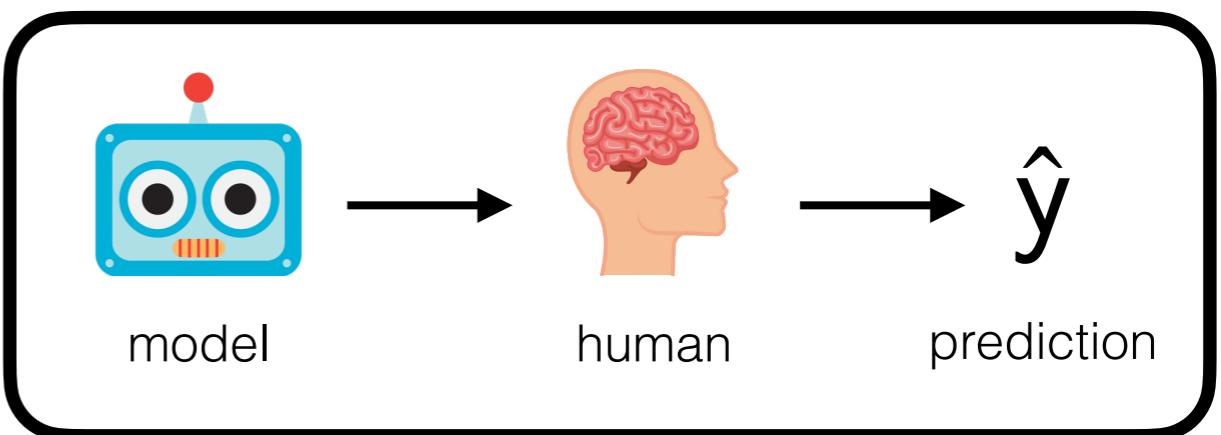
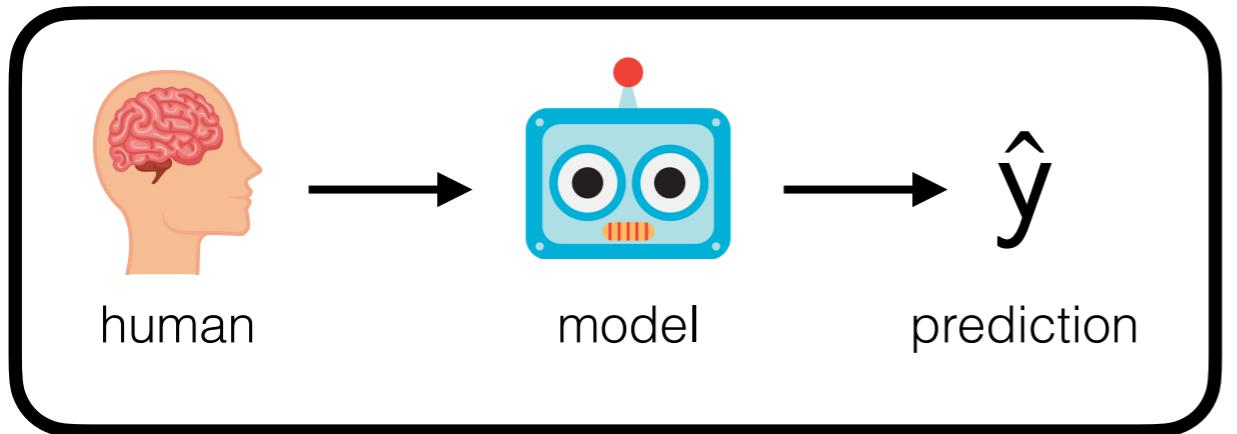
model

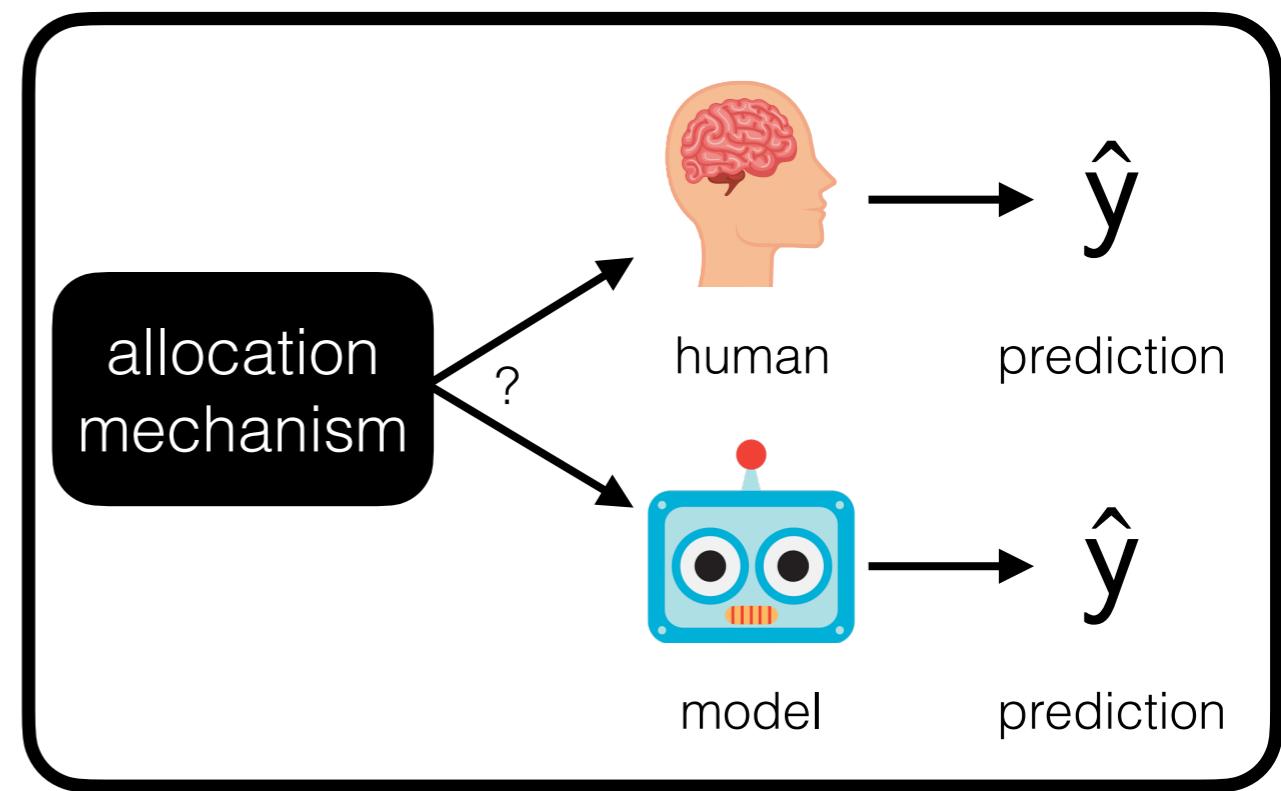
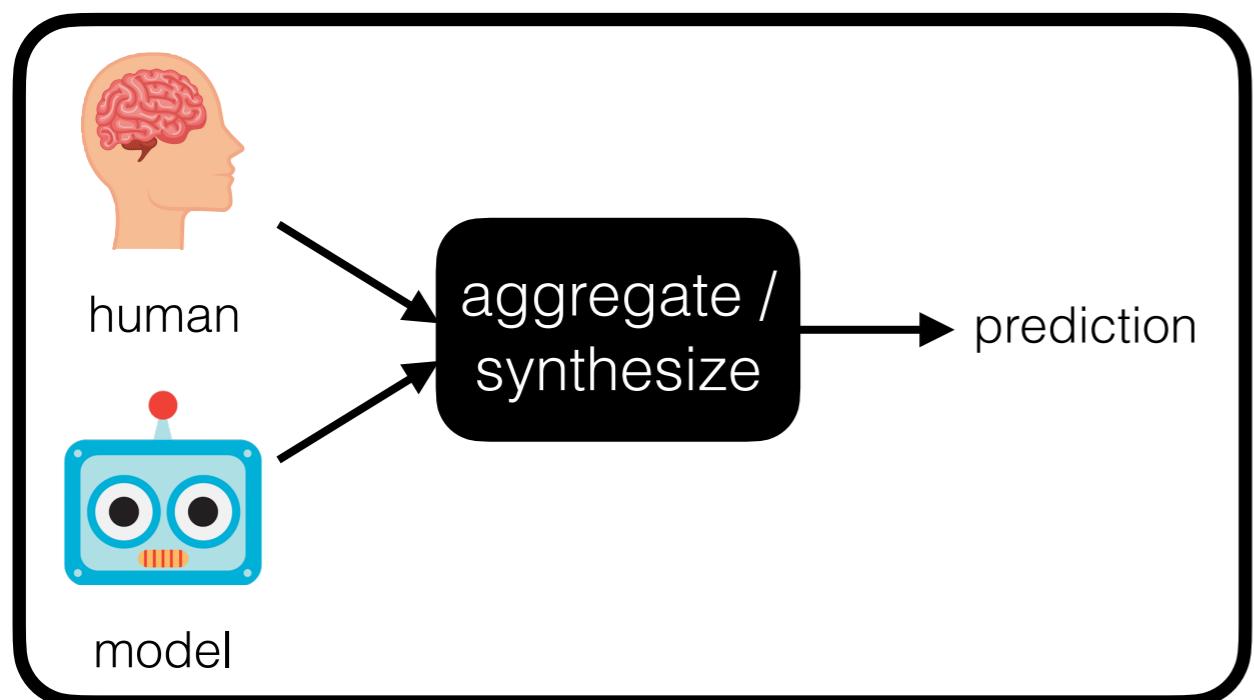
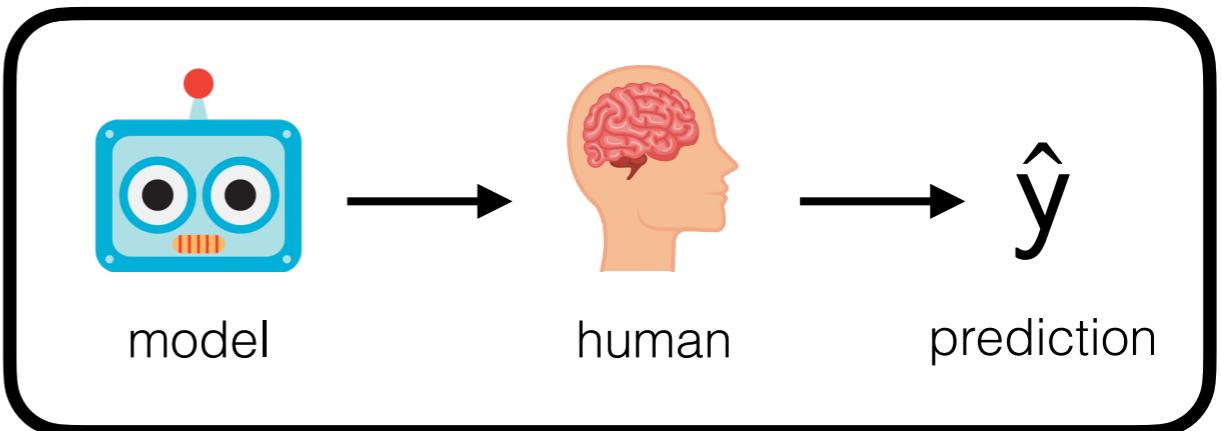
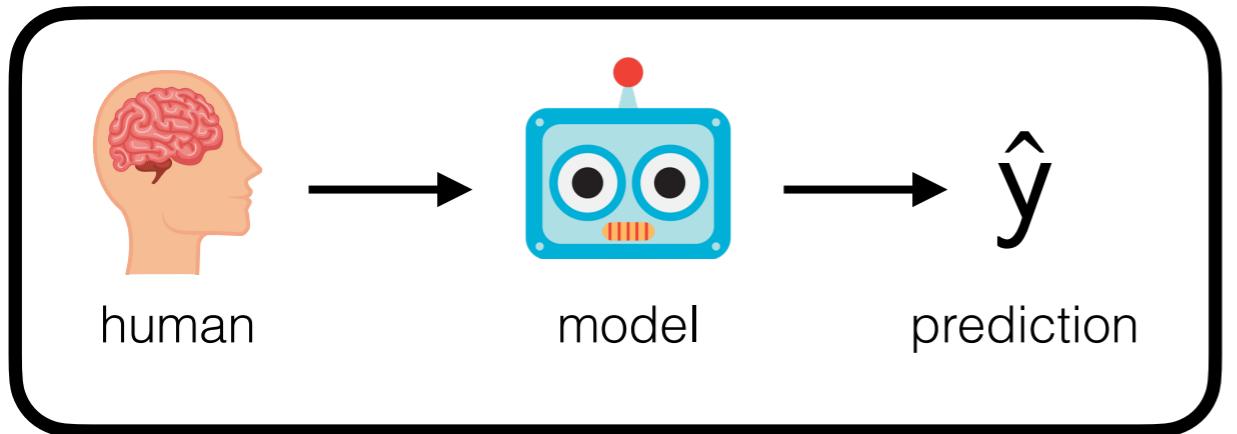


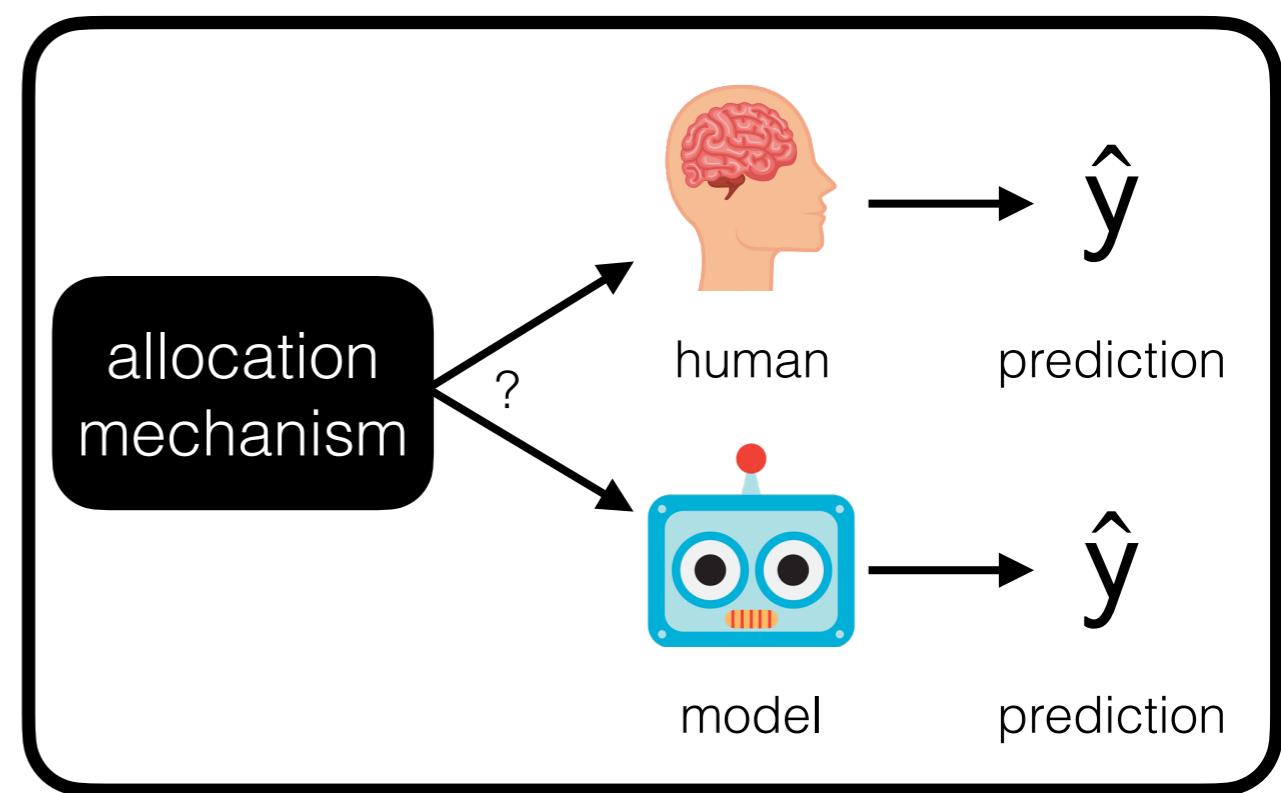
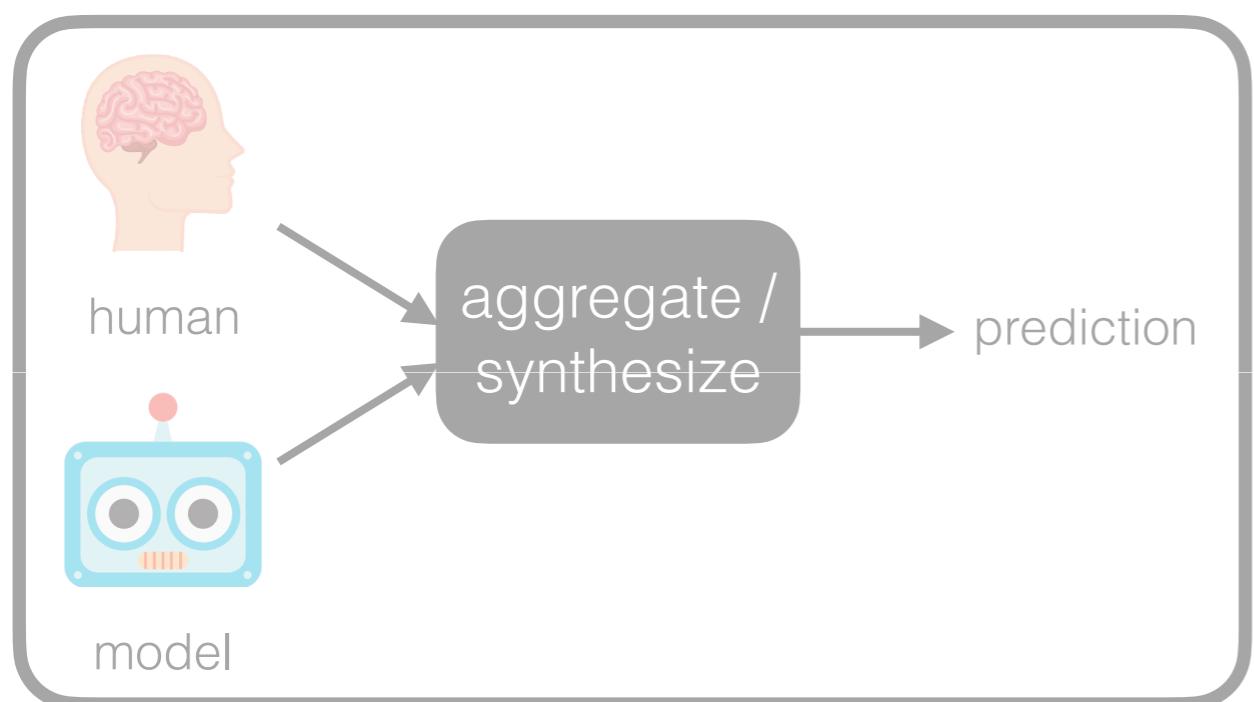
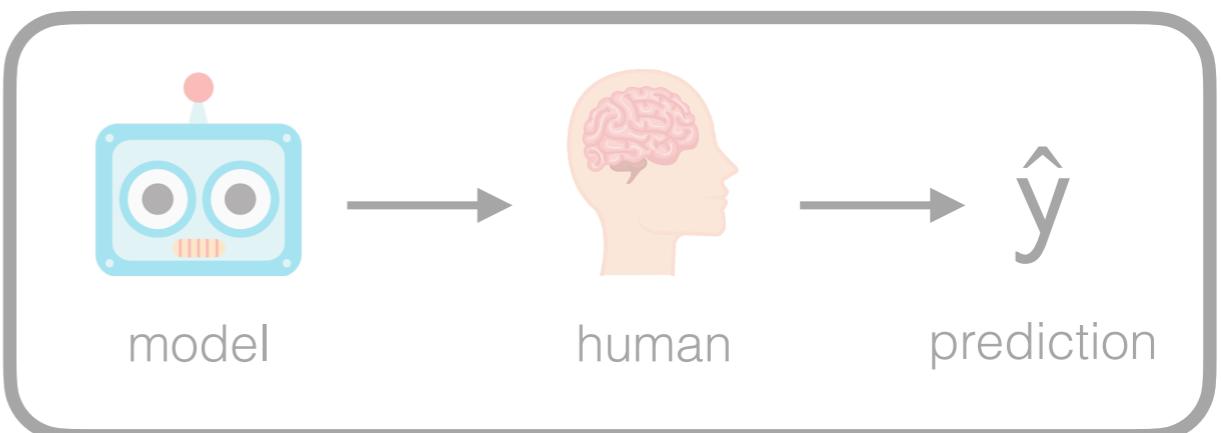
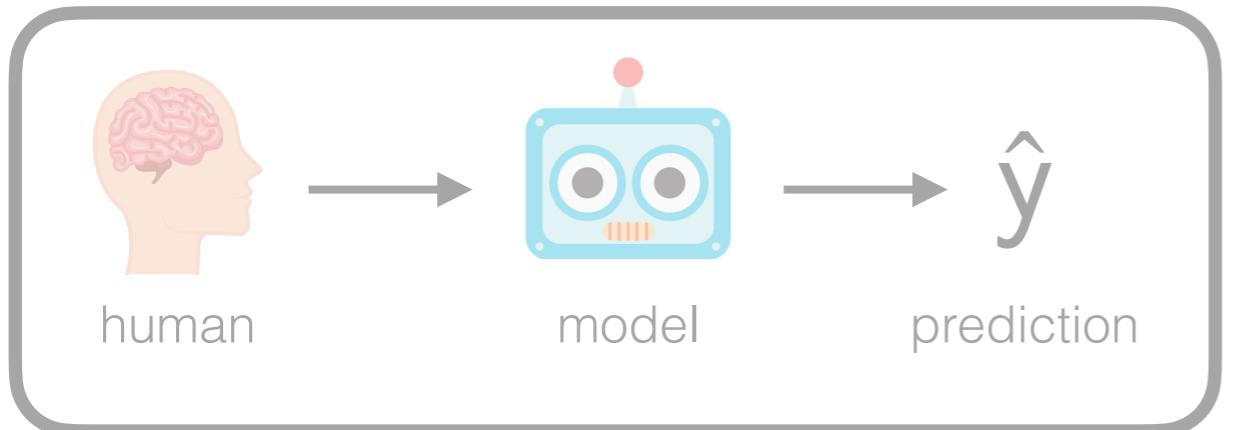
human

\hat{y}

prediction



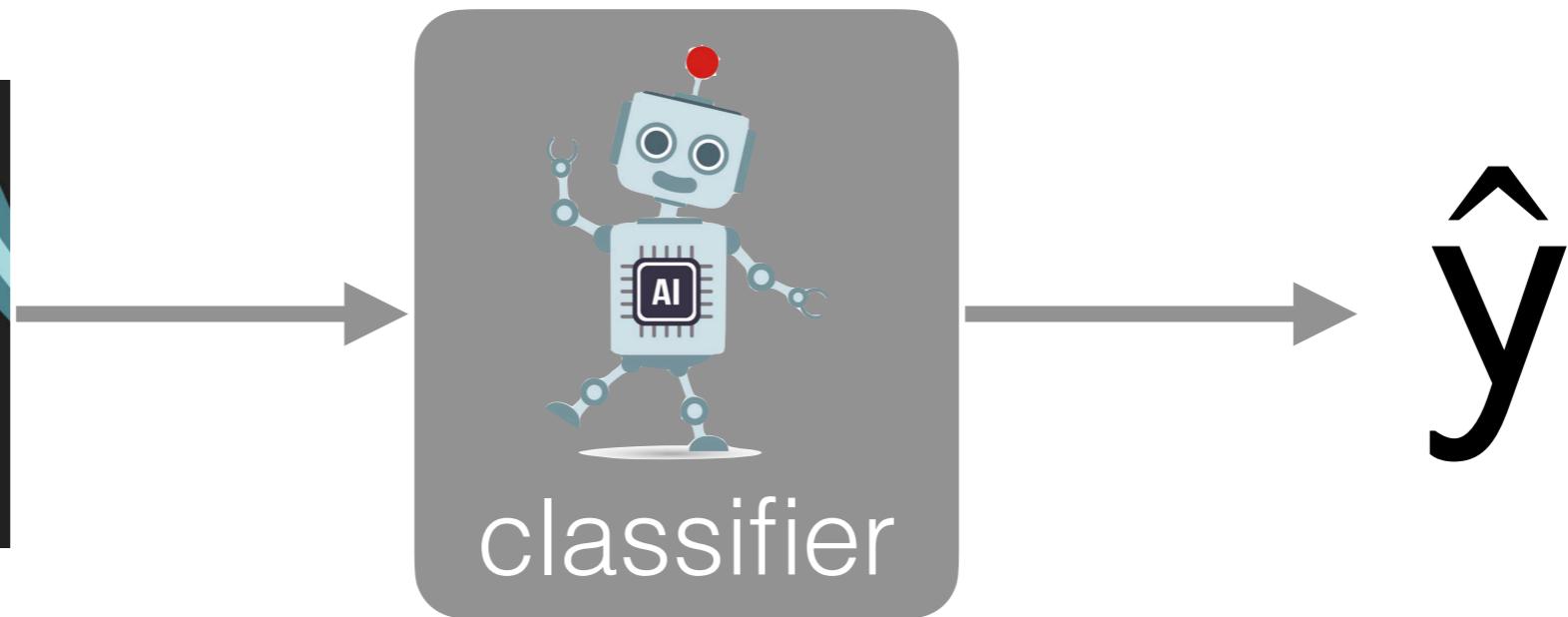




- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

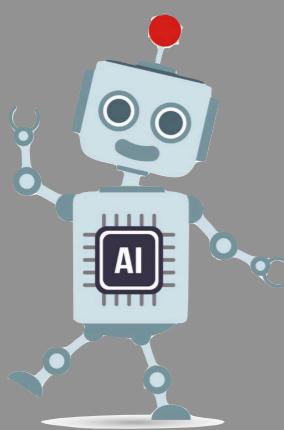
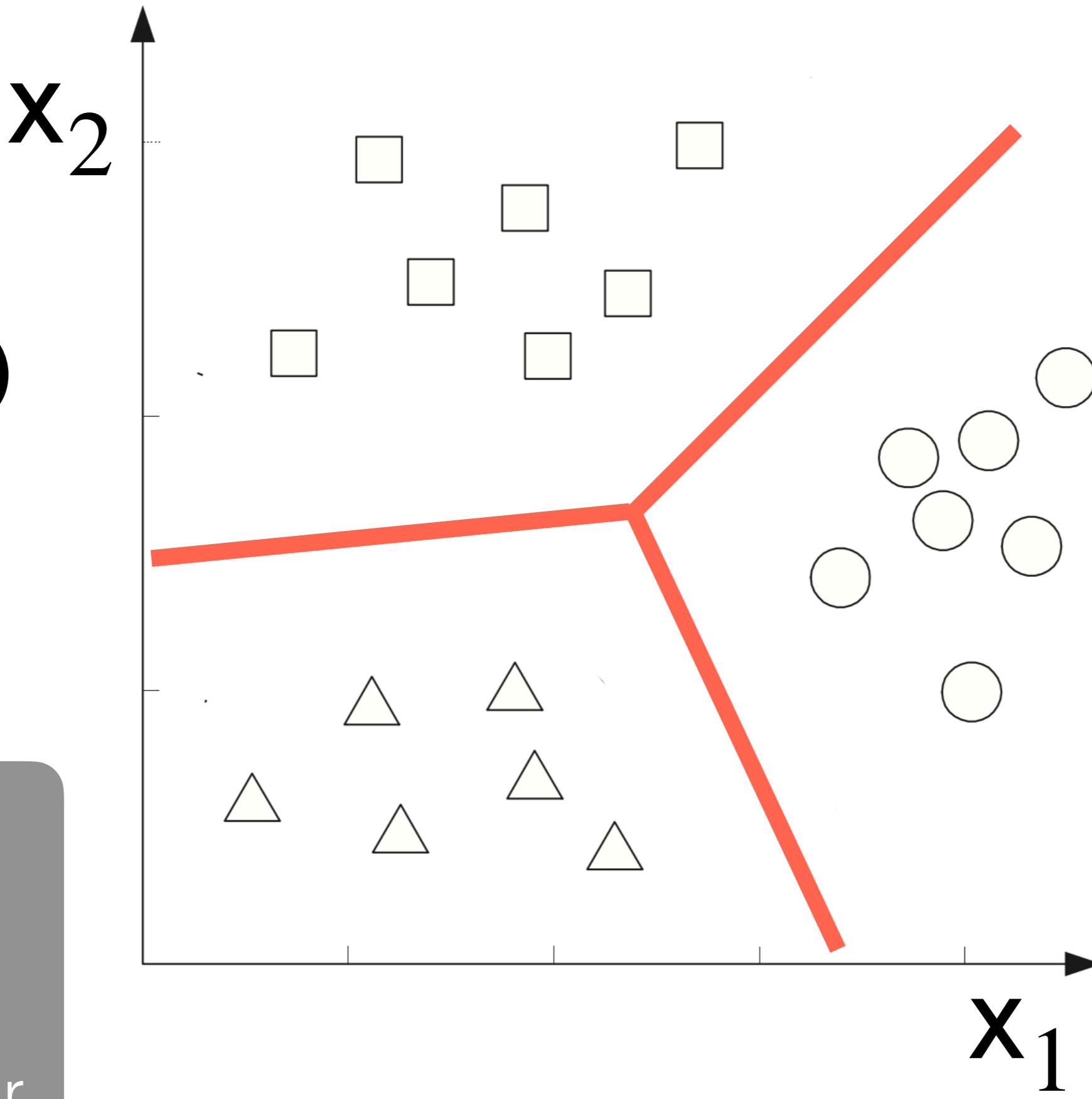
- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

input
features

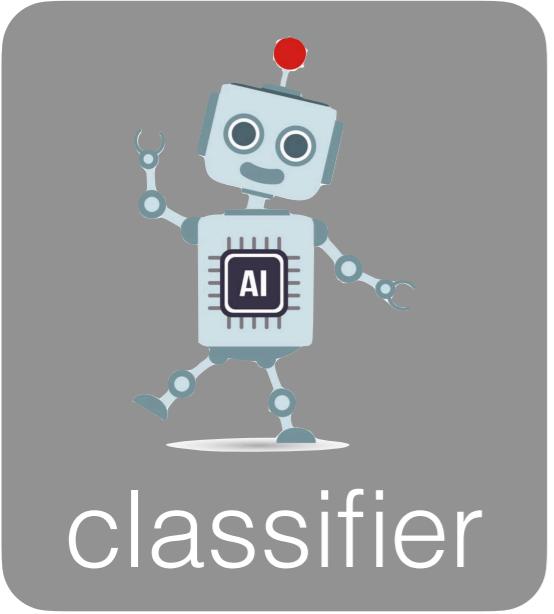


traditional classifier

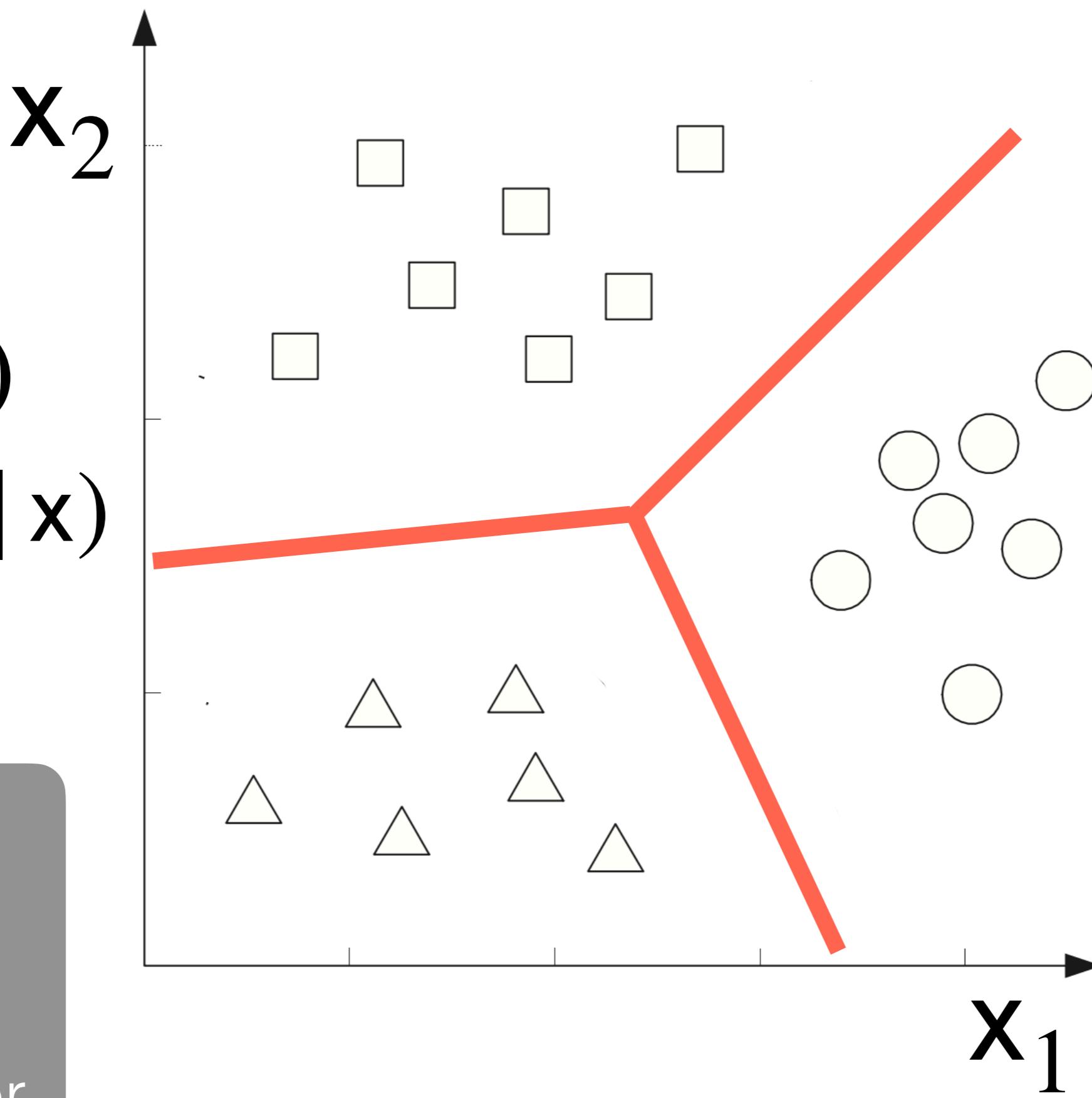
$p(y | x)$



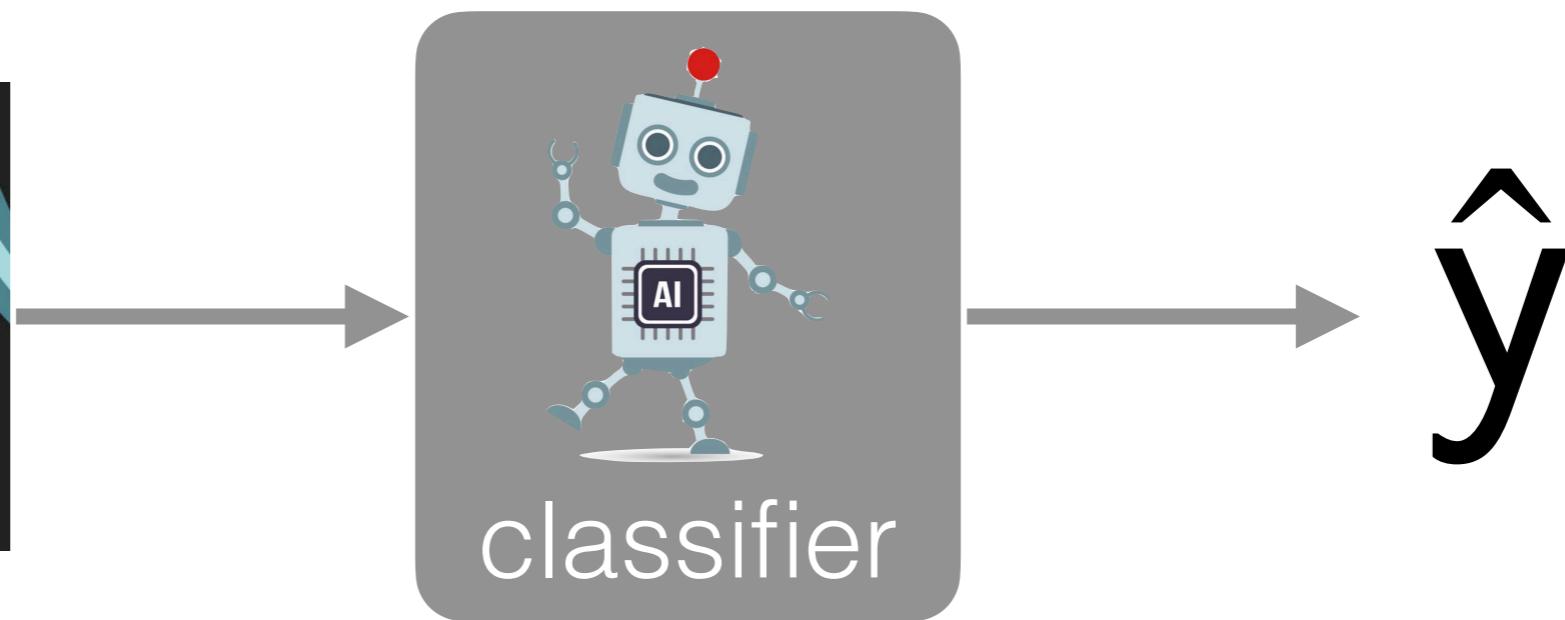
classifier



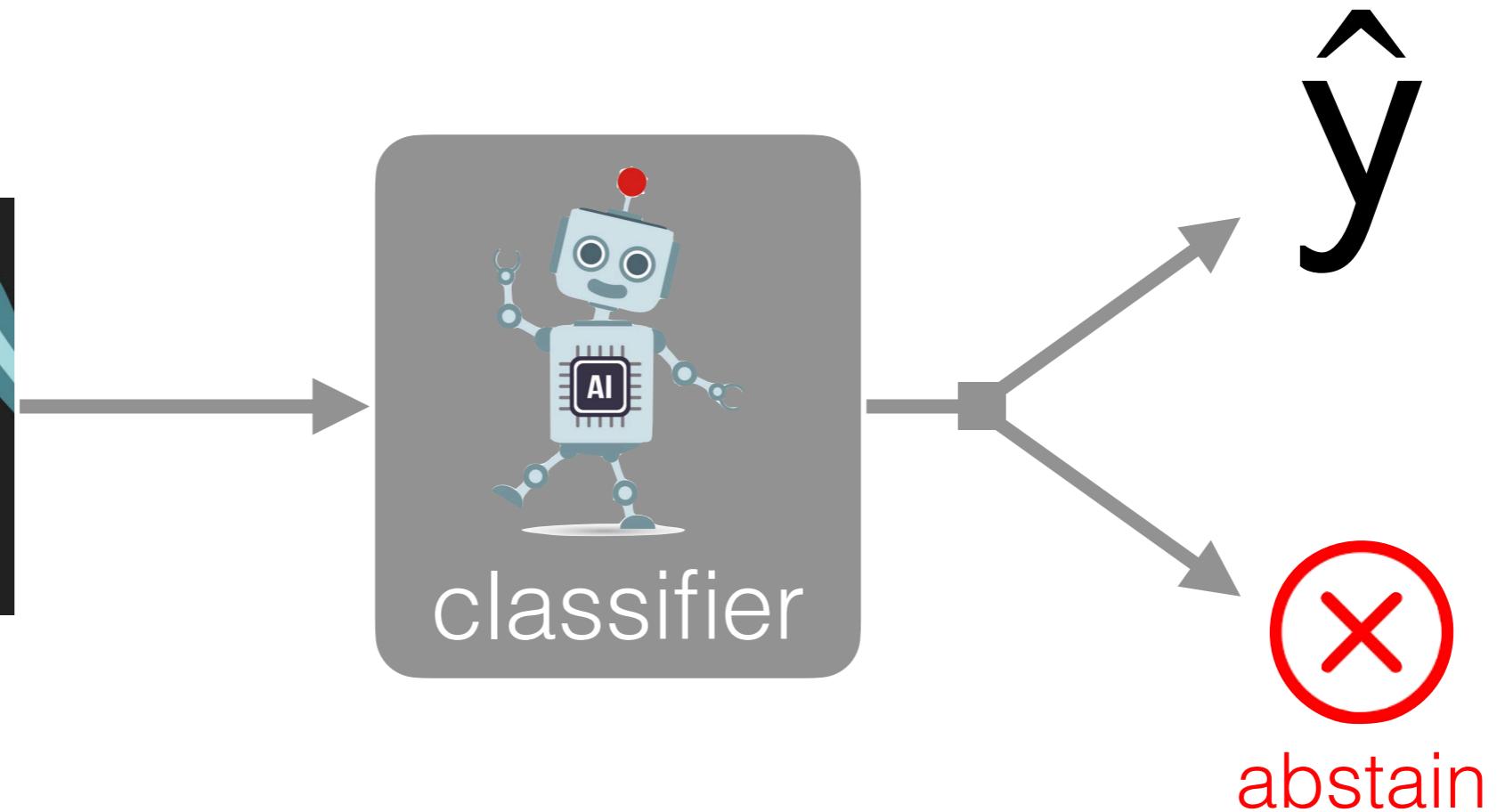
$p(y | x)$
 $\approx \mathbb{P}(y | x)$



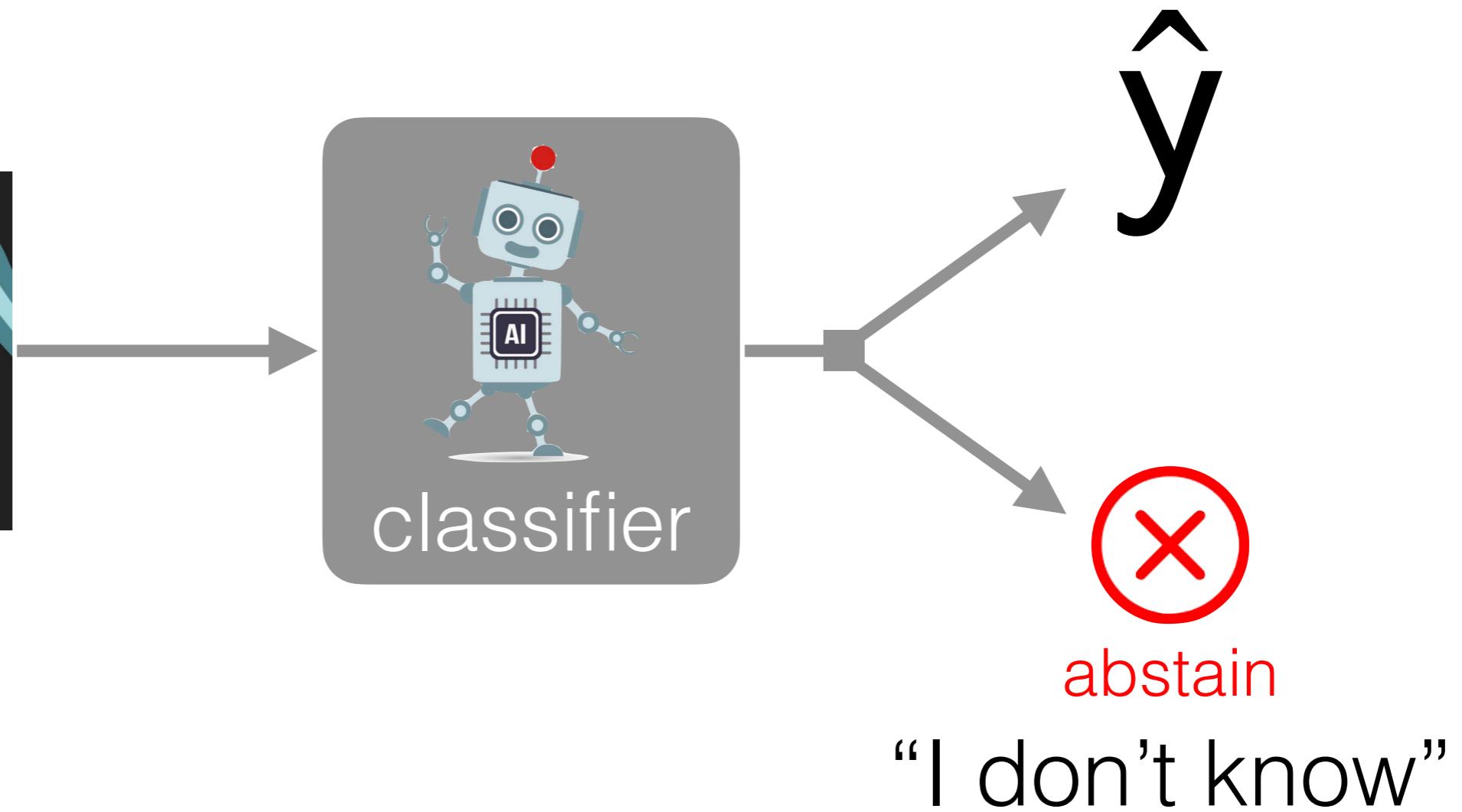
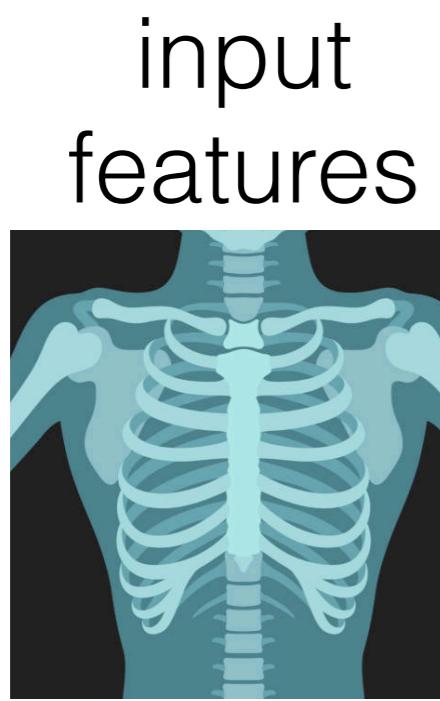
input
features



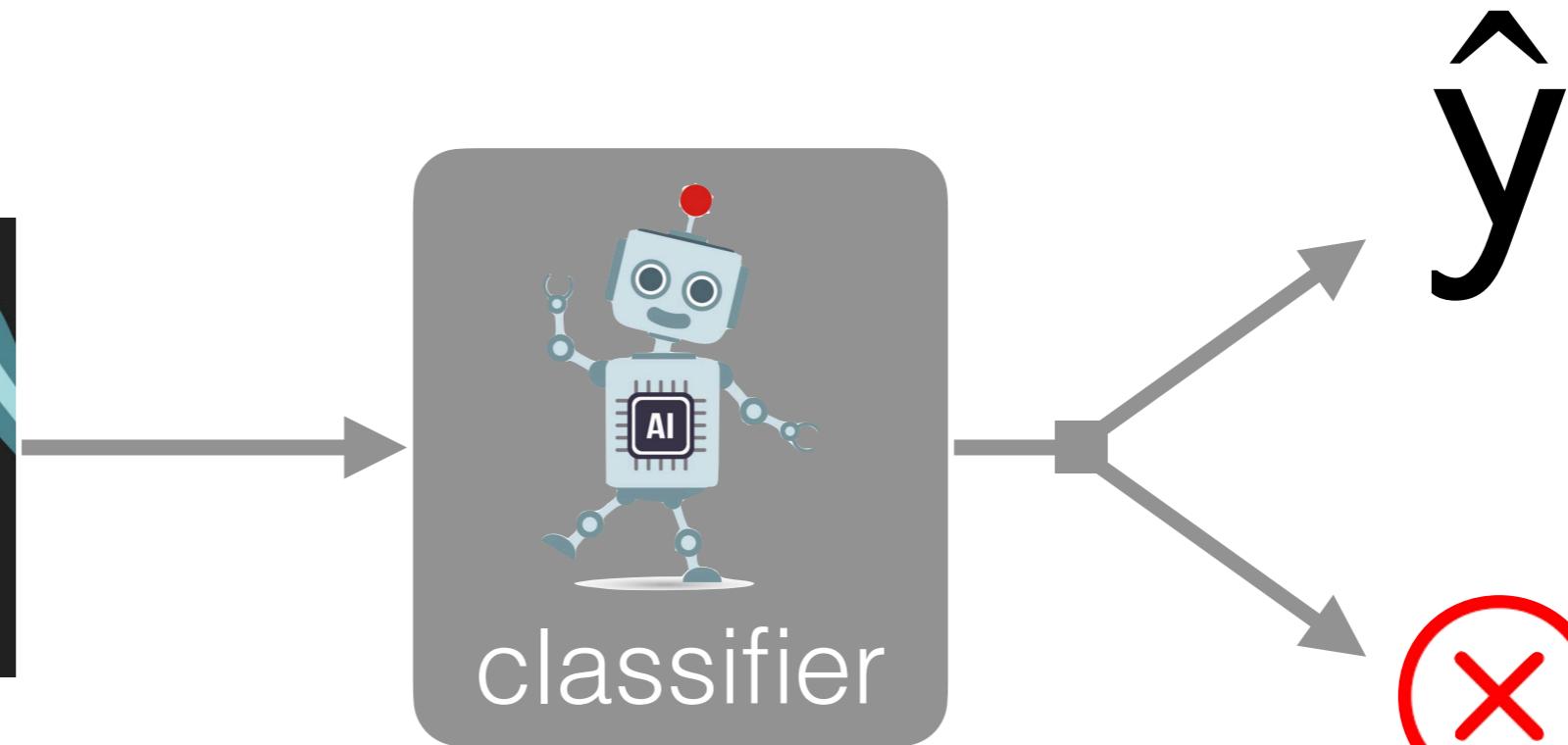
traditional classifier



learning to reject

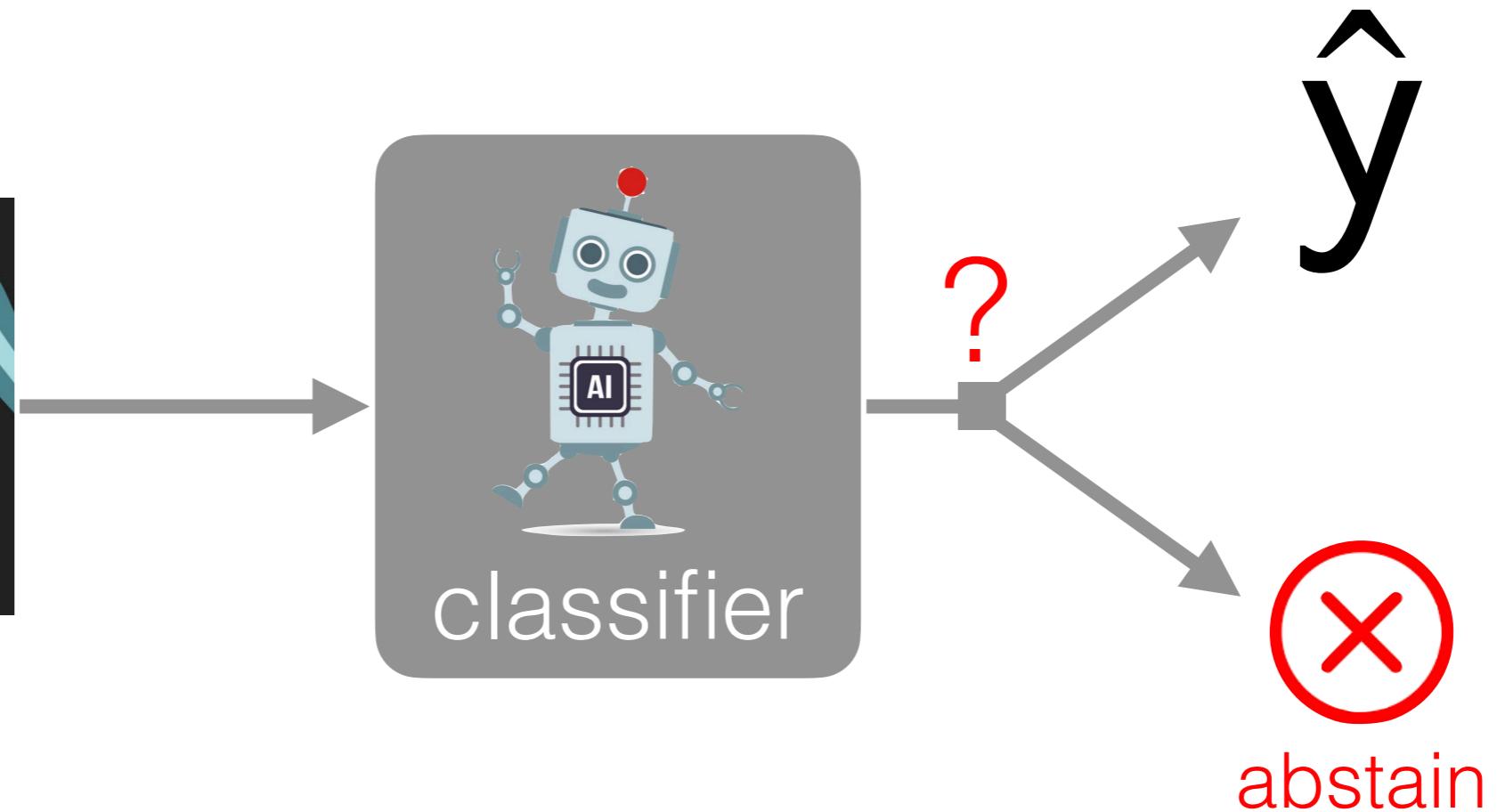


learning to reject



learning to reject

input
features



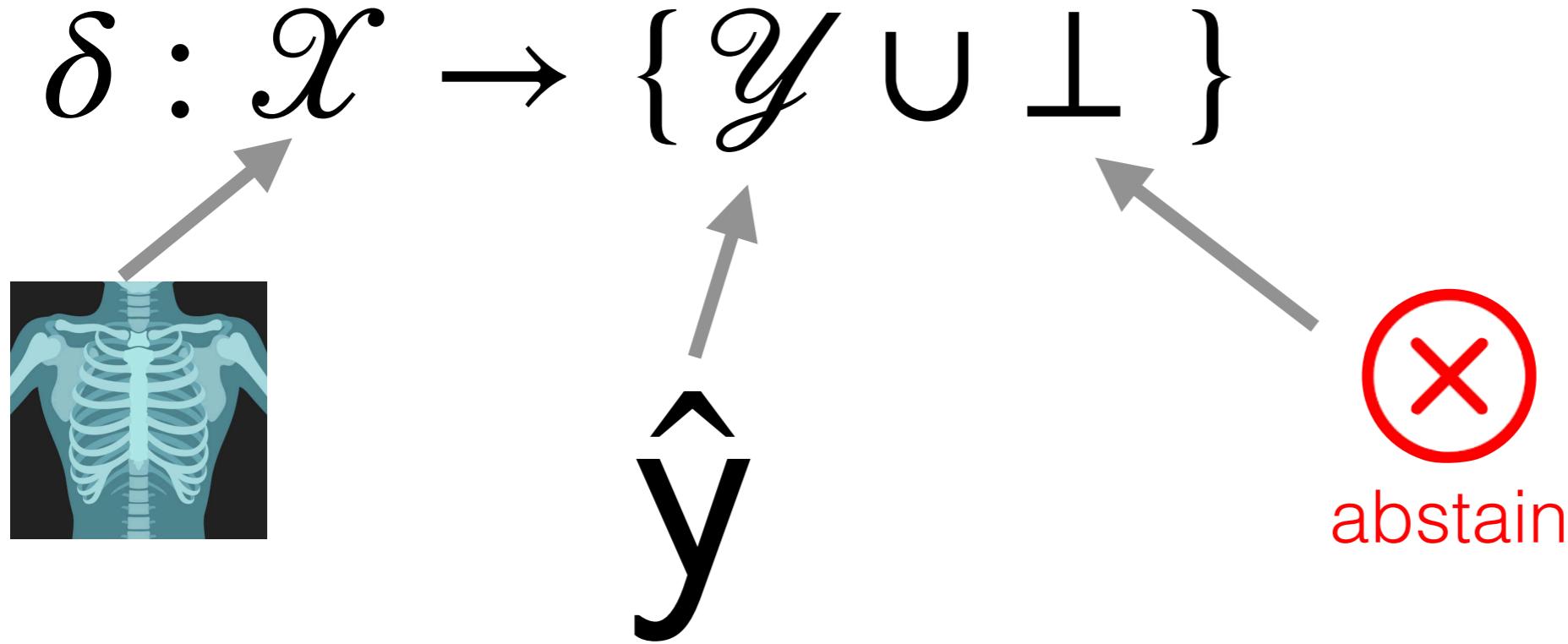
learning to reject

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

Derivation from Bayes decision theory

$$\delta : \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

Derivation from Bayes decision theory



Derivation from Bayes decision theory

$$\delta: \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

$$c(\delta(x), y) = \begin{cases} 0 & \textbf{if } \delta(x) = y \\ 1 & \textbf{if } \delta(x) \neq y \\ \lambda & \textbf{if } \delta(x) = \perp \end{cases}$$

Derivation from Bayes decision theory

$$\delta: \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

$$c(\delta(x), y) = \begin{cases} 0 & \textbf{if } \delta(x) = y \\ 1 & \textbf{if } \delta(x) \neq y \\ \lambda & \textbf{if } \delta(x) = \perp \end{cases}$$



cost of rejecting
 $\lambda \in (0, 1)$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

when we choose to make a prediction...

$$\delta(x) = \hat{y} = \arg \max_y P(y | x)$$

Derivation from Bayes decision theory

$$\mathcal{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

when we choose to make a prediction...

$$\delta(x) = \hat{y} = \arg \max_y P(y | x)$$

$$\mathcal{R}(x, \hat{y}) = \sum_{y \neq \hat{y}} P(y | x) = 1 - P(\hat{y} | x)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

when we choose to reject...

$$\delta(x) = \perp$$

Derivation from Bayes decision theory

$$\mathcal{R}(x, \delta) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot c(\delta(x), y)$$

when we choose to reject...

$$\delta(x) = \perp$$

$$\mathcal{R}(x, \perp) = \sum_{y \in \mathcal{Y}} P(y | x) \cdot \lambda = \lambda$$

Derivation from Bayes decision theory

Thus we need to choose the action with the minimum risk:

$$\mathcal{R}(x, \hat{y}) = 1 - P(\hat{y} | x)$$

vs

$$\mathcal{R}(x, \perp) = \lambda$$

Derivation from Bayes decision theory

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y | x) & \text{if } 1 - \max_y \mathbb{P}(y | x) \leq \lambda \\ \perp & \text{otherwise} \end{cases}$$

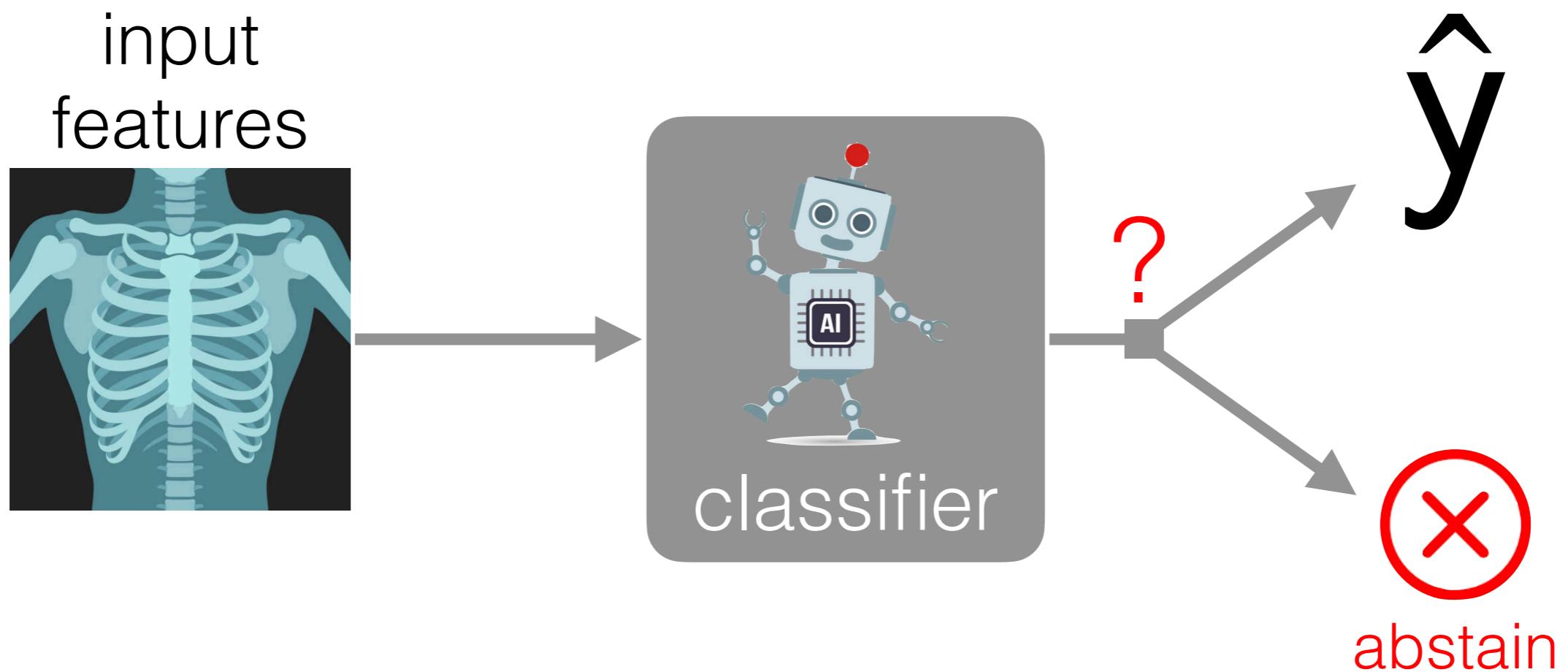
Derivation from Bayes decision theory

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y | x) & \text{if } 1 - \max_y \mathbb{P}(y | x) \leq \lambda \\ \perp & \text{otherwise} \end{cases}$$

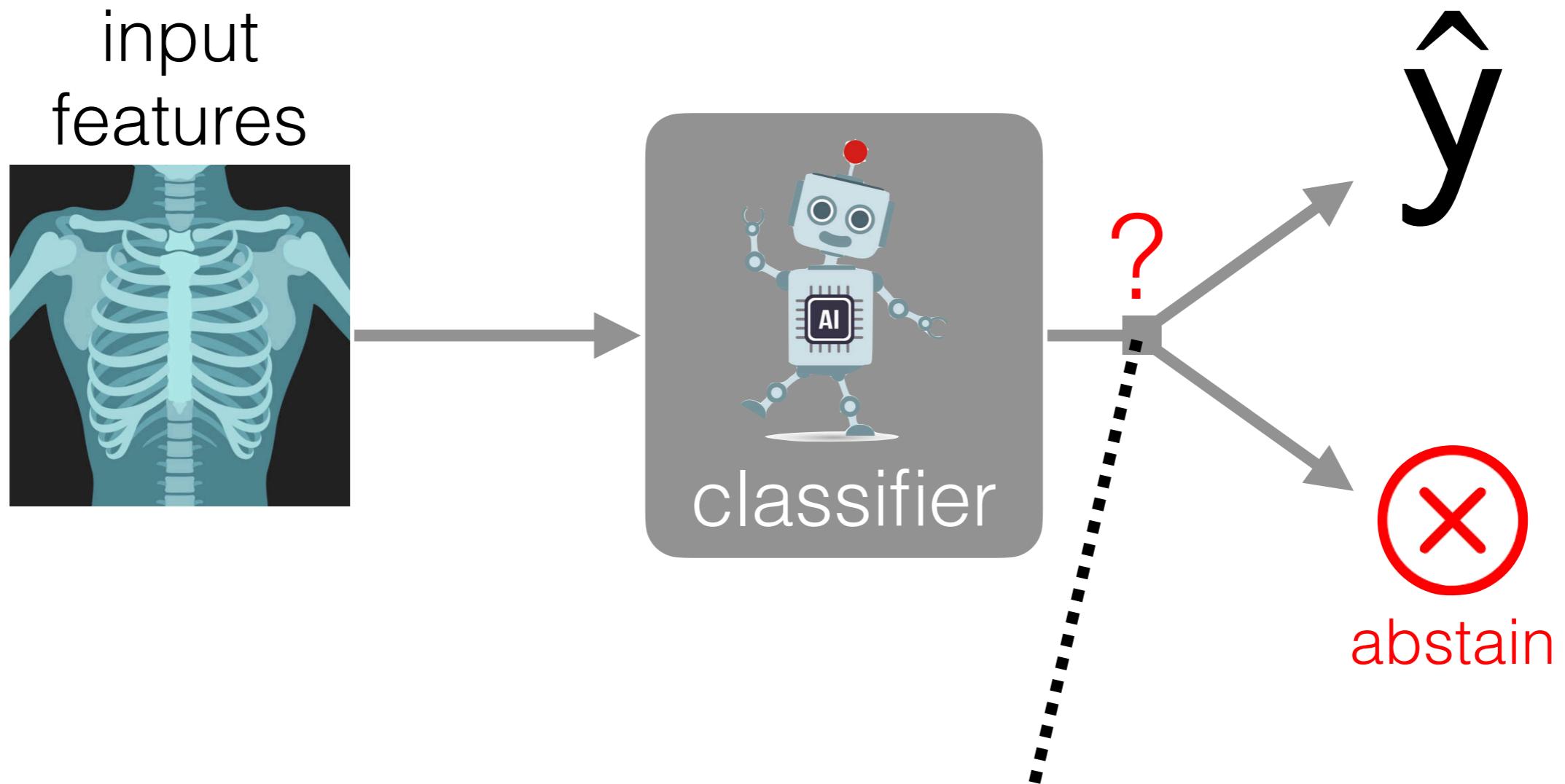
rearranging a bit...

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y | x) & \text{if } \max_y \mathbb{P}(y | x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

learning to reject

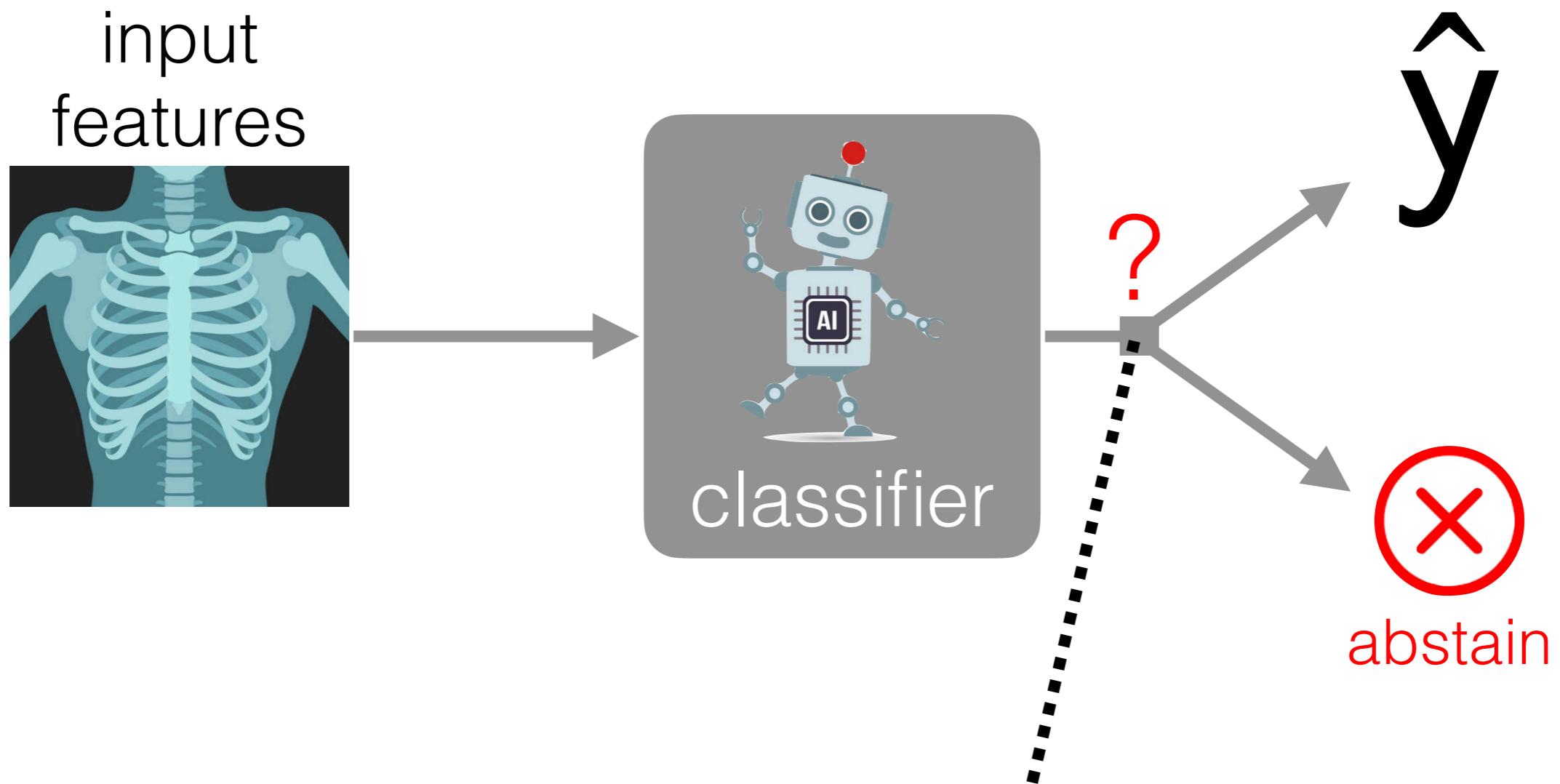


learning to reject



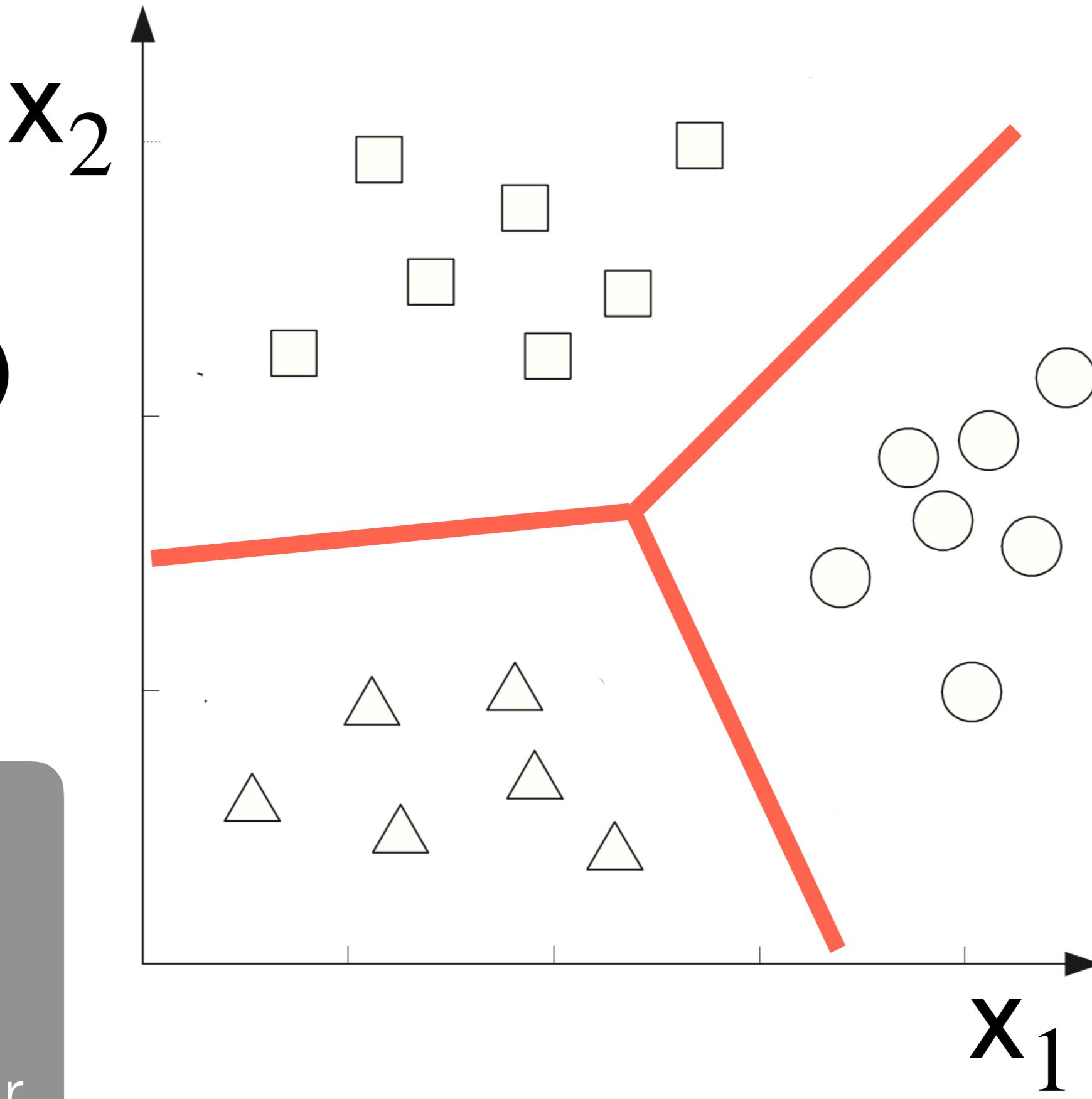
$$\delta^*(x) = \begin{cases} \hat{y} & \text{if } \max_y \mathbb{P}(y|x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

learning to reject

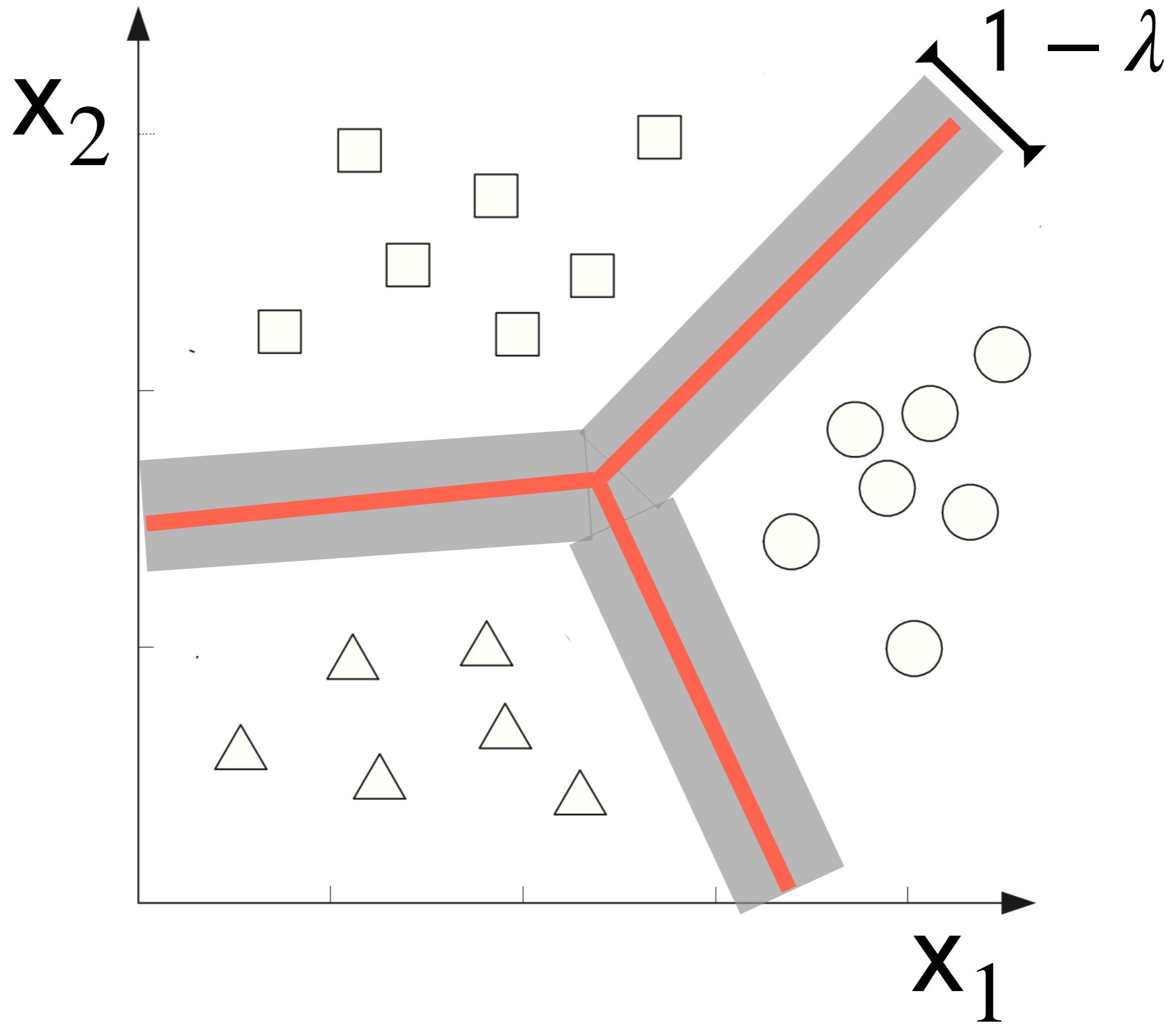


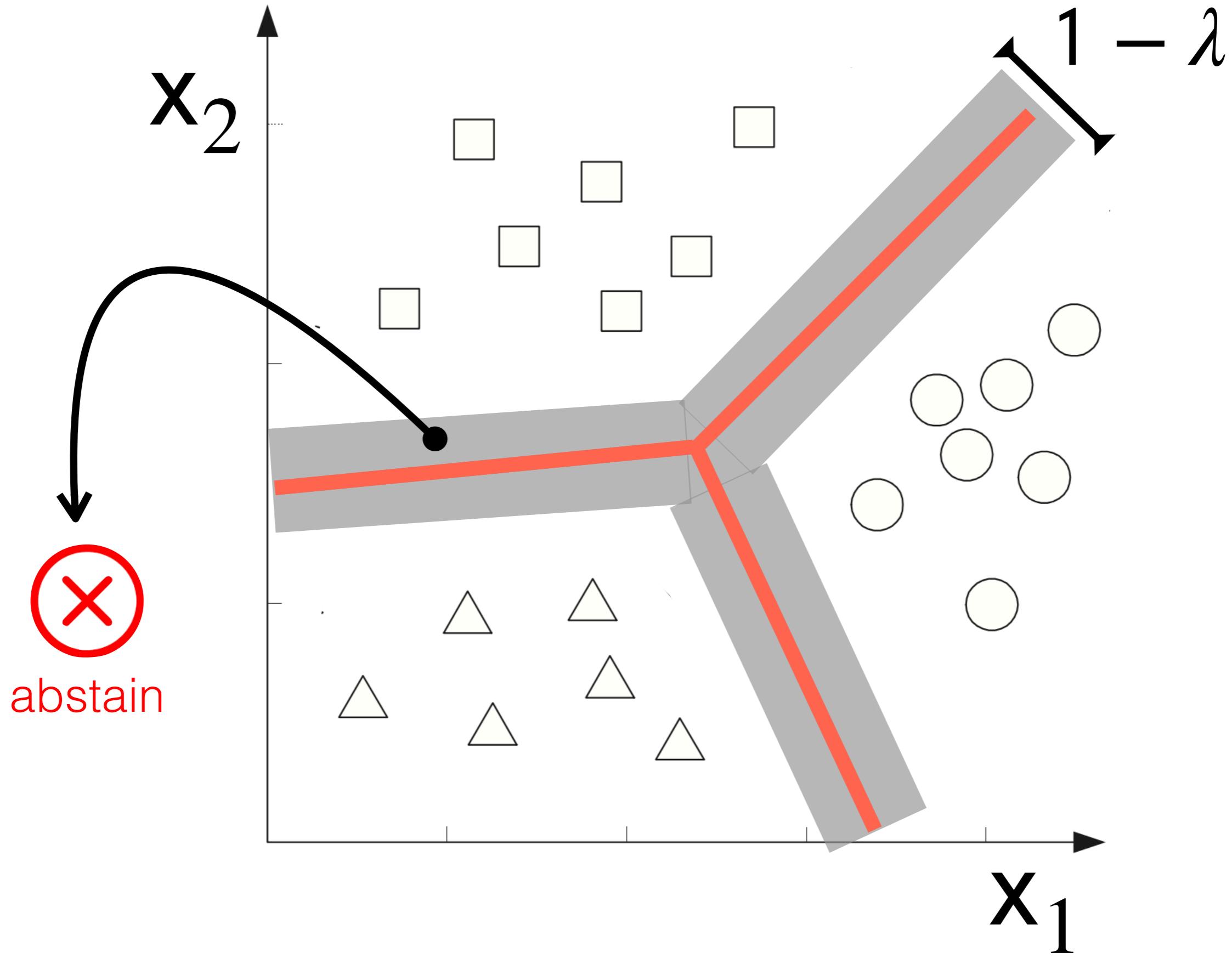
$$\hat{\delta}(x) = \begin{cases} \hat{y} & \text{if } \max_y p(y|x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

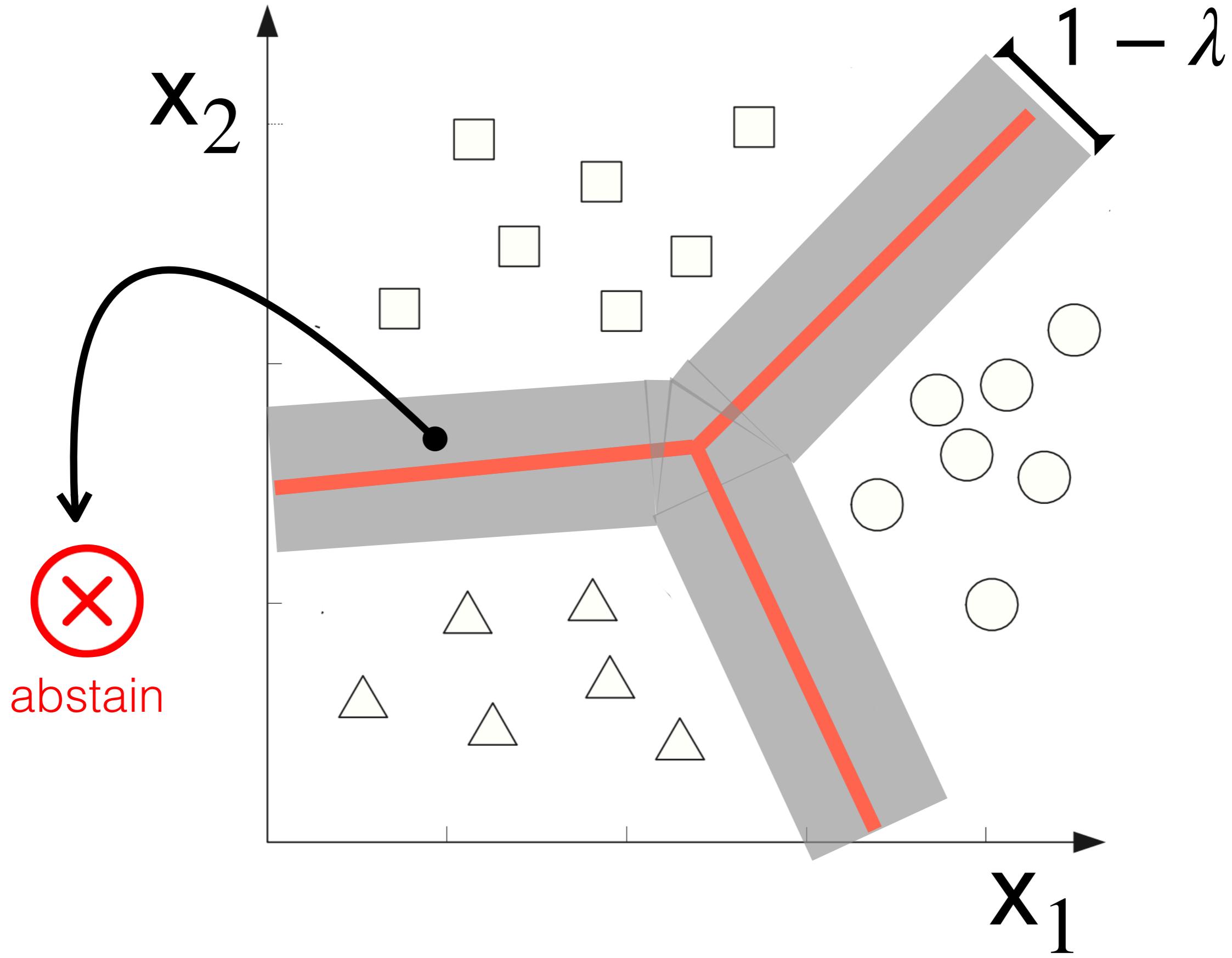
$p(y | x)$



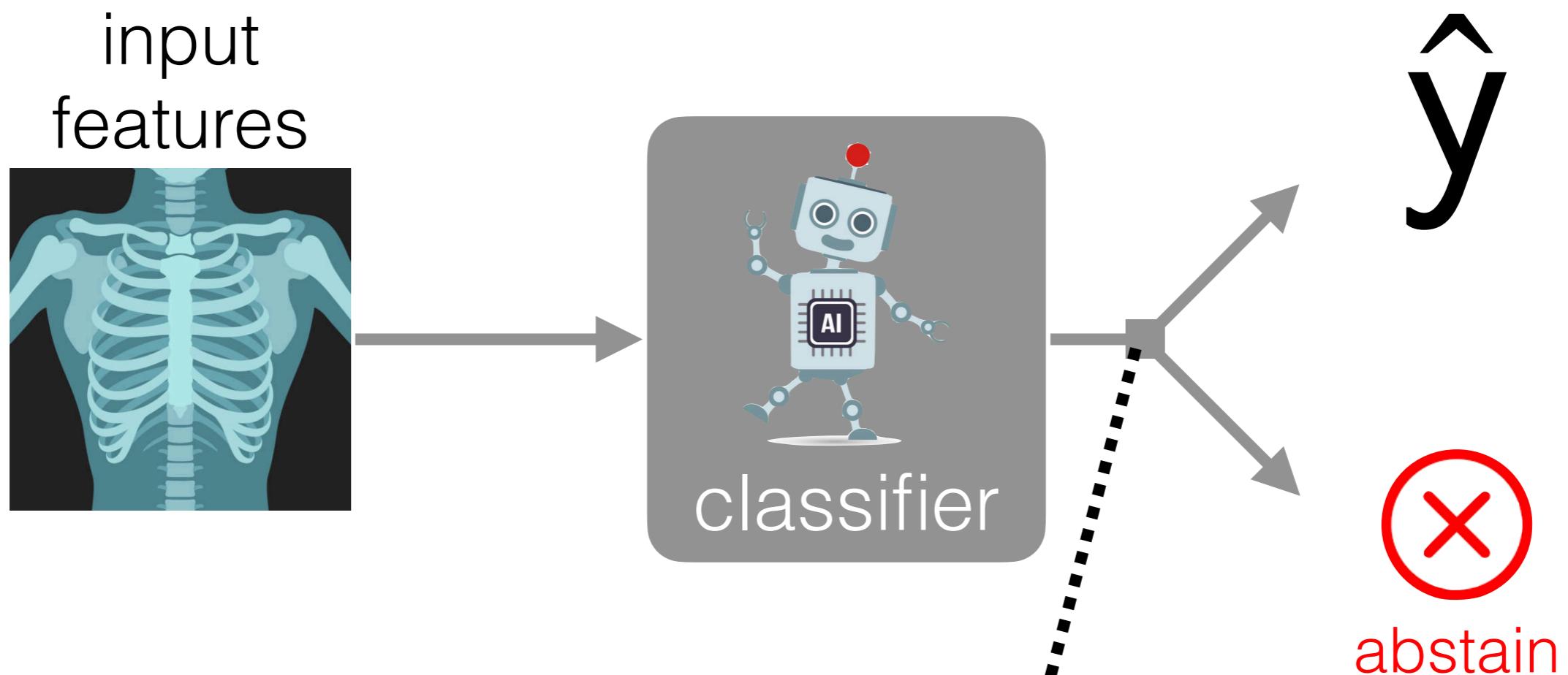
classifier







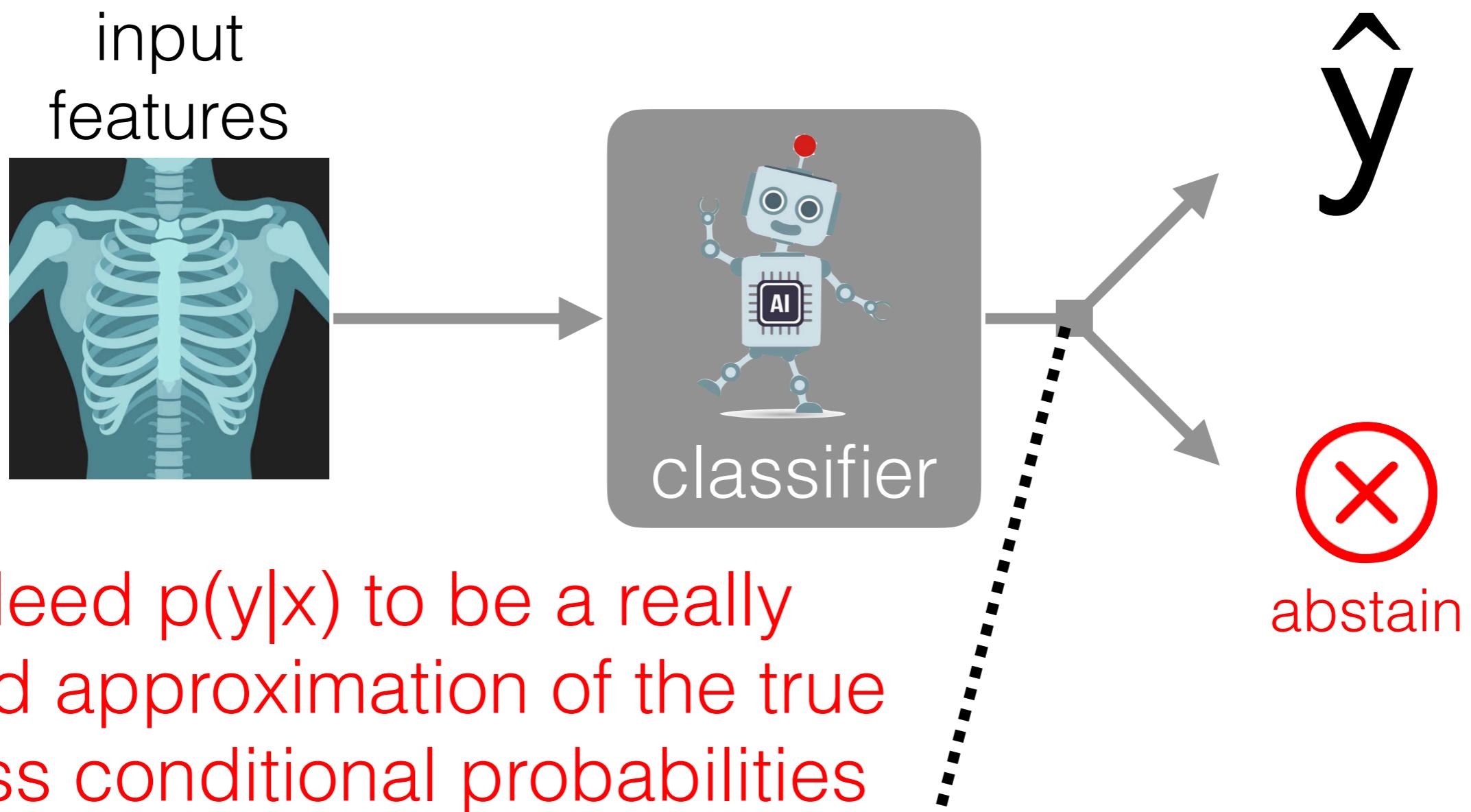
learning to reject



Any Problems?

$$\hat{\delta}(x) = \begin{cases} \hat{y} & \text{if } \max_y p(y|x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

learning to reject



$$\hat{\delta}(x) = \begin{cases} \hat{y} & \text{if } \max_y p(y | x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

developing a surrogate loss

$$\hat{\delta}(x) = \begin{cases} \hat{y} & \text{if } \max_y p(y|x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

key idea: the decision rule above is valid for *any* choice of λ .
If we can choose a specific λ ahead of time, we only need to calibrate the model around that λ .

developing a surrogate loss

$$\hat{\delta}(x) = \begin{cases} \hat{y} & \text{if } \max_y p(y|x) \geq 1 - \lambda \\ \perp & \text{otherwise} \end{cases}$$

key idea: the decision rule above is valid for *any* choice of λ .
If we can choose a specific λ ahead of time, we only need to
calibrate the model around that λ .

if $\mathbb{P}(y|x) \geq 1 - \lambda$, then $p(y|x) \geq 1 - \lambda$

if $\mathbb{P}(y|x) < 1 - \lambda$, then $p(y|x) < 1 - \lambda$

softmax implementation

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

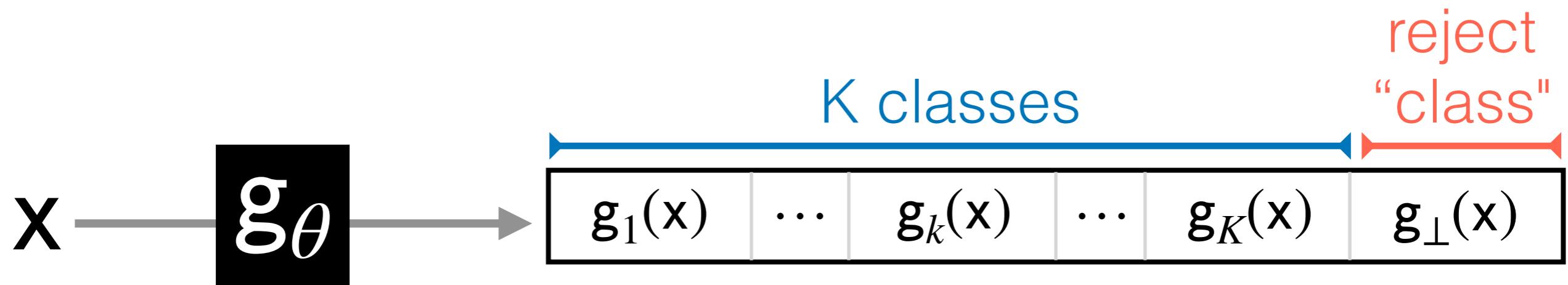
model

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, y_n \right\}_{n=1}^N$$

model



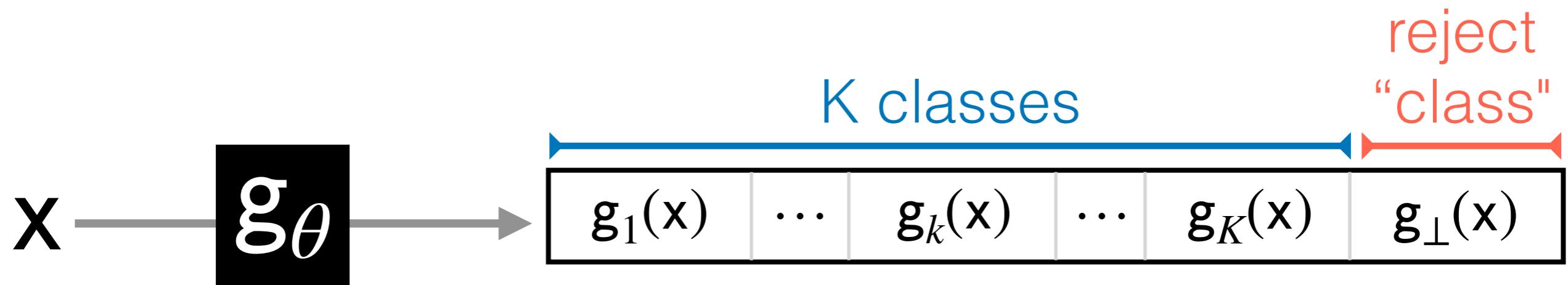
softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, y_n \right\}_{n=1}^N$$

model

$$g_k(\mathbf{x}) \in \mathbb{R}$$

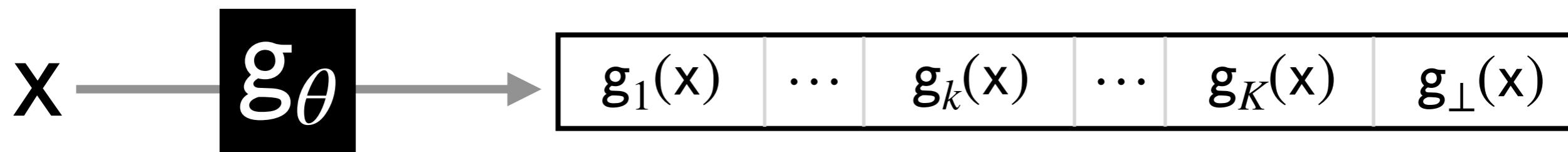


softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

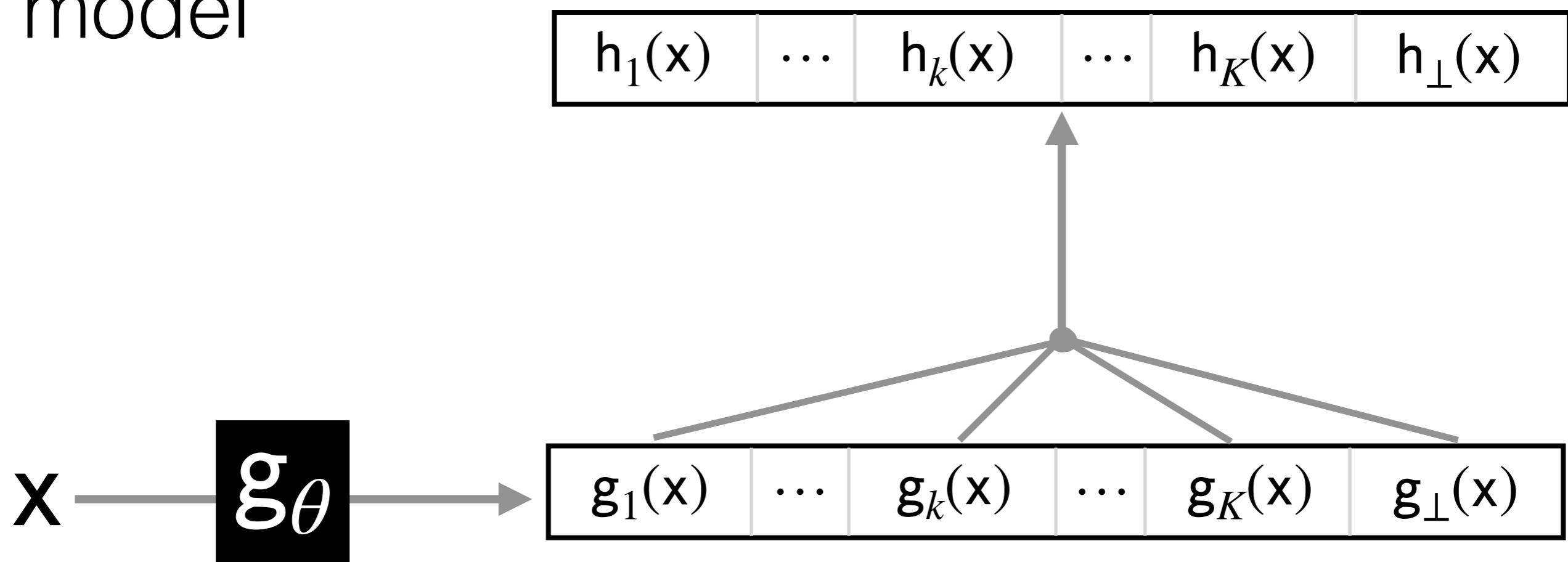


softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model



softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

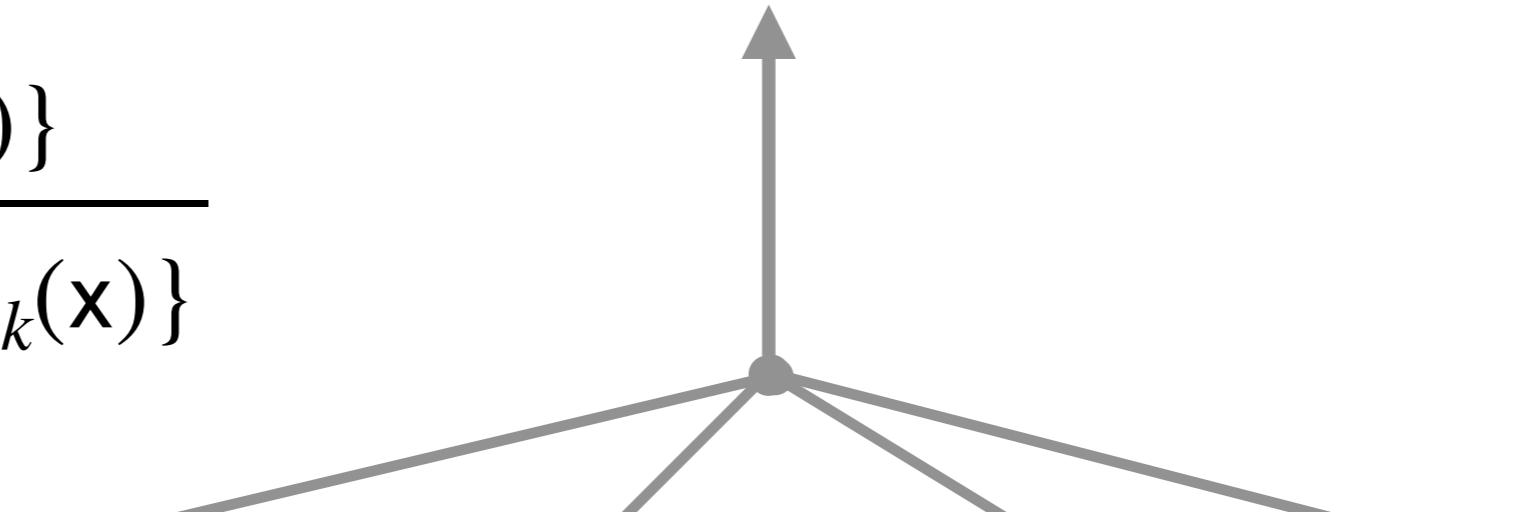
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$h_1(\mathbf{x})$	\dots	$h_k(\mathbf{x})$	\dots	$h_K(\mathbf{x})$	$h_\perp(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	-----------------------

\mathbf{x}

$$\mathbf{g}_\theta$$

$g_1(\mathbf{x})$	\dots	$g_k(\mathbf{x})$	\dots	$g_K(\mathbf{x})$	$g_\perp(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	-----------------------



softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$h_1(\mathbf{x})$	\dots	$h_k(\mathbf{x})$	\dots	$h_K(\mathbf{x})$	$h_{\perp}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	-------------------------

\mathbf{x}

$$\mathbf{g}_{\theta}$$

$g_1(\mathbf{x})$	\dots	$g_k(\mathbf{x})$	\dots	$g_K(\mathbf{x})$	$g_{\perp}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	-------------------------



softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

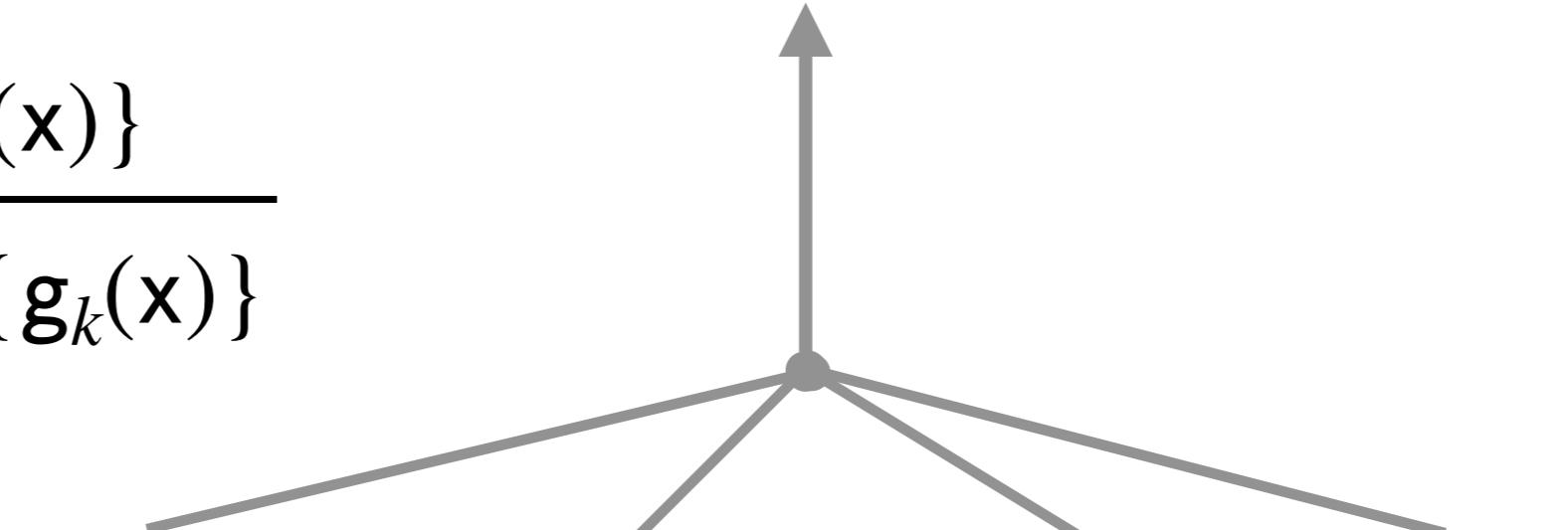
$$h_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

$h_1(x)$	\dots	$h_k(x)$	\dots	$h_K(x)$	$h_{\perp}(x)$
----------	---------	----------	---------	----------	----------------

x

g_{θ}

$g_1(x)$	\dots	$g_k(x)$	\dots	$g_K(x)$	$g_{\perp}(x)$
----------	---------	----------	---------	----------	----------------



softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \lambda) = -\log h_y(\mathbf{x}) - (1 - \lambda) \cdot \log h_\perp(\mathbf{x})$$

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \lambda) = -\log h_y(\mathbf{x}) - (1 - \lambda) \cdot \log h_\perp(\mathbf{x})$$

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \lambda) = -\log h_y(\mathbf{x}) - (1 - \lambda) \cdot \log h_\perp(\mathbf{x})$$

softmax implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n \right\}_{n=1}^N$$

model

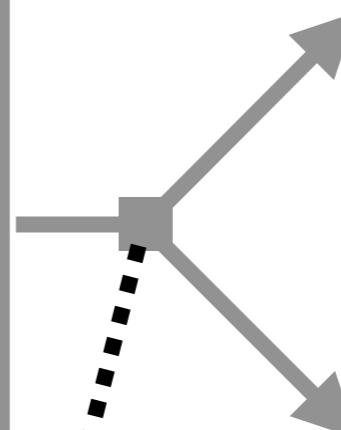
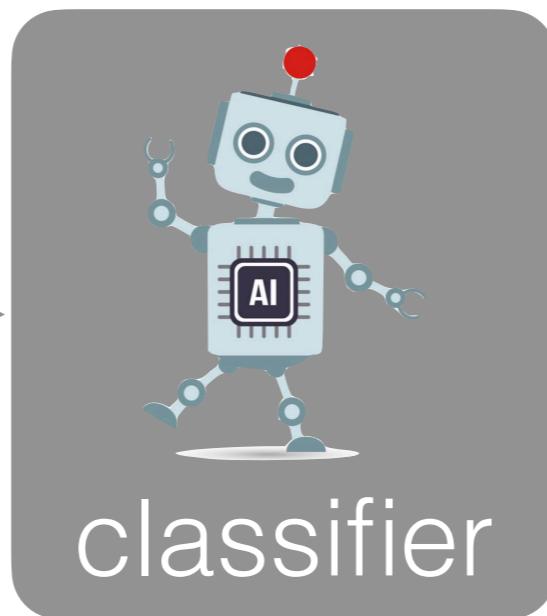
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \lambda) = -\log h_y(\mathbf{x}) - (1 - \lambda) \cdot \log h_\perp(\mathbf{x})$$

learning to reject

input
features



\hat{y}



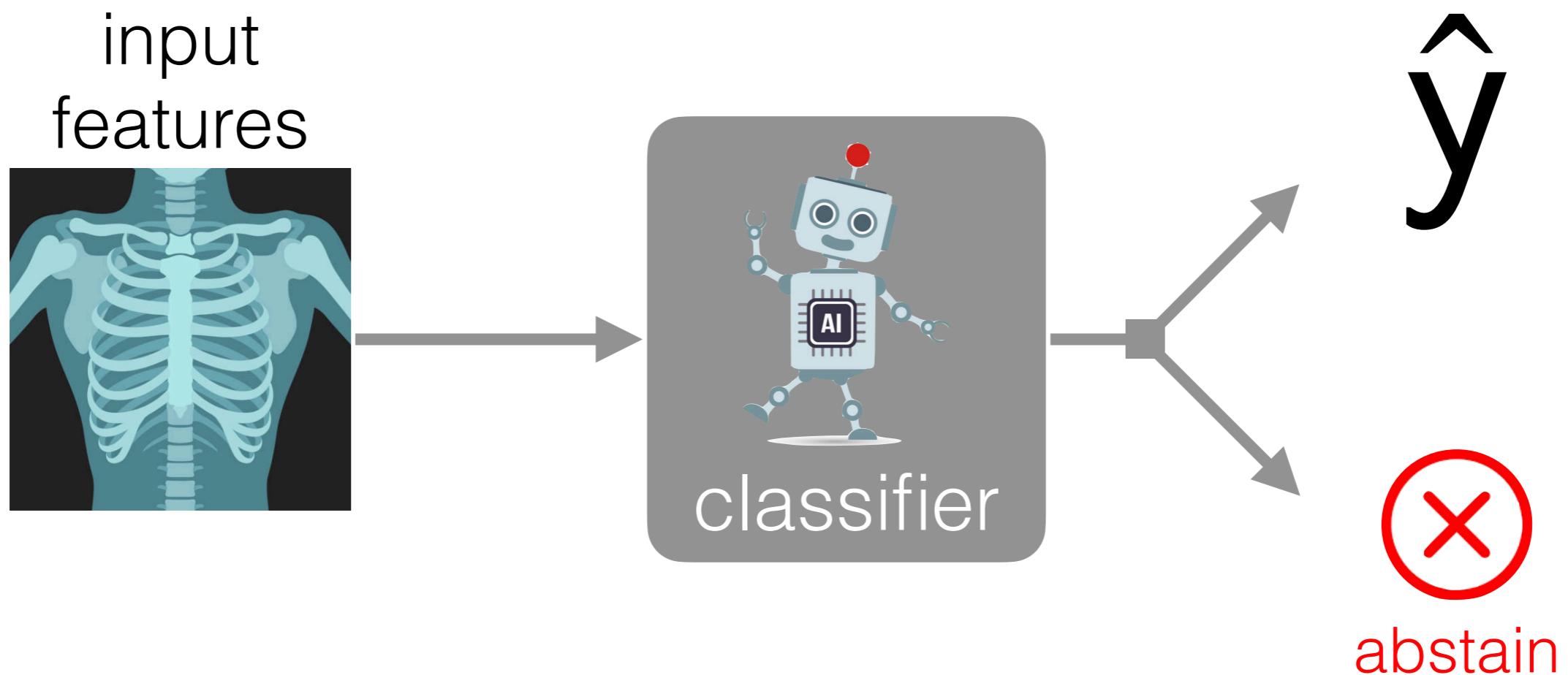
abstain

reject if...

$$\max_{y \in [1, K]} h_y(x) < h_\perp(x)$$

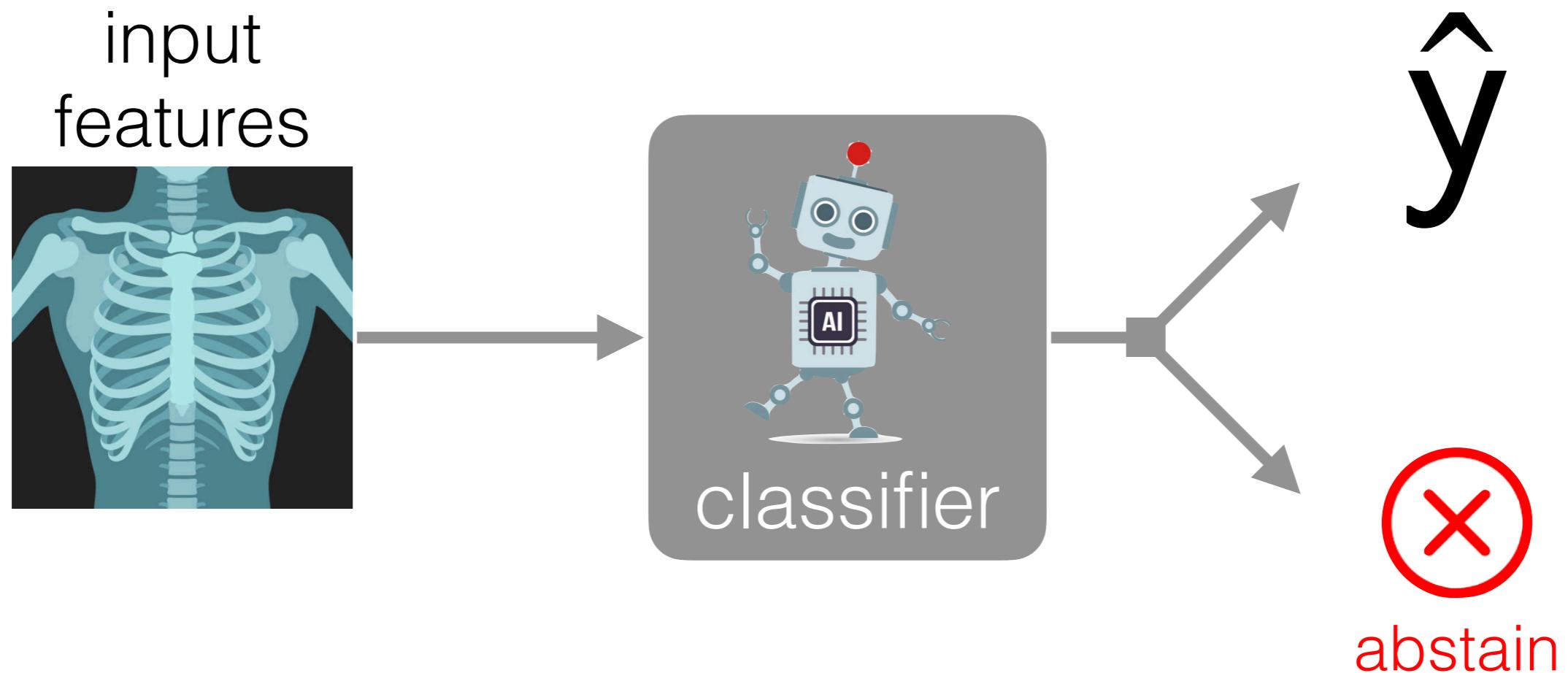
- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

learning to reject



What are the limitations of learning to reject?

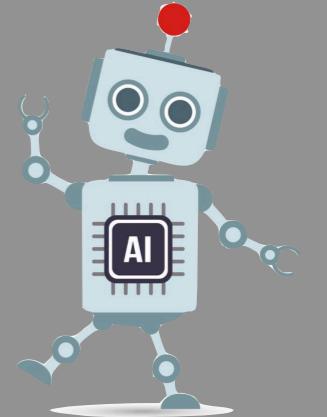
learning to reject



What are the limitations of learning to reject?

What happens after the model rejects?

input
features



classifier



expert

learning to defer (to an expert)

input
features



allocation
mechanism

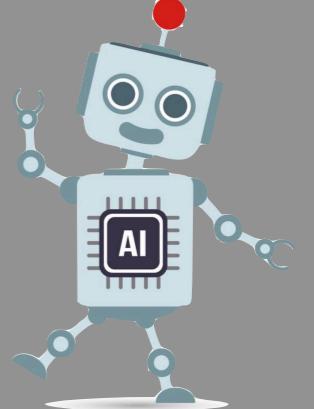
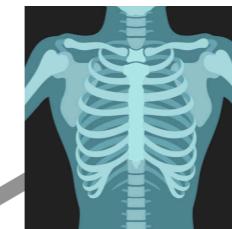


learning to defer (to an expert)

input
features



allocation
mechanism



classifier



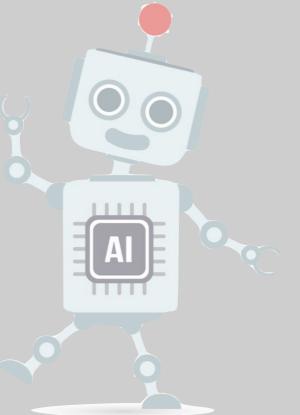
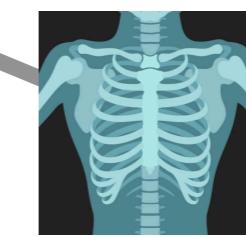
expert

learning to defer (to an expert)

input
features



allocation
mechanism



classifier



expert

learning to defer (to an expert)

input
features



allocation
mechanism



input
features

X



m
prediction

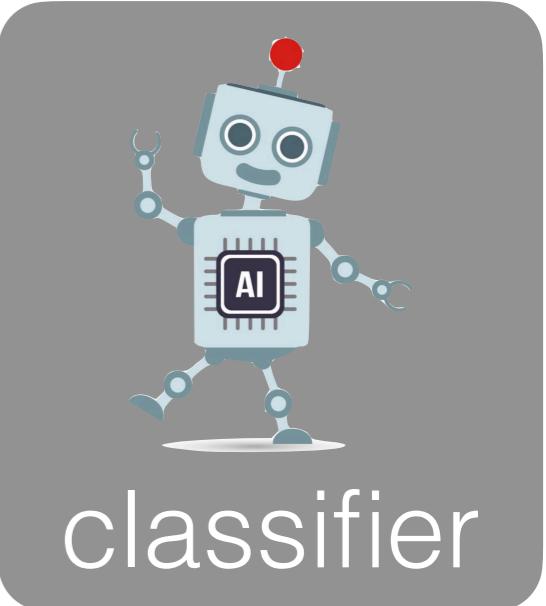
(black box)



input
features



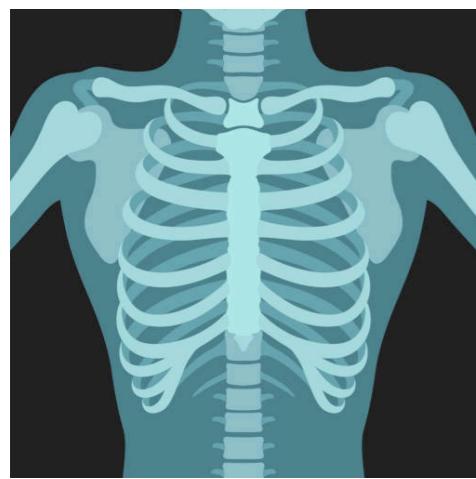
allocation
mechanism



now we will incorporate information about
what will happen if the model doesn't
predict and the human does instead.

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

input
features



allocation
mechanism

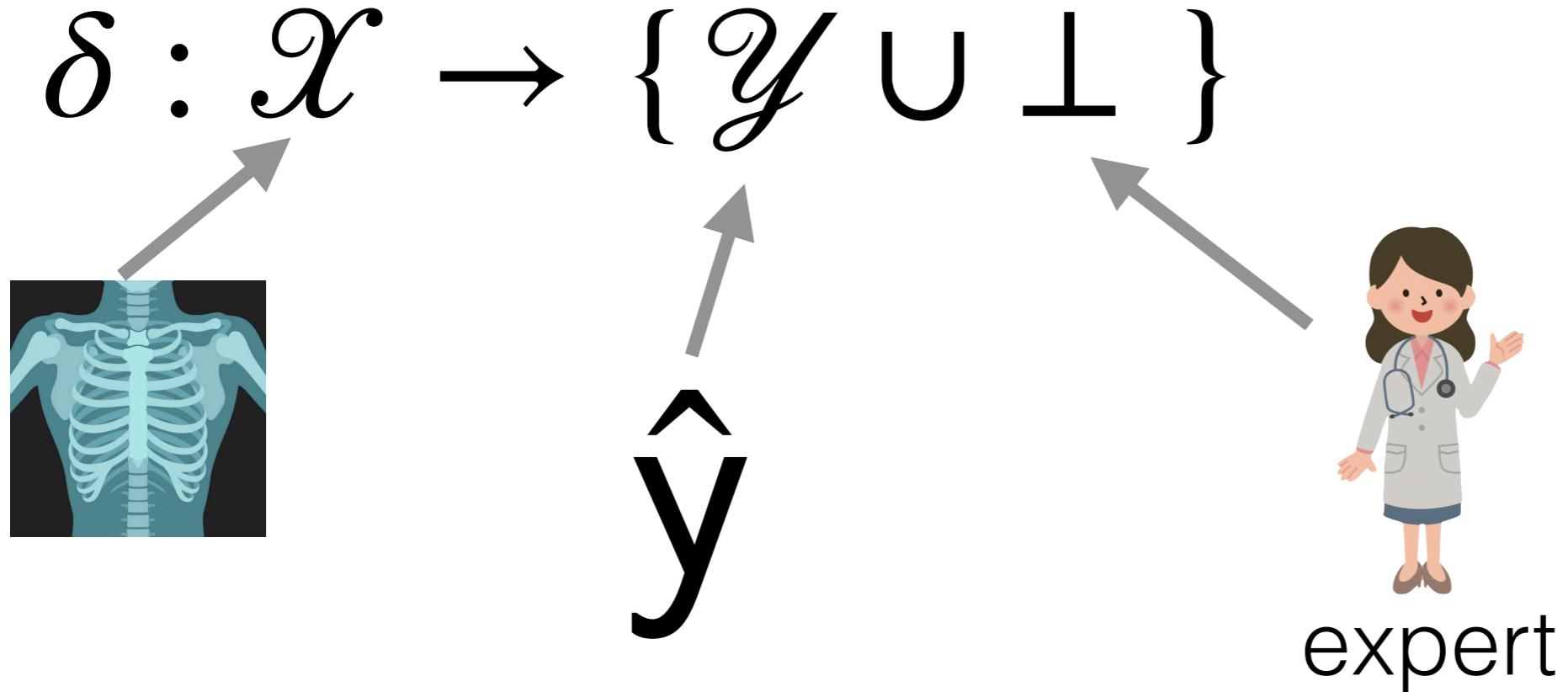
???



Derivation from Bayes decision theory

$$\delta : \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

Derivation from Bayes decision theory



Derivation from Bayes decision theory

$$\delta: \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

$$c(\delta(x), y, m) = \begin{cases} 0 & \text{if } \delta(x) = y \\ 1 & \text{if } \delta(x) \neq y \\ \mathbb{I}[m \neq y] & \text{if } \delta(x) = \perp \end{cases}$$

Derivation from Bayes decision theory

$$\delta: \mathcal{X} \rightarrow \{\mathcal{Y} \cup \perp\}$$

$$c(\delta(x), y, m) = \begin{cases} 0 & \text{if } \delta(x) = y \\ 1 & \text{if } \delta(x) \neq y \\ \mathbb{I}[m \neq y] & \text{if } \delta(x) = \perp \end{cases}$$

0-1 loss applied to
expert's prediction

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

when we choose to make a prediction...

$$\delta(x) = \hat{y} = \arg \max_y \mathbb{P}(y | x)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

when we choose to make a prediction...

$$\delta(x) = \hat{y} = \arg \max_y \mathbb{P}(y | x)$$

$$\mathfrak{R}(x, \hat{y}) = \sum_{y \neq \hat{y}} \mathbb{P}(y | x) = 1 - \mathbb{P}(\hat{y} | x)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

when we choose to defer...

$$\delta(x) = \perp$$

Derivation from Bayes decision theory

$$\mathfrak{R}(x, \delta) = \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot c(\delta(x), y, m)$$

when we choose to defer...

$$\delta(x) = \perp$$

$$\begin{aligned} \mathfrak{R}(x, \perp) &= \sum_{y \in \mathcal{Y}} \sum_{m \in \mathcal{Y}} \mathbb{P}(y | x) \cdot \mathbb{P}(m | x) \cdot \mathbb{I}[m \neq y] \\ &= \mathbb{P}_{m,y}(m \neq y | x) \end{aligned}$$

Derivation from Bayes decision theory

Thus we need to choose the action with the minimum risk:

$$\mathcal{R}(x, \hat{y}) = 1 - P(\hat{y} | x)$$

vs

$$\mathcal{R}(x, \perp) = P_{m,y}(m \neq y | x)$$

Derivation from Bayes decision theory

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y|x) & \text{if } 1 - \max_y \mathbb{P}(y|x) \leq \mathbb{P}_{m,y}(m \neq y|x) \\ \perp & \text{otherwise} \end{cases}$$

Derivation from Bayes decision theory

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y|x) & \text{if } 1 - \max_y \mathbb{P}(y|x) \leq \mathbb{P}_{m,y}(m \neq y|x) \\ \perp & \text{otherwise} \end{cases}$$

rearranging a bit...

$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y|x) & \text{if } \max_y \mathbb{P}(y|x) \geq \mathbb{P}_{m,y}(m = y|x) \\ \perp & \text{otherwise} \end{cases}$$

input
features



allocation
mechanism



$$\delta^*(x) = \begin{cases} \arg \max_y \mathbb{P}(y | x) & \text{if } \max_y \mathbb{P}(y | x) \geq P_{m,y}(m = y | x) \\ \perp & \text{otherwise} \end{cases}$$

input
features



allocation
mechanism

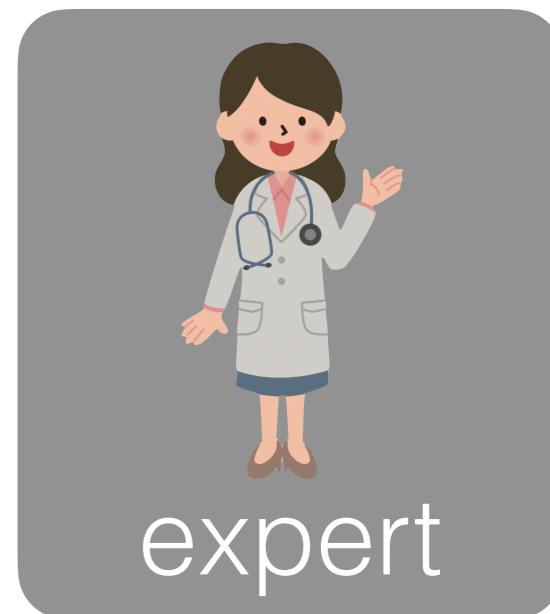


$$\hat{\delta}(x) = \begin{cases} \arg \max_y p(y | x) & \text{if } \max_y p(y | x) \geq p(m = y | x) \\ \perp & \text{otherwise} \end{cases}$$

input
features



allocation
mechanism



Any Problems?

$$\hat{\delta}(x) = \begin{cases} \arg \max_y p(y | x) & \text{if } \max_y p(y | x) \geq p(m = y | x) \\ \perp & \text{otherwise} \end{cases}$$

input
features



allocation
mechanism



Any Problems?

$$\hat{\delta}(x) = \begin{cases} \arg \max_y p(y | x) & \text{if } \max_y p(y | x) \geq p(m = y | x) \\ \perp & \text{otherwise} \end{cases}$$

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ **consistent surrogate loss**
 - ⊗ extension to multiple experts

softmax implementation

[Mozannar & Sontag, 2020]

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

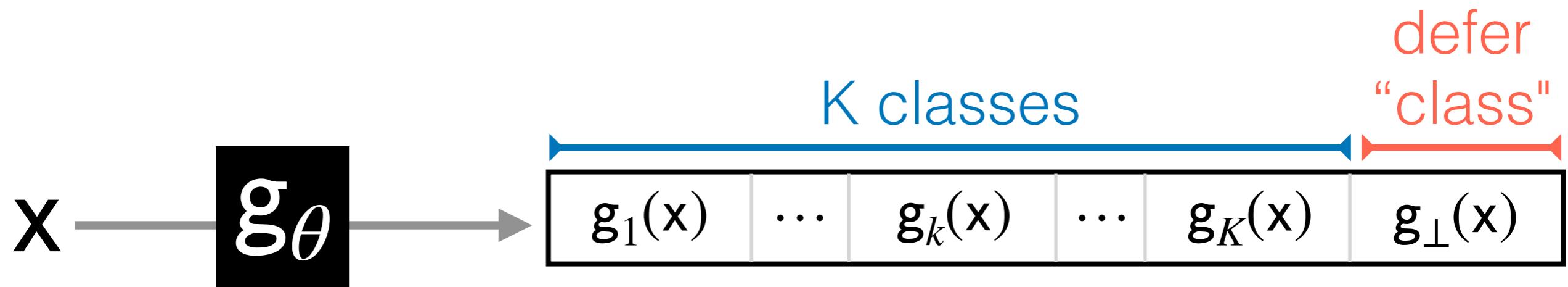
softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model



softmax implementation

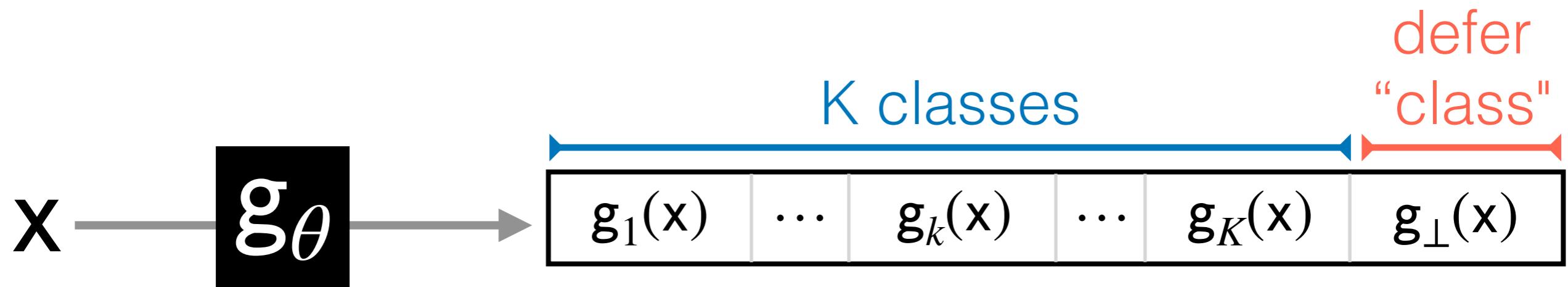
[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$g_k(\mathbf{x}) \in \mathbb{R}$$



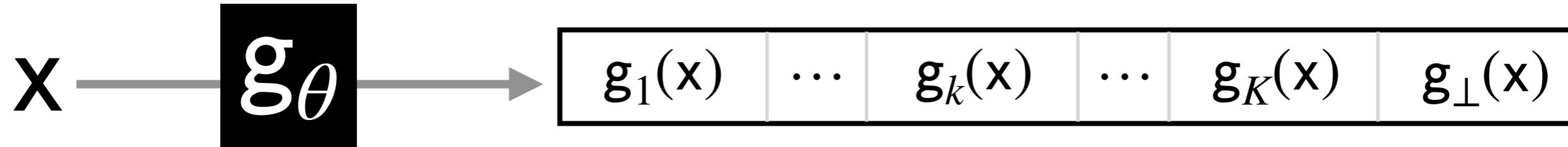
softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model



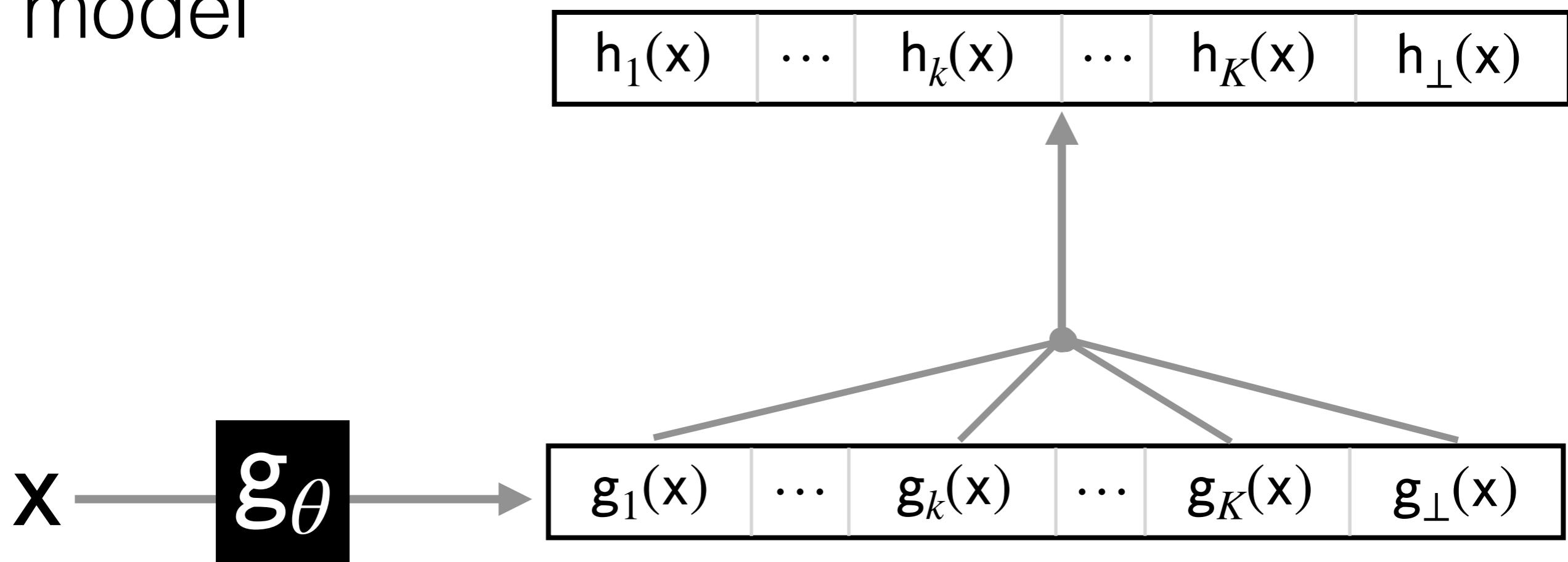
softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model



softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

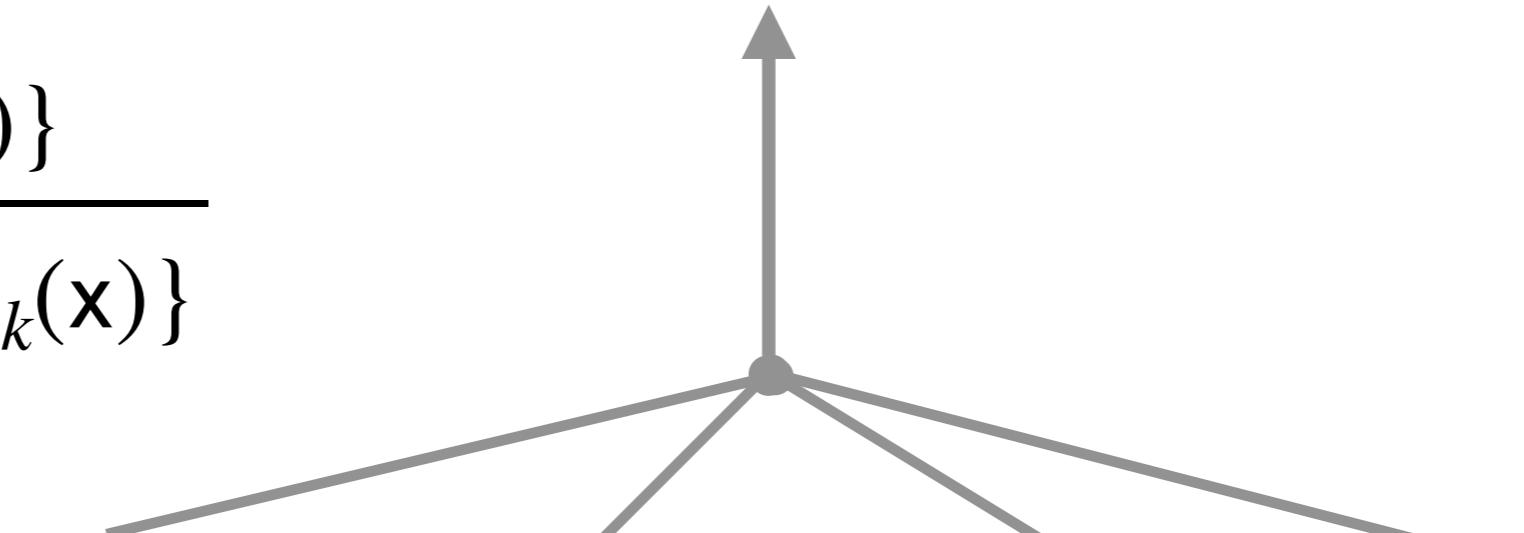
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$$\begin{array}{c} h_1(\mathbf{x}) \quad \cdots \quad h_k(\mathbf{x}) \quad \cdots \quad h_K(\mathbf{x}) \quad h_\perp(\mathbf{x}) \\ \hline \end{array}$$

\mathbf{x}

\mathbf{g}_θ

$$\begin{array}{c} g_1(\mathbf{x}) \quad \cdots \quad g_k(\mathbf{x}) \quad \cdots \quad g_K(\mathbf{x}) \quad g_\perp(\mathbf{x}) \\ \hline \end{array}$$



softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

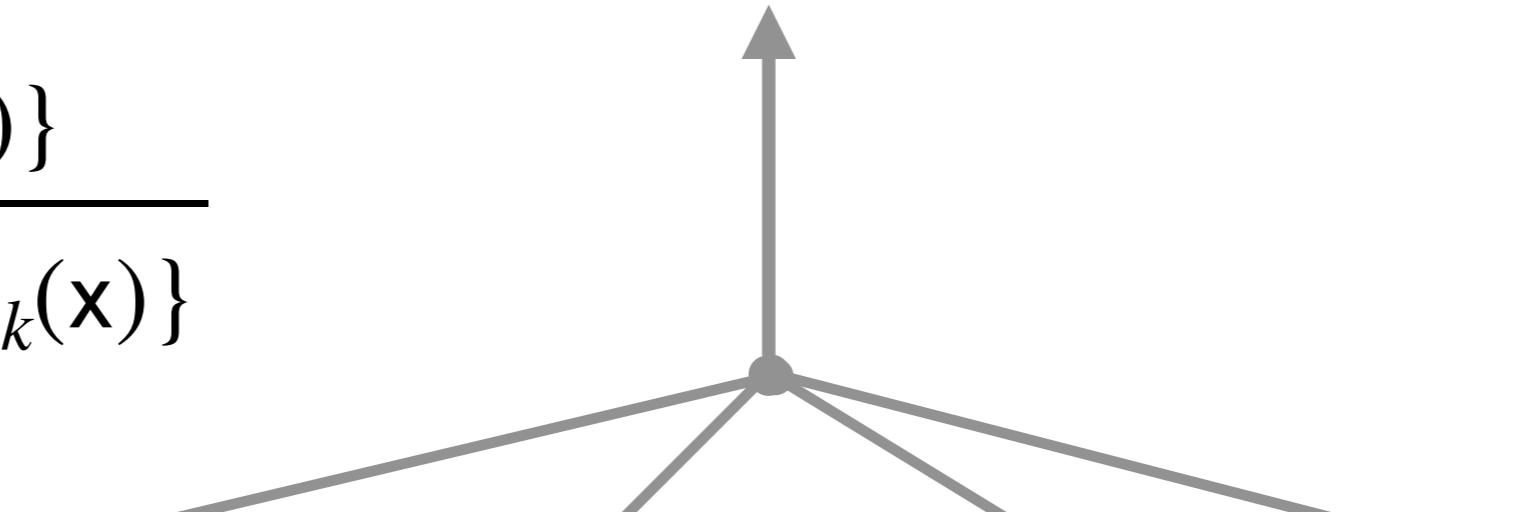
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$$\begin{array}{c} h_1(\mathbf{x}) \quad \cdots \quad h_k(\mathbf{x}) \quad \cdots \quad h_K(\mathbf{x}) \quad h_{\perp}(\mathbf{x}) \\ \hline \end{array}$$

\mathbf{x}

\mathbf{g}_{θ}

$$\begin{array}{c} g_1(\mathbf{x}) \quad \cdots \quad g_k(\mathbf{x}) \quad \cdots \quad g_K(\mathbf{x}) \quad g_{\perp}(\mathbf{x}) \\ \hline \end{array}$$



softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

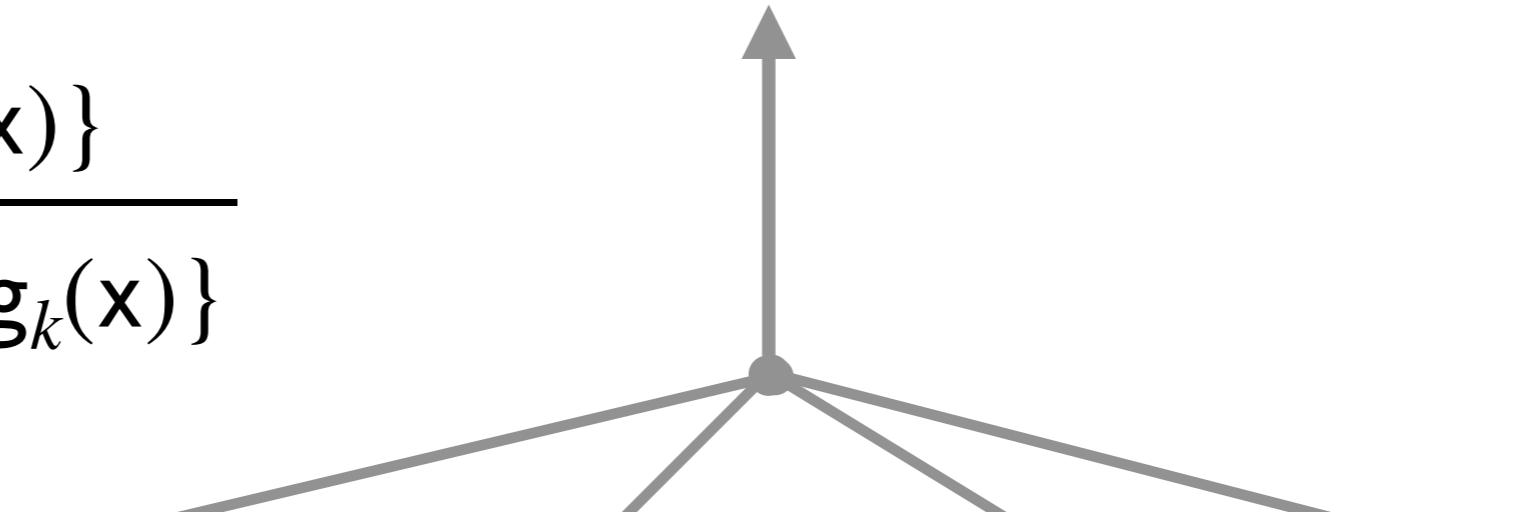
$$h_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

$h_1(x)$	\dots	$h_k(x)$	\dots	$h_K(x)$	$h_{\perp}(x)$
----------	---------	----------	---------	----------	----------------

x

g_{θ}

$g_1(x)$	\dots	$g_k(x)$	\dots	$g_K(x)$	$g_{\perp}(x)$
----------	---------	----------	---------	----------	----------------



softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = m] \cdot \log h_m(\mathbf{x})$$

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = m] \cdot \log h_\perp(\mathbf{x})$$

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = m] \cdot \log h_\perp(\mathbf{x})$$

softmax implementation

[Mozannar & Sontag, 2020]

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = m] \cdot \log h_\perp(\mathbf{x})$$

input
features



allocation
mechanism

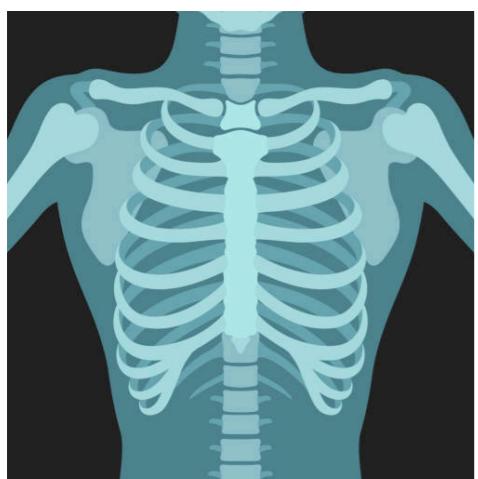


defer to expert if...

$$\max_{y \in [1, K]} h_y(x) < h_\perp(x)$$

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

input
features



allocation
mechanism

classifier

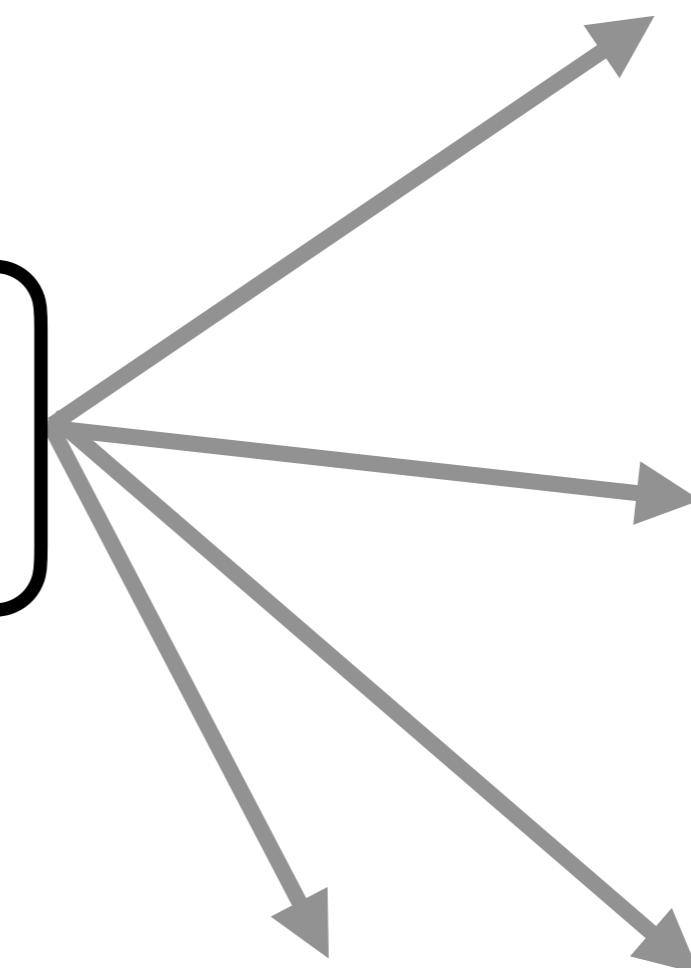


expert

input
features



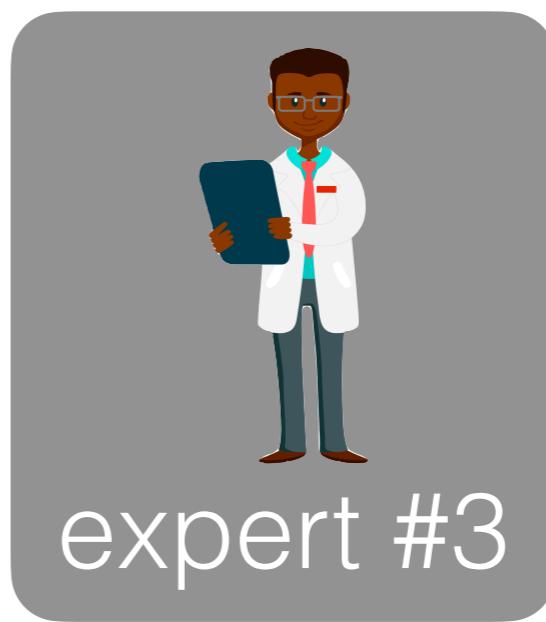
allocation
mechanism



classifier



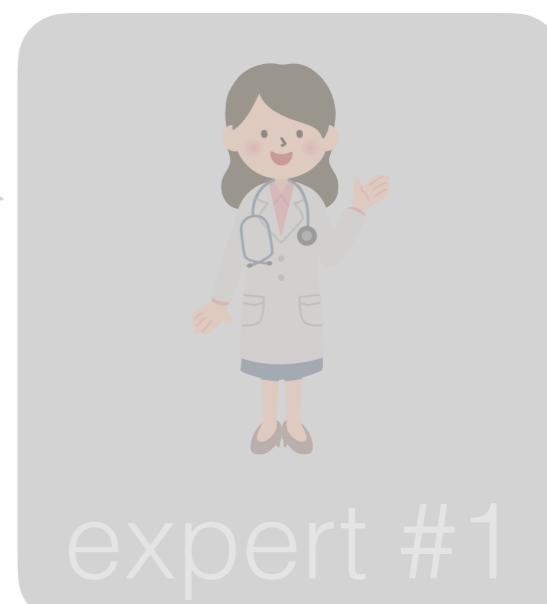
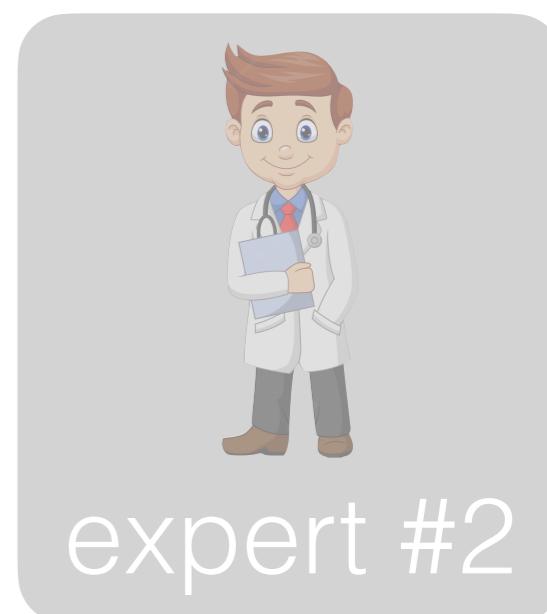
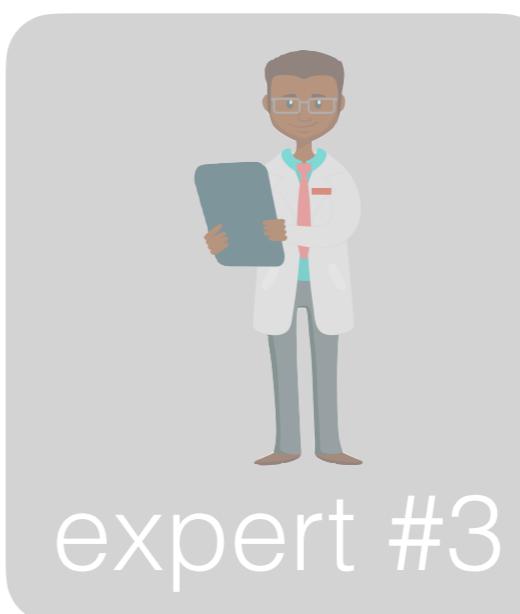
expert #1



input
features



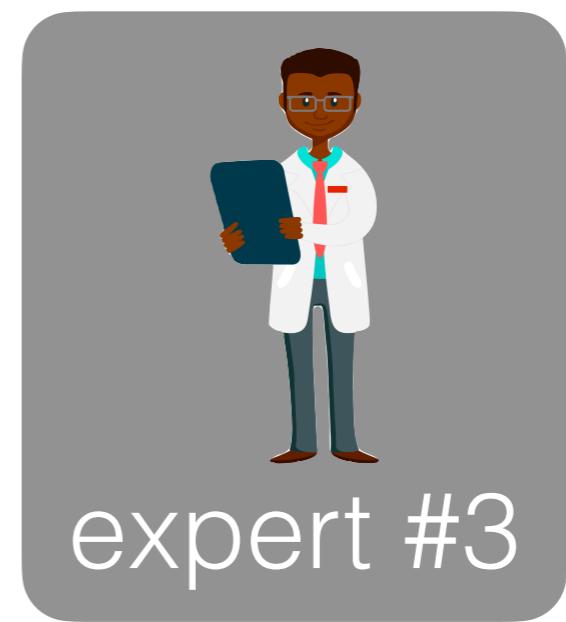
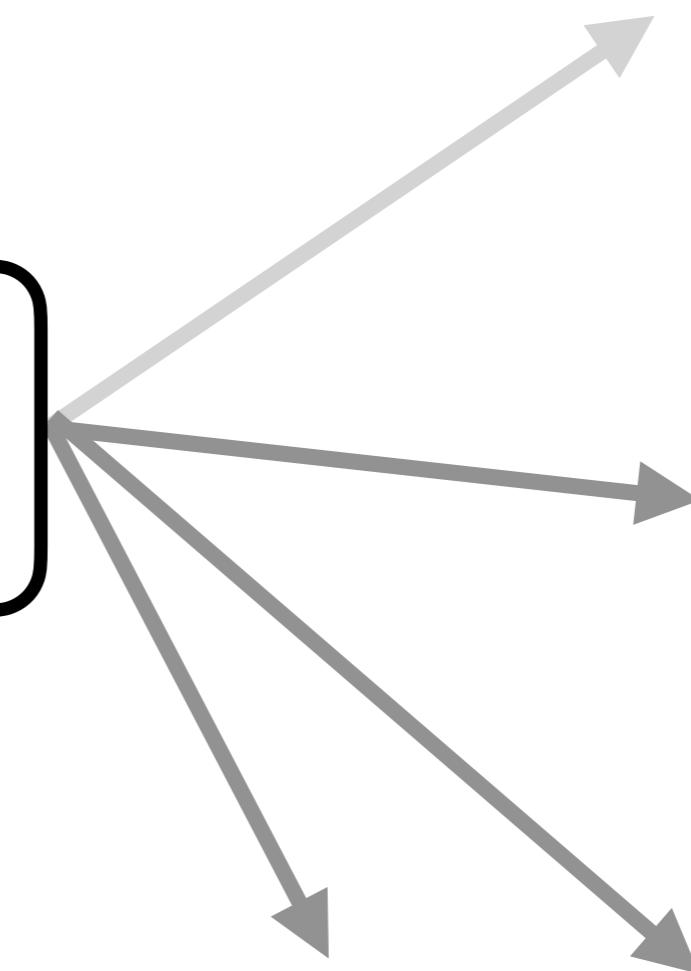
allocation
mechanism



input
features



allocation
mechanism

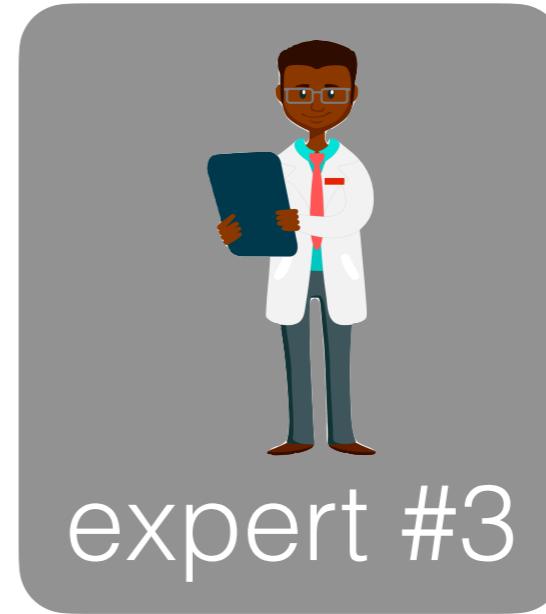


input
features



allocation
mechanism

expert #3



expert #2

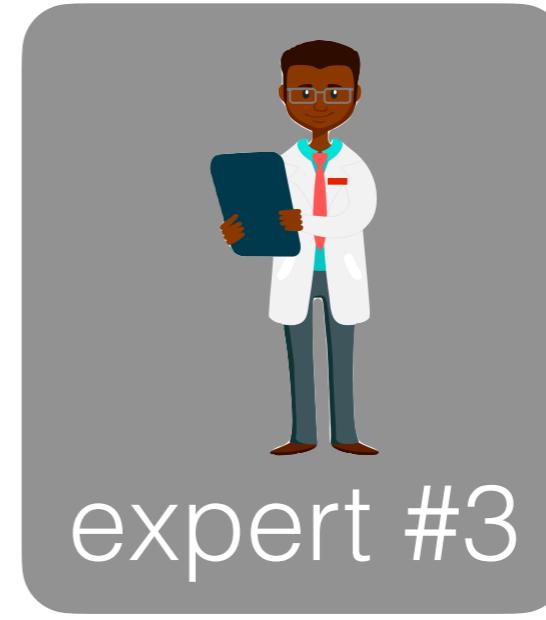


input
features



allocation
mechanism

expert #3



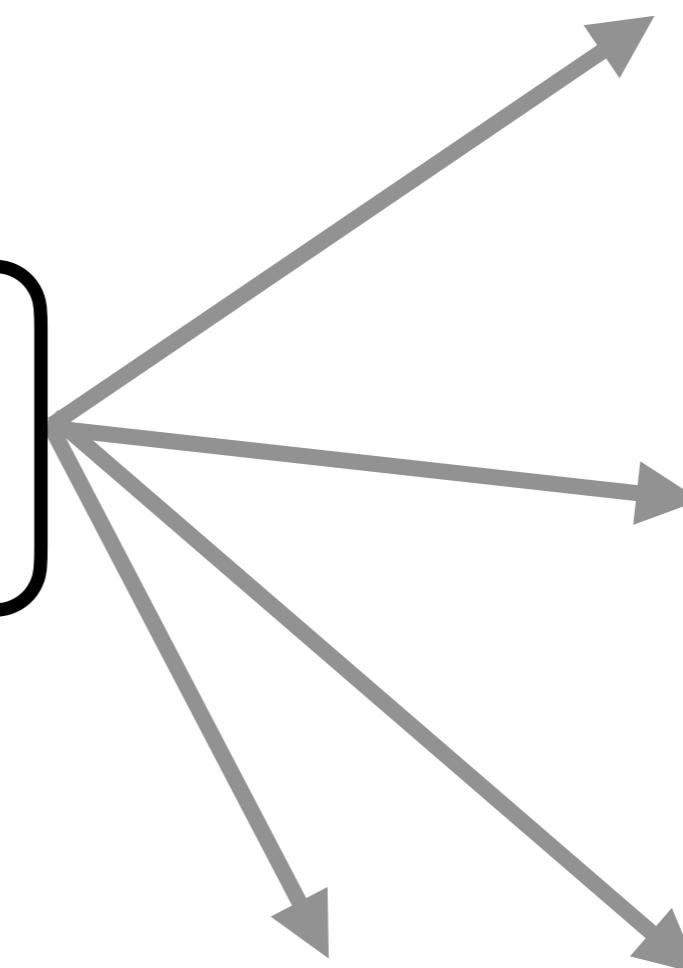
expert #2



input
features



allocation
mechanism

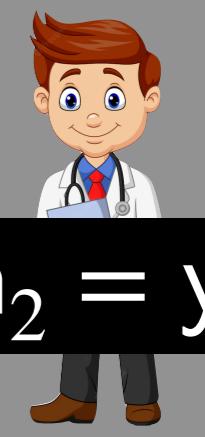


classifier



$$P(m_1 = y | x)$$

expert #1



$$P(m_3 = y | x)$$

expert #3

$$P(m_2 = y | x)$$

expert #2



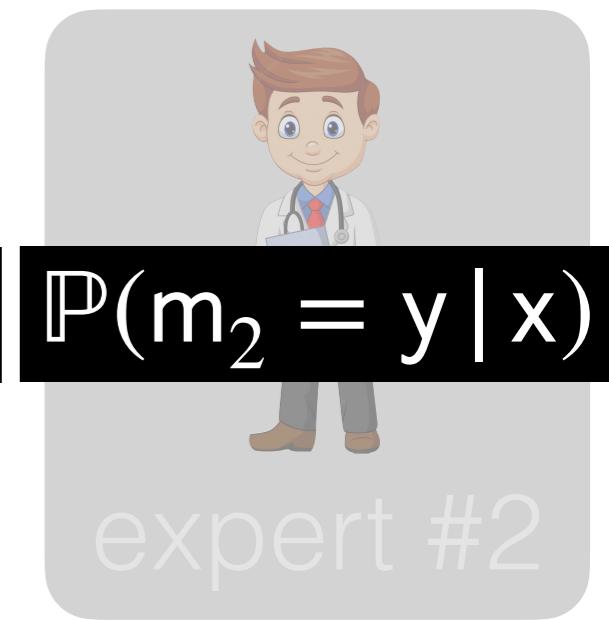
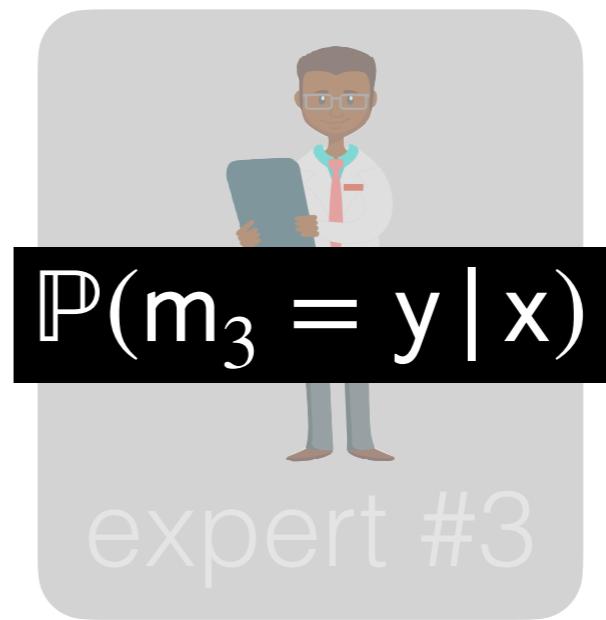
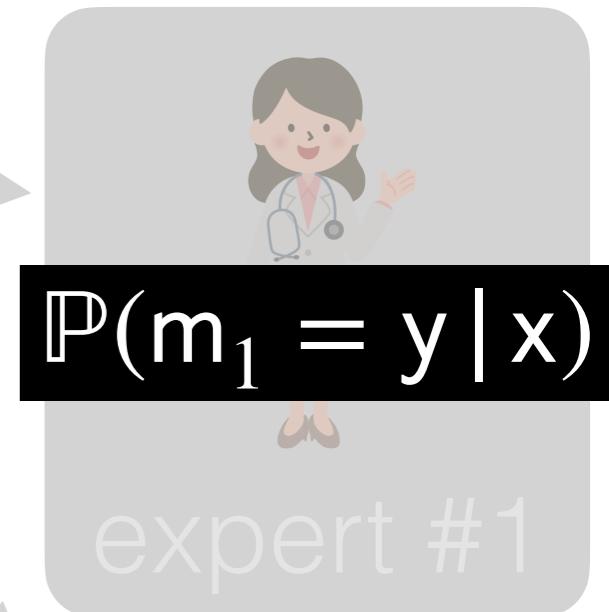
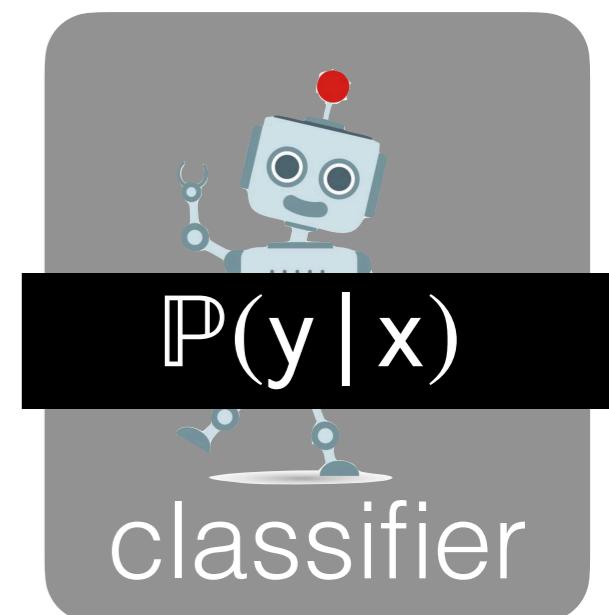
input
features



allocation
mechanism

use classifier if...

$$\max_y \mathbb{P}(y|x) > \mathbb{P}(m_j = y|x), \forall j$$



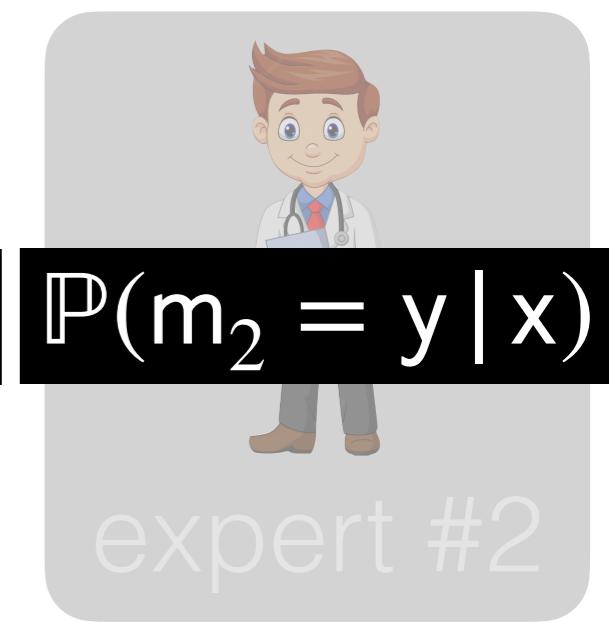
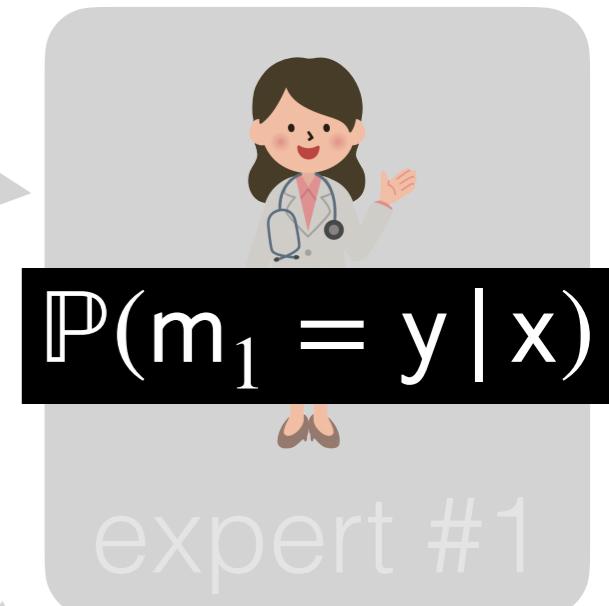
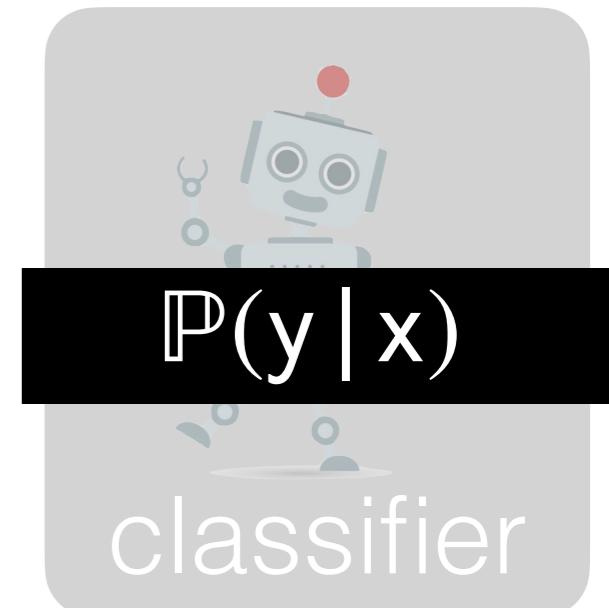
input
features



allocation
mechanism

else, pick best expert:

$$\arg \max_j \mathbb{P}(m_j = y | x)$$



multi-expert implementation

training data

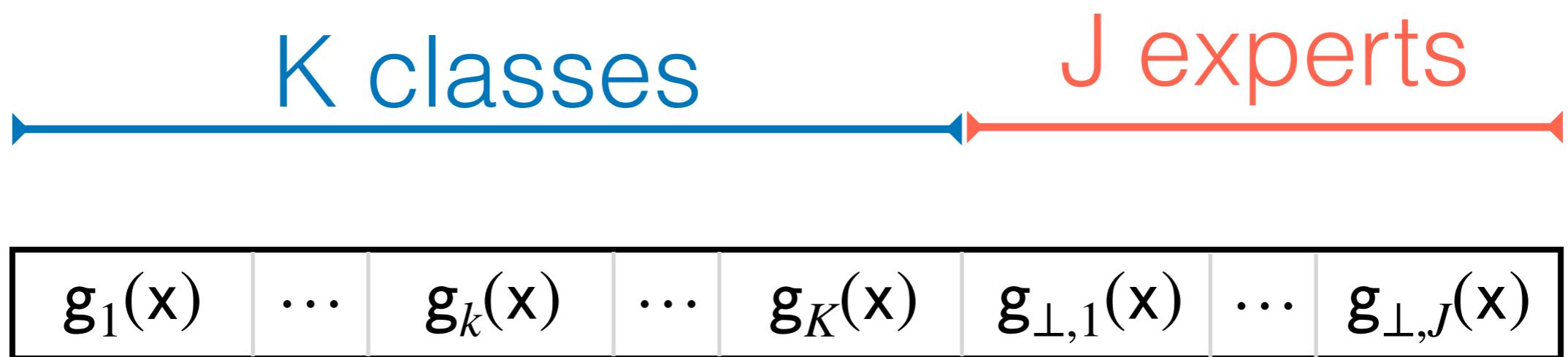
$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model



multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$h_1(\mathbf{x})$	\dots	$h_k(\mathbf{x})$	\dots	$h_K(\mathbf{x})$	$h_{\perp,1}(\mathbf{x})$	\dots	$h_{\perp,J}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	---------------------------	---------	---------------------------

K classes

J experts

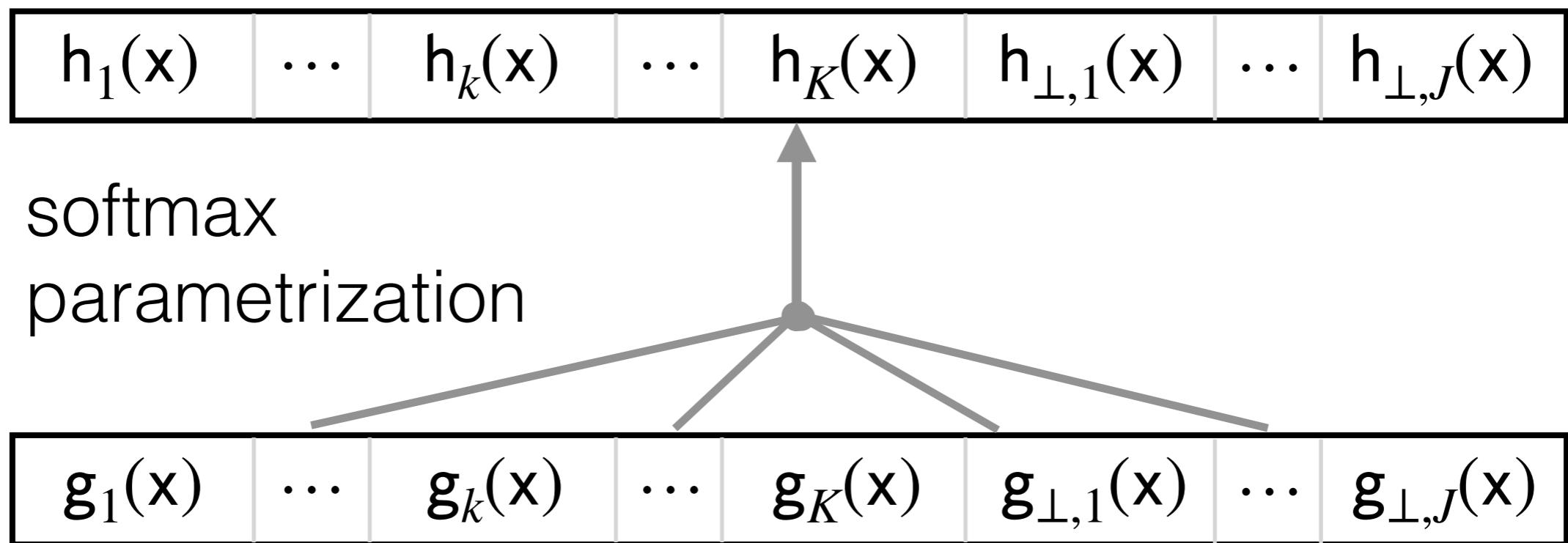
$g_1(\mathbf{x})$	\dots	$g_k(\mathbf{x})$	\dots	$g_K(\mathbf{x})$	$g_{\perp,1}(\mathbf{x})$	\dots	$g_{\perp,J}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	---------------------------	---------	---------------------------

multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model



multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+J} \exp\{g_k(\mathbf{x})\}}$$

multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+J} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \{\mathbf{m}_j\}_{j=1}^J) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = m_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+J} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \{\mathbf{m}_j\}_{j=1}^J) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = m_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

multi-expert implementation

training data

$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+J} \exp\{g_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, \mathbf{y}, \{\mathbf{m}_j\}_{j=1}^J) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = m_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

- ⊗ Learning to reject
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
- ⊗ Learning to defer to an expert
 - ⊗ Bayes decision theory
 - ⊗ consistent surrogate loss
 - ⊗ extension to multiple experts

