

DSC190: Machine Learning with Few Labels

Weak/distant supervision
A “standardized” view of ML

Zhiting Hu

Lecture 8, October 19, 2021

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

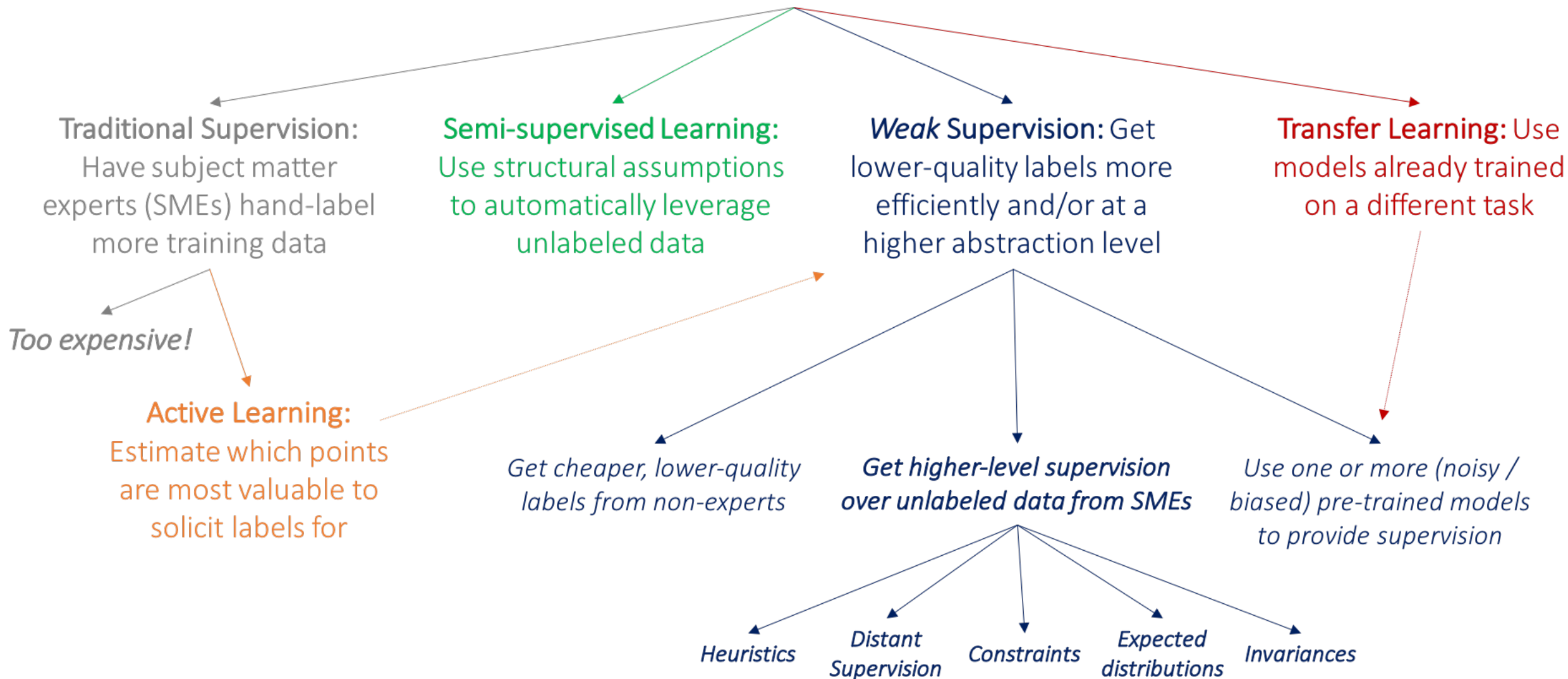
Outline

- Distant supervision
- A “standardized” view of ML

The difficulty with supervised learning

- Annotated data is expensive and costs increase when...
 - *A task requires specialized expertise*
E.g. "Only a trained linguist or a board certified radiologist can label my data"
 - *Labeling examples involves making multiple decisions*
E.g. "Annotate this sentence with a parse tree"
(instead of a single binary decision)

How to get more labeled training data?



Example (I): labeling with heuristics

Task: Build a chest x-ray classifier
(normal/abnormal)



Indication: Chest pain. Findings: Mediastinal contours are within **normal** limits. Heart size is within **normal** limits. **No** focal consolidation, **pneumothorax** or **pleural effusion**. Impression: **No** acute cardiopulmonary abnormality.

Can you use the accompanying medical report (text modality) to label the x-ray (image modality)?

Example (I): labeling with heuristics

Indication: Chest pain. Findings: Mediastinal contours are within **normal** limits. Heart size is within **normal** limits. **No** focal consolidation, **pneumothorax** or **pleural effusion**. Impression: **No** acute cardiopulmonary abnormality.



How do we obtain Y?

Y

CNN

Example (I): labeling with heuristics

Indication: Chest pain. Findings: Mediastinal contours are within **normal** limits. Heart size is within **normal** limits. **No** focal consolidation, **pneumothorax** or **pleural effusion**. Impression: **No** acute cardiopulmonary abnormality.

Normal Report

```
def LF_pneumothorax(c):
    if re.search(r'pneumo.*', c.report.text):
        return "ABNORMAL"

def LF_pleural_effusion(c):
    if "pleural effusion" in c.report.text:
        return "ABNORMAL"

def LF_normal_report(c, thresh=2):
    if len(NORMAL_TERMS.intersection(c.
report.words)) > thresh:
        return "NORMAL"
```

LFs

(labeling functions)

Source: Khandwala et. al 2017, Cross Modal Data Programming for Medical Images

Example (II): Labeling with knowledge bases

Task: relation extraction from text

- Hypothesis: If two entities belong to a certain relation, any sentence containing those two entities is likely to express that relation
- Key idea: use a *knowledge base* of relations to get lots of *noisy* training examples

Example (II): Labeling with knowledge bases

Frequent Freebase relations

Relation name	Size	Example
/people/person/nationality	281,107	John Dugard, South Africa
/location/location/contains	253,223	Belgium, Nijlen
/people/person/profession	208,888	Dusa McDuff, Mathematician
/people/person/place_of_birth	105,799	Edwin Hubble, Marshfield
/dining/restaurant/cuisine	86,213	MacAyo's Mexican Kitchen, Mexican
/business/business_chain/location	66,529	Apple Inc., Apple Inc., South Park, NC
/biology/organism_classification_rank	42,806	Scorpaeniformes, Order
/film/film/genre	40,658	Where the Sidewalk Ends, Film noir
/film/film/language	31,103	Enter the Phoenix, Cantonese
/biology/organism_higher_classification	30,052	Calopteryx, Calopterygidae
/film/film/country	27,217	Turtle Diary, United States
/film/writer/film	23,856	Irving Shulman, Rebel Without a Cause
/film/director/film	23,539	Michael Mann, Collateral
/film/producer/film	22,079	Diane Eskenazi, Aladdin
/people/deceased_person/place_of_death	18,814	John W. Kern, Asheville
/music/artist/origin	18,619	The Octopus Project, Austin
/people/person/religion	17,582	Joseph Chartrand, Catholicism
/book/author/works_written	17,278	Paul Auster, Travels in the Scriptorium
/soccer/football_position/players	17,244	Midfielder, Chen Tao
/people/deceased_person/cause_of_death	16,709	Richard Daintree, Tuberculosis
/book/book/genre	16,431	Pony Soldiers, Science fiction
/film/film/music	14,070	Stavisky, Stephen Sondheim
/business/company/industry	13,805	ATS Medical, Health care

Example (II): Labeling with knowledge bases

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Training data



Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Example (II): Labeling with knowledge bases

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y

Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Example (II): Labeling with knowledge bases

Corpus text

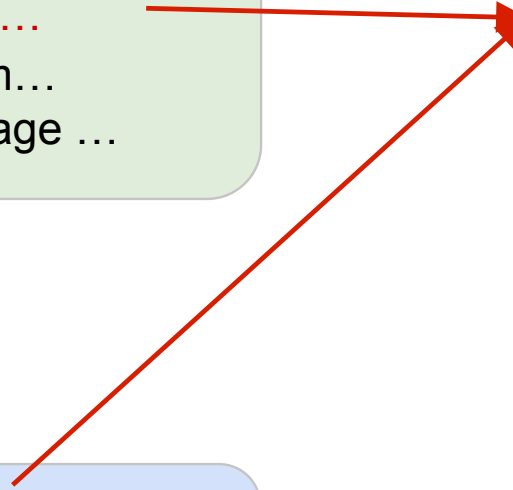
Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)



Example (II): Labeling with knowledge bases

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

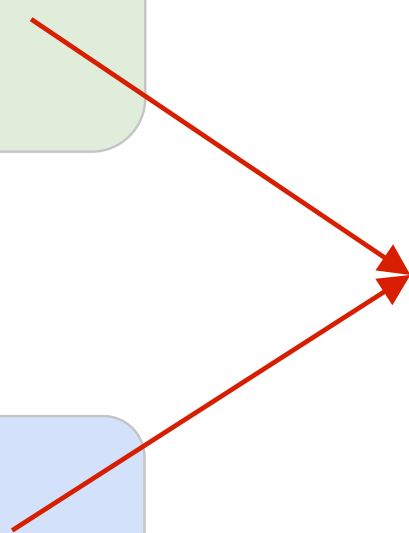
Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y



Example (II): Labeling with knowledge bases

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Freebase

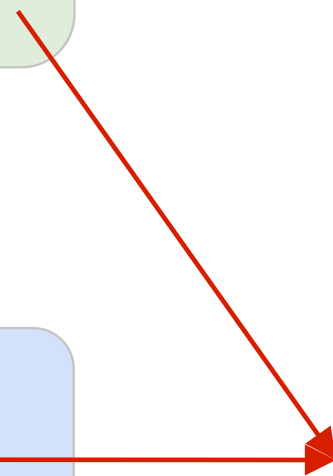
Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y

(Larry Page, Google)
Label: Founder
Feature: Y was founded by X



Example (II): Labeling with knowledge bases

Negative training data

Can't train a classifier with only positive data!
Need negative training data too!

Solution?
Sample 1% of unrelated pairs of entities.

Corpus text

Larry Page took a swipe at Microsoft...
...after Harvard invited Larry Page to...
Google is Bill Gates' worst fear ...

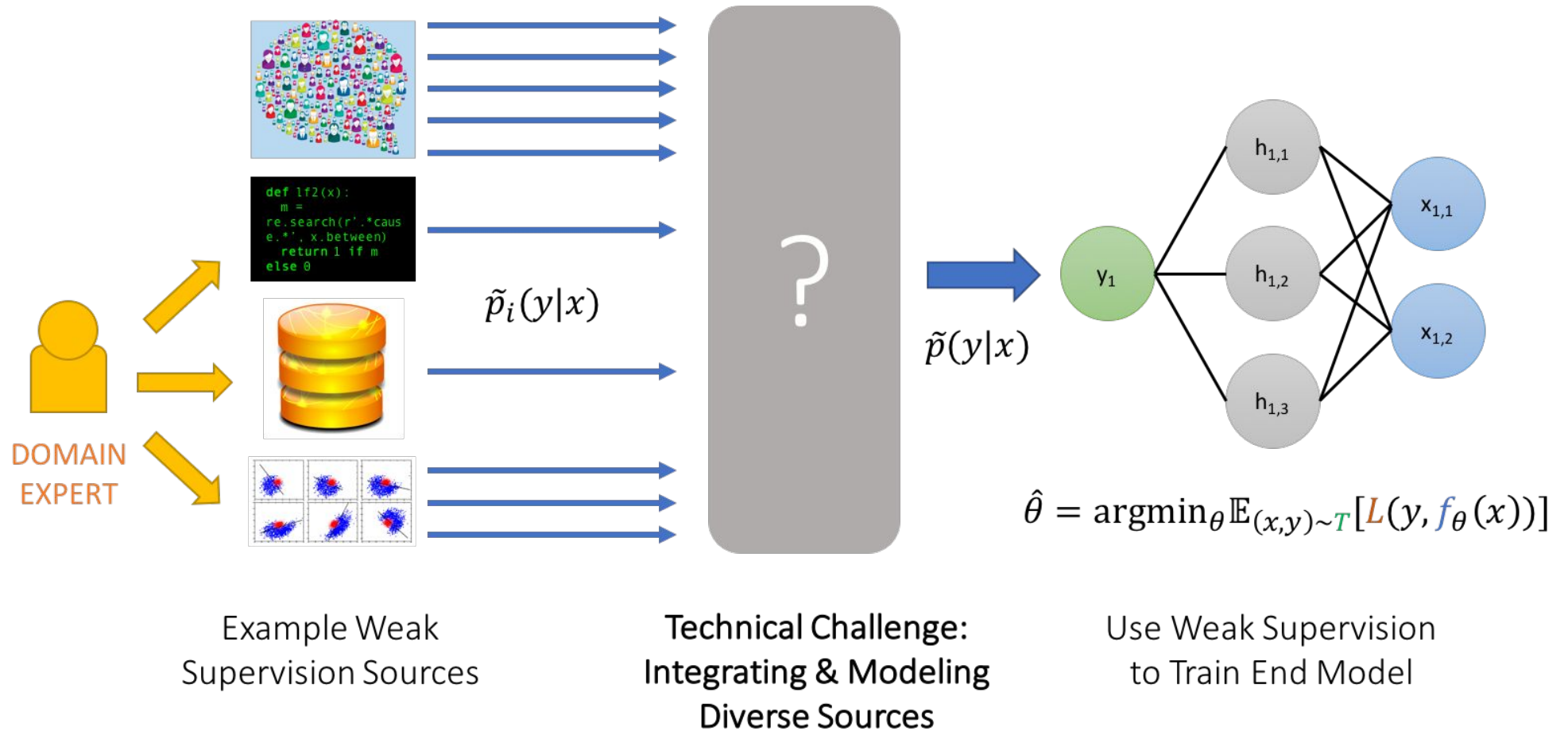
Training data

(Larry Page, Microsoft)
Label: NO_RELATION
Feature: X took a swipe at Y

(Larry Page, Harvard)
Label: NO_RELATION
Feature: Y invited X

(Bill Gates, Google)
Label: NO_RELATION
Feature: Y is X's worst fear

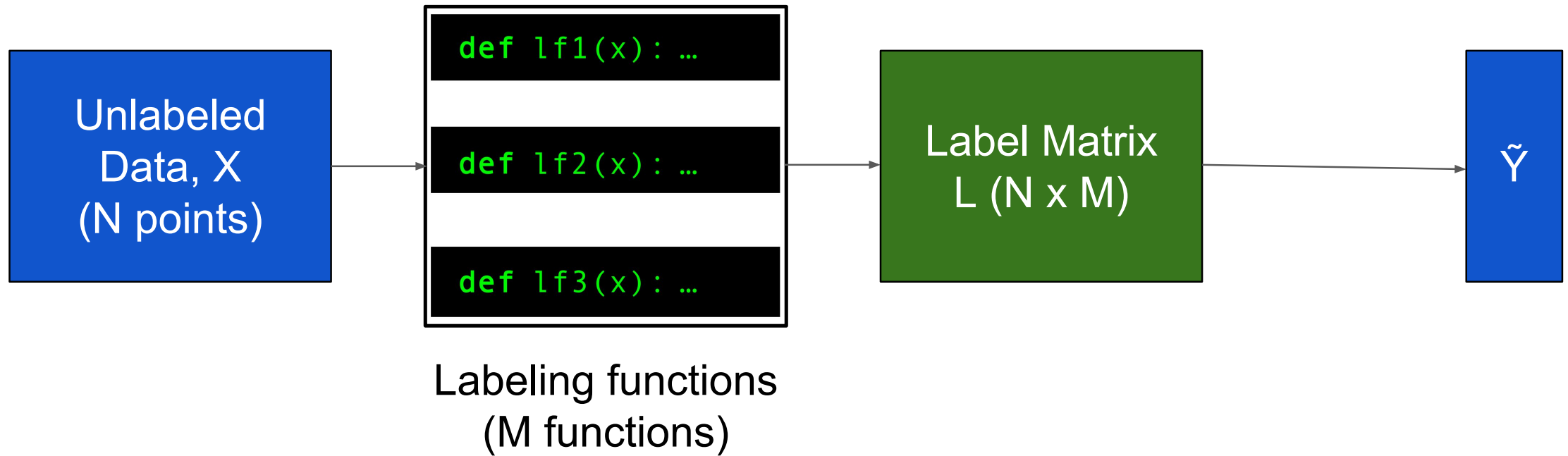
Integrating multiple noisy labels



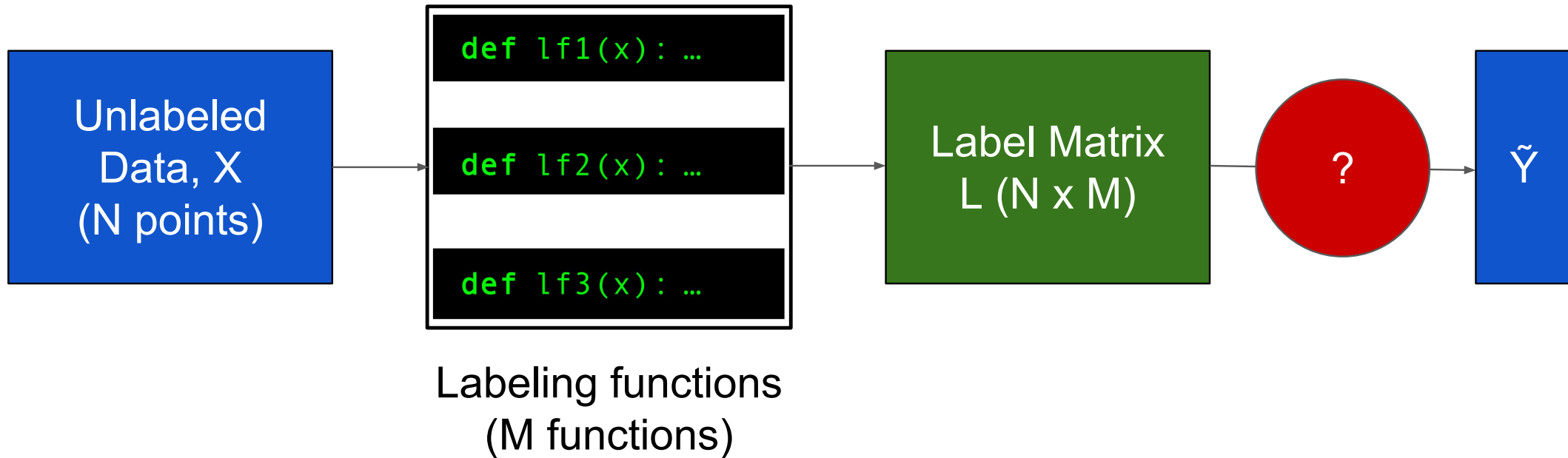
Source: A. Ratner et. al <https://dawn.cs.stanford.edu/2017/07/16/weak-supervision/>

[Credit: http://cs231n.stanford.edu/slides/2018/cs231n_2018_ds07.pdf]

Integrating multiple noisy labels



Integrating multiple noisy labels



Integrating multiple noisy labels

How do we obtain probabilistic labels, \tilde{Y} , from the label matrix, L ?

Approach 1 - Majority Vote

Take the majority vote of the labelling functions (LFs).

Let's say $L = [[0, 1, 0, 1, 0]; [1, 1, 1, 1, 0]]$.

$$\tilde{Y} = [0, 1]$$

Integrating multiple noisy labels

How do we obtain probabilistic labels, \tilde{Y} , from the label matrix, L ?

Approach 1 - Majority Vote

Indication: Chest pain. Findings: Mediastinal contours are within **normal** limits. Heart size is within **normal** limits. **No** focal consolidation, **pneumothorax** or **pleural effusion**. Impression: **No** acute cardiopulmonary abnormality.

Normal Report

Majority vote fails:

```
def LF_pneumothorax(c):  
    if re.search(r'pneumo.*', c.report.text):  
        return "ABNORMAL"  
  
def LF_pleural_effusion(c):  
    if "pleural effusion" in c.report.text:  
        return "ABNORMAL"  
  
def LF_normal_report(c, thresh=2):  
    if len(NORMAL_TERMS.intersection(c.  
report.words)) > thresh:  
        return "NORMAL"
```

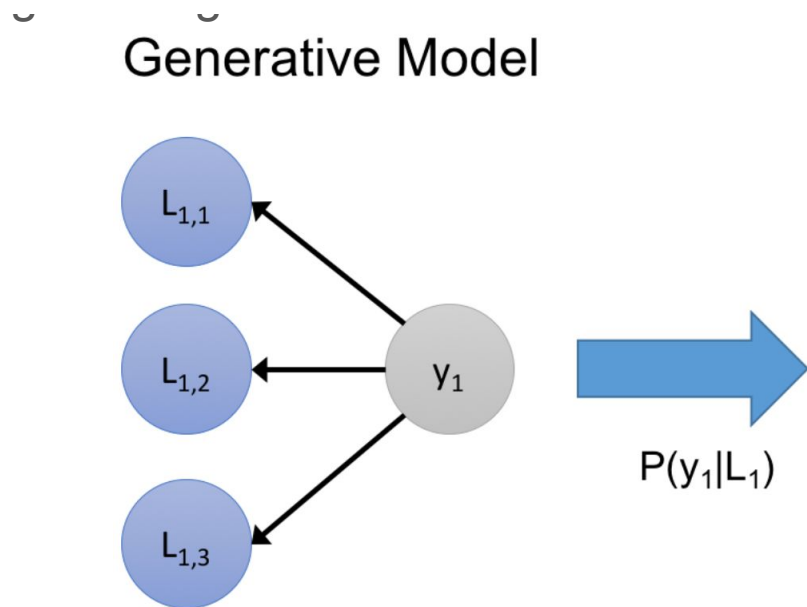
LFs

Integrating multiple noisy labels

How do we obtain probabilistic labels, $\tilde{\mathbf{Y}}$, from the label matrix, \mathbf{L} ?

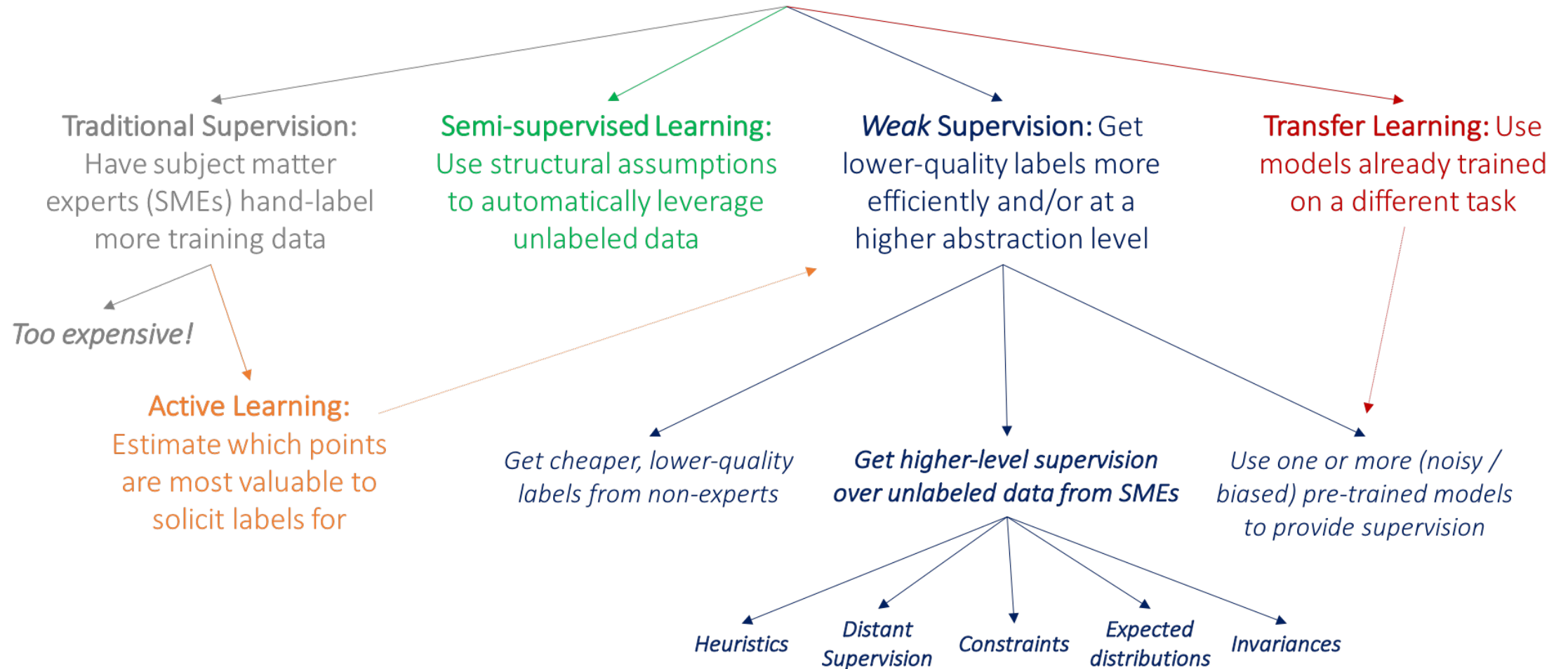
Approach 2

Train a generative model over $\mathbf{P}(\mathbf{L}, \mathbf{Y})$ where \mathbf{Y} are the **(unknown)** true labels



Summary: Weak/distant supervision

How to get more labeled training data?



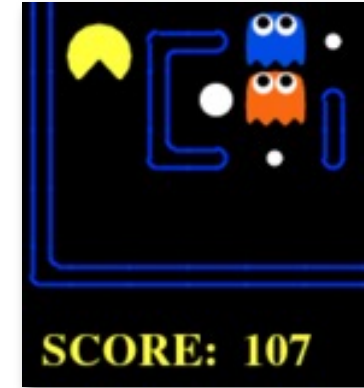
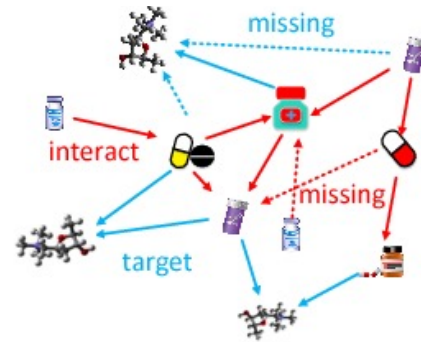
Summary: Weak/distant supervision

- Noisy labels from heuristics, knowledge bases, constraints, ...
- Integrating multiple noisy labels
 - Majority vote
 - Generative modeling
 - ...
- Not all information/experiences can easily be converted into labels
 - “Every part of speech sequence should have a verb”
 - “In a sentence with word ‘but’, the sentiment of text after ‘but’ dominates”
 - “Every image patch that is recognized as a bicycle should have at least one patch that is recognized as a wheel”
 - I have a “discriminator” model that can tell me whether a model-generated image is good or not
- Need a more flexible framework to incorporate all forms of experiences

Experiences of all kinds



Type-2 diabetes is 90% more common than type-1



Data examples

Rules/Constraints

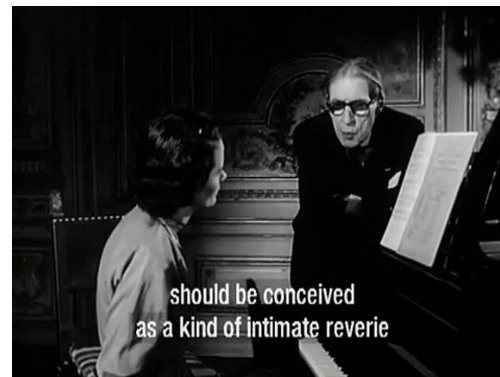
Knowledge graphs

Rewards

Auxiliary agents



Adversaries

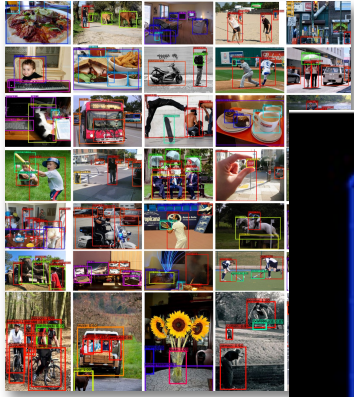


Master classes

...

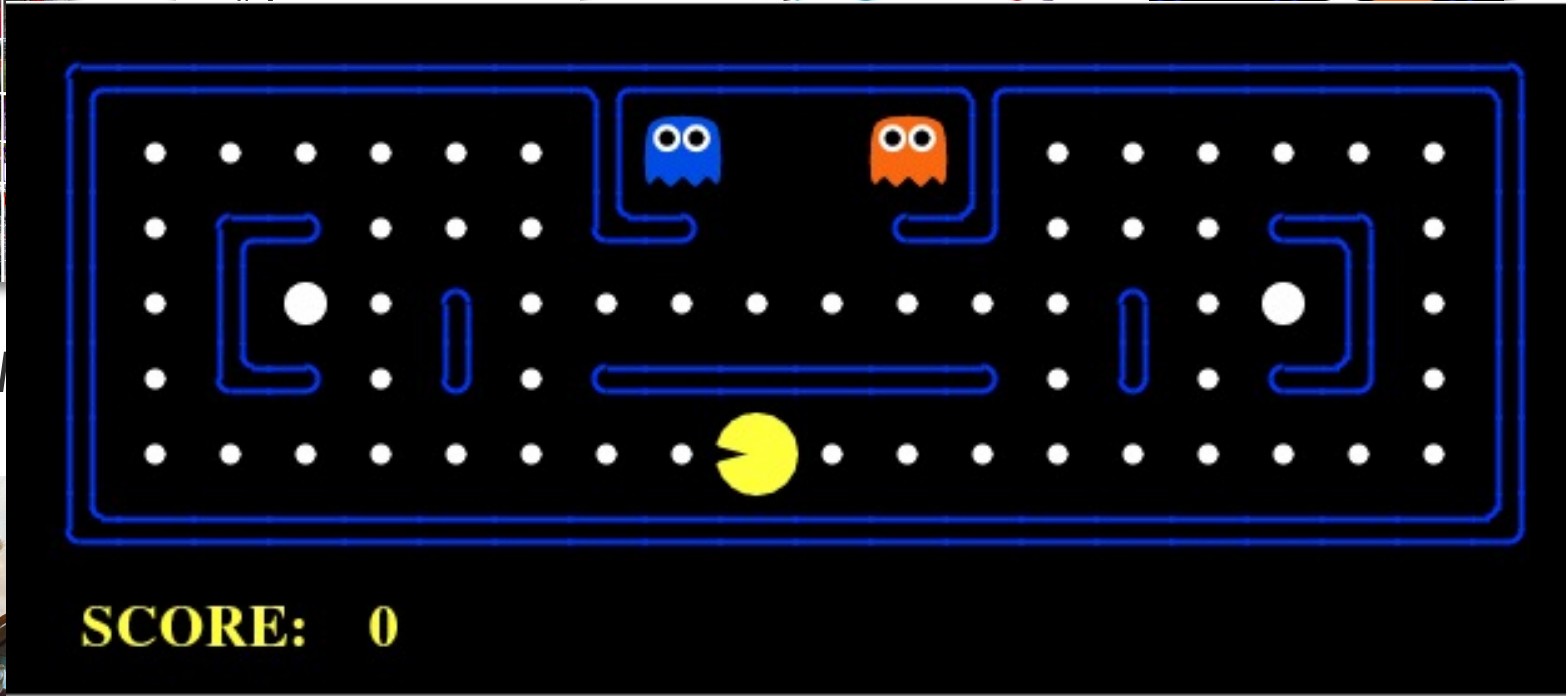
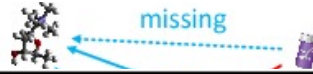
And all combinations thereof

Experiences of all kinds



Data examples

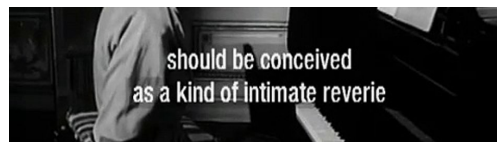
Type-2



Auxiliary agents



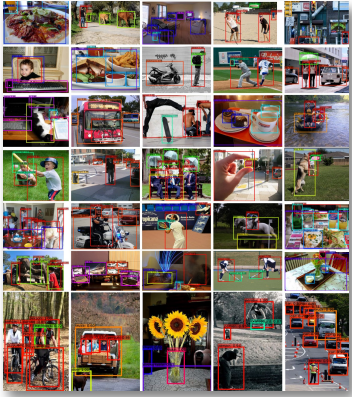
Adversaries



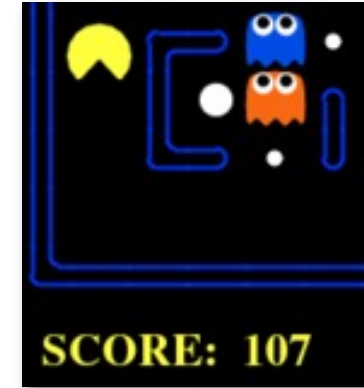
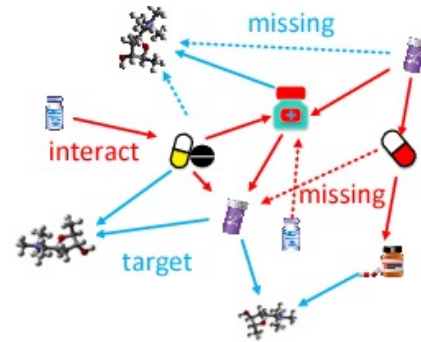
Master classes

ations thereof

Experiences of all kinds



Type-2 diabetes is 90% more common than type-1



Data examples

Rules/Constraints

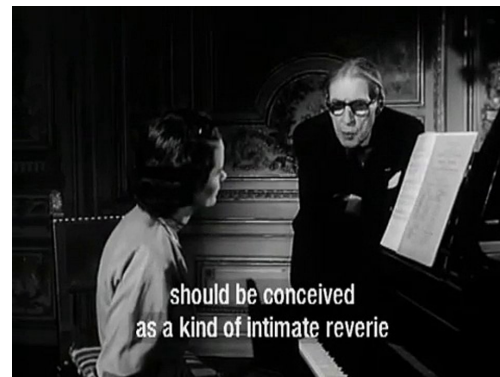
Knowledge graphs

Rewards

Auxiliary agents



Adversaries

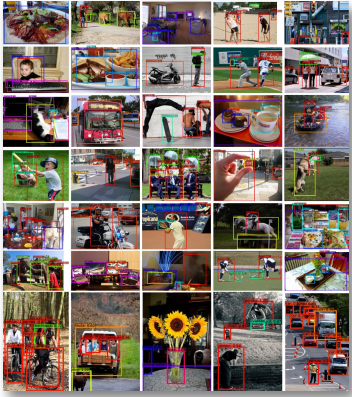


Master classes

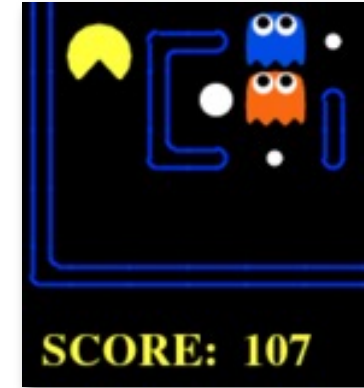
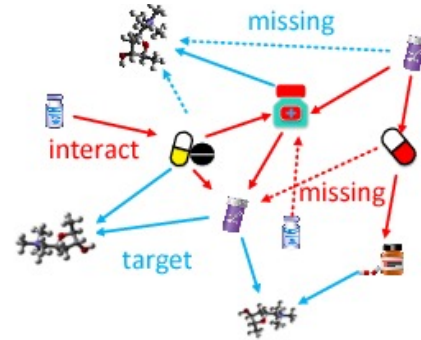
...

And all combinations thereof

Can we incorporate all types of experiences in learning?



Type-2 diabetes is 90% more common than type-1



Data examples

Rules/Constraints

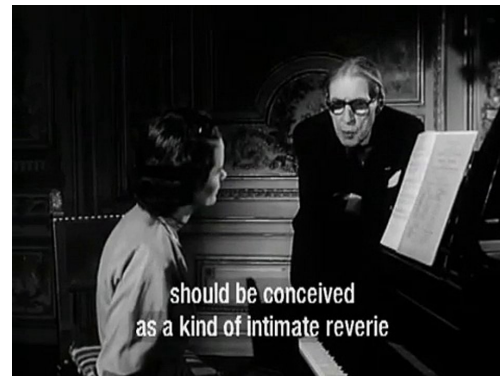
Knowledge graphs

Rewards

Auxiliary agents



Adversaries



Master classes

...

And all combinations thereof

Algorithm marketplace

Designs Driven by: experience, task, loss function, training procedure ...



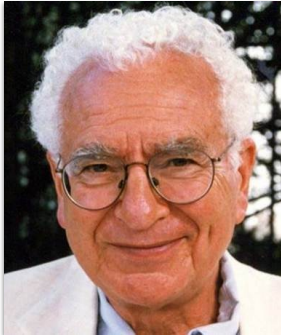
maximum likelihood estimation reinforcement learning as inference
data re-weighting inverse RL active learning
policy optimization
data augmentation reward-augmented maximum likelihood
label smoothing imitation learning softmax policy gradient
actor-critic adversarial domain adaptation
GANs posterior regularization
knowledge distillation intrinsic reward constraint-driven learning
prediction minimization generalized expectation
regularized Bayes learning from measurements
energy-based GANs
weak/distant supervision

Algorithm marketplace

Designs Driven by: experience, task, loss function, training procedure ...



Need a unifying perspective



You don't need something more in order to get something more.

-- Murray Gell-Mann (1929–2019), Physicist, Nobel laureate

ence

g

ood

ion

ing

"Standard equations" in Physics

Maxwell's Eqns:
original form

$e + \frac{df}{dx} + \frac{dg}{dy} + \frac{dh}{dz} = 0$	(1) Gauss' Law
$\mu\alpha = \frac{dH}{dy} - \frac{dG}{dz}$ $\mu\beta = \frac{dF}{dz} - \frac{dH}{dx}$ $\mu\gamma = \frac{dG}{dx} - \frac{dF}{dy}$	(2) Equivalent to Gauss' Law for magnetism
$P = \mu \left(\gamma \frac{dy}{dt} - \beta \frac{dz}{dt} \right) - \frac{dF}{dt} - \frac{d\Psi}{dz}$ $Q = \mu \left(\alpha \frac{dz}{dt} - \gamma \frac{dx}{dt} \right) - \frac{dG}{dt} - \frac{d\Psi}{dy}$ $R = \mu \left(\beta \frac{dx}{dt} - \alpha \frac{dy}{dt} \right) - \frac{dH}{dt} - \frac{d\Psi}{dx}$	(3) Faraday's Law (with the Lorentz Force and Poisson's Law)
$\frac{dy}{dx} - \frac{d\beta}{dz} = 4\pi p'$ $\frac{d\alpha}{dz} - \frac{d\gamma}{dx} = 4\pi q'$ $\frac{d\beta}{dx} - \frac{d\alpha}{dy} = 4\pi r'$	(4) Ampère-Maxwell Law
$P = -\xi p \quad Q = -\xi q \quad R = -\xi r$	Ohm's Law
$P = kf \quad Q = kg \quad R = kh$	The electric elasticity equation ($\mathbf{E} = \mathbf{D}/\epsilon$)
$\frac{de}{dt} + \frac{dp}{dx} + \frac{dq}{dy} + \frac{dr}{dz} = 0$	Continuity of charge

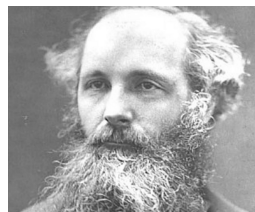
Maxwell's Eqns simplified w/ rotational symmetry

$$\nabla \cdot \mathbf{D} = \rho_V$$

$$\nabla \cdot \mathbf{B} = 0$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

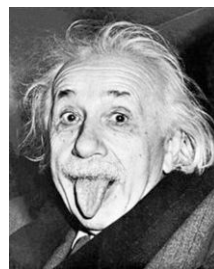
$$\nabla \times \mathbf{H} = \frac{\partial \mathbf{D}}{\partial t} + \mathbf{J}$$



Maxwell's Eqns further simplified w/ symmetry of special relativity

$$\epsilon^{uvk\lambda} \partial_v F_{k\lambda} = 0$$

$$\partial_v F^{uv} = \frac{4\pi}{c} j^u$$



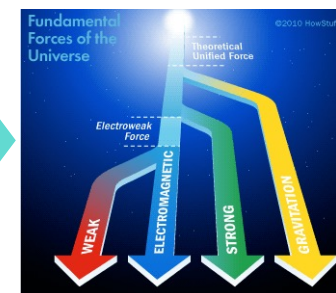
Standard Model w/ Yang-Mills theory and US(3) symmetry

$$\mathcal{L}_{gf} = -\frac{1}{2} \text{Tr}(F^2)$$

$$= -\frac{1}{4} F^{\alpha\mu\nu} F_{\mu\nu}^{\alpha}$$



Unification of fundamental forces?



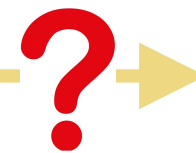
Diverse electro-magnetic theories



1861

1910s

1970s



A “Standardized” View of ML

Recap: MLE

- The most classical learning algorithm

- Supervised:

- Observe data $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$
- Solve with SGD

$$\min_{\theta} - \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} \left[\log p_{\theta}(\mathbf{y}^* | \mathbf{x}^*) \right]$$

- Unsupervised:

- Observe $\mathcal{D} = \{(\mathbf{x}^*)\}$, \mathbf{y} is latent variable
- Posterior $p_{\theta}(\mathbf{y} | \mathbf{x})$
- Solve with EM, etc

$$\min_{\theta} - \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} \left[\log \int_{\mathbf{y}} p_{\theta}(\mathbf{x}^*, \mathbf{y}) \right]$$

Recap: MLE as entropy maximization

- Duality between Supervised MLE and maximum entropy, when p is exponential family

$$\begin{aligned} \min_{p(\mathbf{x}, \mathbf{y})} H(p) & \xrightarrow{\text{Shannon entropy } H} \\ \text{s.t. } \mathbb{E}_p[T(\mathbf{x}, \mathbf{y})] &= \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}}[T(\mathbf{x}, \mathbf{y})] \xrightarrow{\text{features } T(\mathbf{x}, \mathbf{y})} \end{aligned}$$

data as constraints

* Proof with Lagrangian method

Recap: MLE as entropy maximization

- **Unsupervised MLE** can be achieved by maximizing the negative free energy:
 - Introduce **auxiliary** distribution $q(\mathbf{y}|\mathbf{x}^*)$ (and then play with its entropy and cross entropy, etc.)

$$\begin{aligned}\log \int_{\mathbf{y}} p_{\theta}(\mathbf{x}^*, \mathbf{y}) &= \mathbb{E}_{q(\mathbf{y}|\mathbf{x}^*)} \left[\log \frac{p_{\theta}(\mathbf{x}^*, \mathbf{y})}{q(\mathbf{y}|\mathbf{x}^*)} \right] + \text{KL}(q(\mathbf{y}|\mathbf{x}^*) || p_{\theta}(\mathbf{y}|\mathbf{x}^*)) \\ &\geq H(q(\mathbf{y}|\mathbf{x}^*)) + \mathbb{E}_{q(\mathbf{y}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{x}^*, \mathbf{y})]\end{aligned}$$

Bayesian Inference

- Posterior

$$p(\mathbf{z}|\mathcal{D}) = \frac{\pi(\mathbf{z}) \prod_{\mathbf{x}^* \in \mathcal{D}} p(\mathbf{x}^*|\mathbf{z})}{p(\mathcal{D})}$$

- Connecting to maximum entropy, as an optimization problem [Zellner, 1988]:

$$\min_{q(\mathbf{z})} -\mathbf{H}(q(\mathbf{z})) + \log p(\mathcal{D}) - \mathbb{E}_{q(\mathbf{z})} \left[\log \pi(\mathbf{z}) + \sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^*|\mathbf{z}) \right]$$

$$s. t. q(\mathbf{z}) \in \mathcal{P}$$

(the normality constraint of a probability distribution)

Posterior Regularization

- Under the optimization viewpoint of Bayesian inference, it's natural to consider other types of constraints that encode richer problem structures and domain knowledge
- Posterior regularization [Ganchev et al., 2010], or regularized Bayes (Zhu et al., 2014)

$$\min_{q, \xi} -\mathbf{H}(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} \left[\sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z}) \right] + U(\xi)$$

$$s.t. \quad q(\mathbf{z}) \in \mathcal{Q}(\xi)$$

$$\xi \geq 0,$$

- ξ : slack variables
- $U(\xi)$: a penalty function (e.g., L1 norm of ξ)
- $\mathcal{Q}(\xi)$: a subset of valid distributions over \mathbf{z} that satisfy the constraints determined by ξ

Posterior Regularization

$$\min_{q, \xi} -\mathbf{H}(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} \left[\sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z}) \right] + U(\xi)$$

$$s.t. \quad q(\mathbf{z}) \in \mathcal{Q}(\xi)$$

$$\xi \geq 0,$$

- Ex: let $T(\mathbf{x}^*, \mathbf{z})$ be a feature vector of data instance \mathbf{x}^* , a constrained posterior set $\mathcal{Q}(\xi)$ with "feature expectation" constraints can be defined as

$$\mathcal{Q}(\xi) := \{q(\mathbf{z}) : \mathbb{E}_q [T(\mathbf{x}^*; \mathbf{z})] \leq \xi\}$$

- i.e., bounds the feature expectations with ξ
- Assuming $U(\xi) = \sum \xi_i$, rewrite without slack variables

$$\min_{q, \xi} -\mathbf{H}(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} \left[\sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z}) \right] - c \cdot \mathbb{E}_{q(\mathbf{z})} [T(\mathbf{x}^*; \mathbf{z})]$$

Posterior Regularization

$$\min_{q, \xi} -H(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} \left[\sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z}) \right] + U(\xi)$$

$$s.t. \quad q(\mathbf{z}) \in \mathcal{Q}(\xi) := \{q(\mathbf{z}) : \mathbb{E}_q [T(\mathbf{x}^*; \mathbf{z})] \leq \xi\}$$

$$\xi \geq 0,$$

- $U(\xi) = \sum \xi_i$
- solution for $q(\mathbf{z})$:
$$q(\mathbf{z}) = \exp\{\log p(\mathbf{x}, \mathbf{z}) + T(\mathbf{x}^*, \mathbf{z})\} / Z$$
$$= p(\mathbf{x}, \mathbf{z}) \exp\{T(\mathbf{x}^*, \mathbf{z})\} / Z$$

** Proof with Lagrangian method*

Key Takeaways

- Recap: Supervised MLE and maximum entropy
- Recap: Unsupervised MLE and maximum entropy
- Bayesian inference and maximum entropy
 - Bayesian inference as optimization
- Posterior regularization:
 - Constrained Bayesian inference => constrained optimization

Questions?