DSC291: Machine Learning with Few Labels

Supervised / Unsupervised Learning

Zhiting Hu Lecture 3, April 8, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Logistics

- Office Hour this week:
 - Thursday 2pm PT
- Office Hour in future weeks:
 - Tuesday 2pm PT
- Office: HDSI 442
- TA's Office Hour TBA
- Will announce on Piazza later today



Machine learning solutions given few data (labels)

- (1) How can we make more efficient use of data?
 - Clean but small-size, Noisy, Out-of-domain Ο
- (2) Can we incorporate other types of experience in learning?



diabetes is 90% more common than type-1

Type-2



Rules/Constraints Knowledge graphs



Auxiliary agents



Adversaries



And all combinations thereof

Rewards

Machine learning solutions given few data (labels)

- (1) How can we make more efficient use of data?
 - Clean but small-size, Noisy, Out-of-domain, ...
- Algorithms
- **Supervised learning**: MLE, maximum entropy principle
- Unsupervised learning: EM, variational inference, VAEs
- Self-supervised learning: successful instances, e.g., BERT, GPTs, contrastive learning;

applications to downstream tasks

- Distant/weakly supervised learning: successful instances
- Data manipulation: augmentation, re-weighting, curriculum learning, ...
- Meta-learning

Mostly first half of the course

Machine learning solutions given few data (labels)

- (2) Can we incorporate other types of experience in learning?
 - Learning from auxiliary models, e.g., adversarial models.
 - Generative adversarial learning (GANs and variants), co-training, ...
 - Learning from structured knowledge
 - Posterior regularization, constraint-driven learning, ...
 - Learning from rewards
 - Reinforcement learning: model-free vs model-based, policy-based vs valuebased, on-policy vs off-policy, extrinsic reward vs intrinsic reward, ...,
 - Learning in dynamic environment (not covered)
 - Online learning, lifelong/continual learning, ...

Second half of the course

1eraign



Dseef







aries Maste





Algorithm marketplace

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





Where we are now? Where we want to be?

• Alchemy vs chemistry





maximum likelihood estimation reinforcement learning as inference

inverse RL active learning



Physics in the 1800's

- Electricity & magnetism:
 - Coulomb's law, Ampère, Faraday, ...
- Theory of light beams:
 Particle theory: Isaac Newton, Laplace, Plank
 Wave theory: Grimaldi, Chris Huygens, Thomas Young, Maxwell
 - Law of gravity
 - Aristotle, Galileo, Newton, ...









A "standard model" of ML



- Panoramically learn from all types of experience
- Subsumes many existing algorithms as special cases

Will discuss in later in the class

Lecture Schedule (tentative)

Supervised Learning

KL Divergence

"iscrete x • Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ in $u_{\mathcal{C}}$. and $q(\mathbf{x})$

 $q(\mathbf{x}) \log$

- a.k.a. Relative entropy \bigcirc
- $KL \ge 0$ (Jensen's inequality) \rightarrow homework \bigcirc
- **Questions:** Ο
 - If q is high and p is high in a region, then KL divergence is _____ in this region.

 $\mathrm{KL}(q(\mathbf{x}) || p(\mathbf{x})) =$

- If q is high and p is low in a region, then KL divergence is _____h in this region.
- If q is low in a region, then KL divergence is ______ in this region.

KL Divergence

• Kullback-Leibler (KL) divergence: measures the closeness of two distributions p(x) and q(x)

 $q(\mathbf{x})\log(\mathbf{x})$

Wasserstern Vistaice

- o a.k.a. Relative entropy
- \circ KL >= 0 (Jensen's inequality)
- \circ Intuitively:
 - If q is high and p is high, then we are happy (i.e. low KL divergence)

new Vide Mace

- If q is high and p is low then we pay a price (i.e. high KL divergence).
- If q is low then we don't care (i.e. also low KL divergence, regardless of p)
- not a true "distance":
 - not commutative (symmetric) KL(p||q) ! = KL(q||p)

 $\mathrm{KL}(q(\mathbf{x}) || p(\mathbf{x}))$

doesn't satisfy triangle inequality



Supervised Learning

- Model to be learned $p_{\theta}(x)$
- Observe full data $\mathcal{D} = \{ x_i \}_{i=1}^N$
 - e.g., x_i includes both input (e.g., image) and output (e.g., image label)
 - $\circ \ \mathcal{D}$ defines an empirical data distribution $\widetilde{p}({m x})$
 - $\boldsymbol{x} \sim \mathcal{D} \iff \boldsymbol{x} \sim \tilde{p}(\boldsymbol{x})$
- Maximum Likelihood Estimation (MLE)
 - The most classical learning algorithm
- Question: Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

 $\min_{\alpha} - \mathbb{E}_{x \sim \tilde{p}(x)}$

 $\log p_{\theta}(\boldsymbol{x})$

18



Supervised Learning

- Model to be learned $p_{\theta}(x)$
- Observe full data $\mathcal{D} = \{ \mathbf{x}_i \}_{i=1}^N$
 - e.g., x_i includes both input (e.g., image) and output (e.g., image label)
 - $\circ \mathcal{D}$ defines an empirical data distribution $\widetilde{p}(\mathbf{x})$
 - $\boldsymbol{x} \sim \mathcal{D} \iff \boldsymbol{x} \sim \tilde{p}(\boldsymbol{x})$
- Maximum Likelihood Estimation (MLE)
 - The most classical learning algorithm

• Question: Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

$$KL(\tilde{p}(\boldsymbol{x}) || p_{\theta}(\boldsymbol{x})) = -\mathbb{E}_{\tilde{p}(\boldsymbol{x})}[\log p_{\theta}(\boldsymbol{x})] + H(\tilde{p}(\boldsymbol{x}))$$

Cross entropy

 $\min_{\theta} - \mathbb{E}_{\boldsymbol{x} \sim \tilde{p}(\boldsymbol{x})} \quad \log p_{\theta}(\boldsymbol{x})$

- Each data instance is partitioned into two parts:
 - \circ observed variables x
 - \circ latent (unobserved) variables z
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



- Each data instance is partitioned into two parts:
 - observed variables \boldsymbol{x} Ο
 - latent (unobserved) variables z \bigcirc
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



= (x1, 12)

- Each data instance is partitioned into two parts:
 - \circ observed variables \pmb{x}
 - latent (unobserved) variables Z
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



Each data instance is partitioned into two parts:

Budgets

1.2,3.4%

= 0.001 24

Ginty

Mone.

X: doc: (I, am, good.

Education

tapic id

Children

observed variables \boldsymbol{x} Ο

model

- latent (unobserved) variables z \bigcirc
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$

Arts

- Each data instance is partitioned into two parts:
 - \circ observed variables x
 - \circ latent (unobserved) variables z
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



- Each data instance is partitioned into two parts:
 - observed variables x \bigcirc
 - latent (unobserved) variables zΟ
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



Why is Unsupervised Learning Harder?

• Complete log likelihood: if both x and z can be observed, then

 $\ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta) = \log p(\mathbf{z} | \theta_z) + \log p(\mathbf{x} | \mathbf{z}, \theta_x)$

• Decomposes into a sum of factors, the parameter for each factor can be estimated separately

Now z is not observed:

• Incomplete (or marginal) log likelihood: with z unobserved, our objective becomes the log of a marginal probability:

 $\ell(\theta; \boldsymbol{x}) = \log p(\boldsymbol{x}|\theta)$

Why is Unsupervised Learning Harder?

• Complete log likelihood: if both x and z can be observed, then

 $\ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z}|\theta) = \log p(\mathbf{z}|\theta_z) + \log p(\mathbf{x}|\mathbf{z}, \theta_x)$

 Decomposes into a sum of factors, the parameter for each factor can be estimated separately

Now z is not observed:

 Incomplete (or marginal) log likelihood: with z unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{z} p(\mathbf{x}, \mathbf{z}|\theta)$$

- All parameters become coupled together
- In other models when z is complex (continuous) variables (as we'll see later), marginalization over z is intractable.

Expectation Maximization (EM): Intuition

Expectation Maximization (EM): Intuition

 $\max_{\boldsymbol{\theta}} \ell_c(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{z}) = \log p(\boldsymbol{x}, \boldsymbol{z} | \boldsymbol{\theta})$

- Supervised MLE is easy:
 - Observe both x and z
- Unsupervised MLE is hard:
 Observe only *x*

E-step: $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}, \theta)$

We don't actually observe q, let's estimate it

 $\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{x}} p(\mathbf{x}, \mathbf{z}|\theta)$

M-step: $\max_{\theta} \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})} [\log p(\boldsymbol{x}, \boldsymbol{z}|\theta)]$

Let's "pretend" we also observe Z (its distribution)

Expectation Maximization (EM): Intuition

- Supervised MLE is easy:
 - Observe both x and z
- Unsupervised MLE is hard:
 Observe only *x*

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{z} p(\mathbf{x}, \mathbf{z}|\theta)$$

 $\max_{\boldsymbol{\alpha}} \ell_c(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{z}) = \log p(\boldsymbol{x}, \boldsymbol{z} | \boldsymbol{\theta})$

• EM, intuitively:

E-step:
$$q^{t+1}(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z}|\boldsymbol{x}, \theta^{t})$$

M-step: $\max_{\theta} \mathbb{E}_{q^{t+1}(\boldsymbol{z}|\boldsymbol{x})}[\log p(\boldsymbol{x}, \boldsymbol{z}|\theta)]$
This is an iterative

We don't actually observe q, let's estimate it

Let's "pretend" we also observe Z (its distribution)

process

Questions?