

Week 5: Machine Learning

1. Introduction

1.1. Background

1.2. Different Ways

1.2.1 Idealistic

1.2.2. Pragmatic

1.3. Idealistic (Bayesian) Perspective

1.3.1. Basic Ideas

1.3.2. Problems

1.4. Pragmatic Perspective

1.5. Classification

Supervised Learning

Unsupervised Learning

Semi-supervised Learning

Active-Learning

Reinforcement Learning

1.6. Generalization

Approaches

overfitting

2. Maximum Likelihood

2.1. Basic Idea

2.2. Example

2.3. Comparison between Maximum Likelihood and Bayesian

Bayesian learning:

Maximum likelihood (ML)

Maximum a posteriori(MAP):

2.4. Example: Learning Bayesian Network with Maximum Likelihood

3. Learning with hidden variables

3.1. Necessity

3.2. Problems

3.3. EM algorithm

Intuition

Basic Idea

Two Steps

4. Two examples of EM algorithm

4.1. EM for the Bayesian Network

4.2. EM for HMMs

1. Introduction

1.1. Background

Why we need learning?

- **No model available:** Human designers can not provide models for all possible situations an intelligent agent may encounter
- **Adaptivity:** environment may change over time
- Human don't understand the task well enough, **not possible to manually program**

1.2. Different Ways

1.2.1 Idealistic

- maintain a **belief** over true way the world works (**possible hypotheses**)
- **statistical learning**, always use **Bayes' rule**

1.2.2. Pragmatic

- anything that improves performance
- always use **optimization**

1.3. Idealistic (Bayesian) Perspective

1.3.1. Basic Ideas

- maintain **belief** over how the world works: hypotheses H
 - (**learning**) using bayes rule to update its beliefs

$$P(H|d) = P(d|H)P(H)/P(d)$$

- Then **act** in a bayesian way:

$$V(a) = \sum_h P(h|d)u(a|h)$$

1.3.2. Problems

How to specify the class of hypothesis?

- World is complex, need a huge class H
- Limit the class to allow for tractable inference, but the true model may not in H
- So there are always two **tradeoff**
 - expressiveness of a hypothesis space and the complexity of finding a good hypothesis within that space

1.4. Pragmatic Perspective

From pragmatic perspective, learning is **optimization**.

There are always two ways:

- **end-to-end learning:**
 - parametrize ‘actions’ using some parameters θ
 - directly optimize $V(\theta)$
- other typical approach, for example **maximum likelihood**

1.5. Classification

Supervised Learning

General Process:

- bag of **training data** $d = \langle x_i, y_i \rangle_{i=1 \dots N}$
- assumption: labels generated by ‘**true**’ function $y = f(x)$
- goal: find hypothesis $h(x) \approx f(x)$

Instantiations:

regression, classification

Unsupervised Learning

- Learn patterns in the input **without explicit feedback (no labels)**
- Most common task is **clustering**

Semi-supervised Learning

- large bag of data, only **a few are labeled**

Active-Learning

- large bag of data, only **a few are labeled**
- what point should we ask an annotator to label

Reinforcement Learning

Learns from a series of reinforcements: **rewards or punishments**

1.6. Generalization

How do we know $h \approx f$

Approaches

- use theorems: computational/statistical learning theory
 - For example Error
- use experiments:
 - on a new set of data: test data (test data is **never used**, also not for model selection)

overfitting

2. Maximum Likelihood

2.1. Basic Idea

- instead of computing posterior $P(H|d)$, we optimize

$$h_{ML} = \max_h P(d|h)$$

- to do this optimization, we usually use log likelihood

$$L(h) = \log P(d|h)$$

$$\begin{aligned} h_{ML} &= \max_h \log P(d | h) \\ &= \max_h \log \prod_i P(d_i | h) = \max_h \sum_i \log P(d_i | h) \end{aligned}$$

2.2. Example

Example: a coin toss...

- We have a coin and toss it N times...
 - k heads, l=N-k tails
→ what is the prob. of heads?
- Lets call P(head) = θ
- So.. likelihood:
 $P(d|\theta) = \theta^k (1-\theta)^l$

- ▷ $h_{ML} = \max_h \log P(d|h)$
 $= \max_h \log \prod_i P(d_i|h) = \max_h \sum_i \log$
- ▷ usually much easier to optimize

Maximum likelihood Bernoulli (20.2.1)

$$\begin{aligned} L(\theta) &= \log \prod_{i=1}^k \theta \prod_{i=1}^l (1-\theta) \\ &= \log \theta^k (1-\theta)^l \\ &= k \log \theta + l \log (1-\theta) \end{aligned}$$

Its derivative:

$$\frac{d}{d\theta} L(\theta) = \frac{k}{\theta} - \frac{l}{(1-\theta)}$$

equating with 0 and solving to find the maximum:

$$\begin{aligned} \frac{k}{\theta} - \frac{l}{(1-\theta)} &= 0 \\ \Leftrightarrow \frac{k}{\theta} &= \frac{l}{(1-\theta)} \\ \Leftrightarrow k(1-\theta) &= l\theta \\ \Leftrightarrow k &= (l+k)\theta \\ \Leftrightarrow \frac{k}{l+k} &= \theta = \frac{k}{N} \end{aligned}$$

2.3. Comparison between Maximum Likelihood and Bayesian

Bayesian learning:

- computing posterior $P(H|d)$
 - uses prior information $P(h)$
- use in weighted manner to select best action: $V(a) = \sum_h P(h|d) u(a, h)$

Maximum likelihood (ML)

- $h_{ML} = \max_h P(d|h)$
- select action according $V(a) = u(a, h_{ML})$
- prone to “overfitting”

Maximum a posteriori(MAP):

ML + priors

- ▷ $h_{MAP} = \max_h P(d|h) P(h)$
- ▷ can still overfit

2.4. Example: Learning Bayesian Network with Maximum Likelihood

$$L(\theta) = \log P(d; \theta)$$

$$= \log \prod_i P(\langle \text{party}_i, \text{rested}_i, \text{studied}_i, \text{pass}_i \rangle; \theta)$$

$$= \sum_i \log P(\langle \text{party}_i, \text{rested}_i, \text{studied}_i, \text{pass}_i \rangle; \theta)$$

$$= \sum_i \log P(\text{party}_i; \theta) P(\text{rested}_i | \text{party}_i; \theta) \\ P(\text{studied}_i; \theta) P(\text{pass}_i | \text{studied}_i, \text{rested}_i; \theta)$$

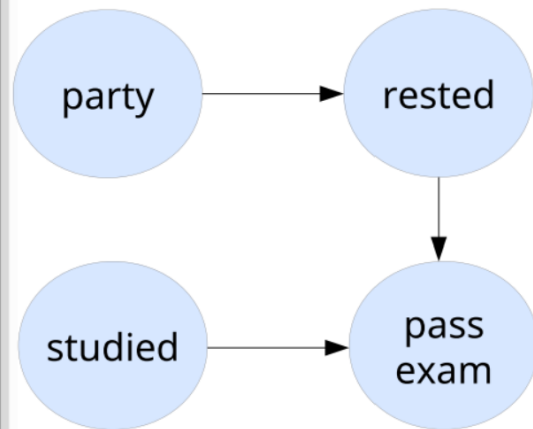
$$= \sum_i \log P(\text{party}_i; \theta)$$

$$+ \sum_i \log P(\text{rested}_i | \text{party}_i; \theta)$$

$$+ \sum_i \log P(\text{studied}_i; \theta)$$

$$+ \sum_i \log P(\text{pass}_i | \text{studied}_i, \text{rested}_i; \theta)$$

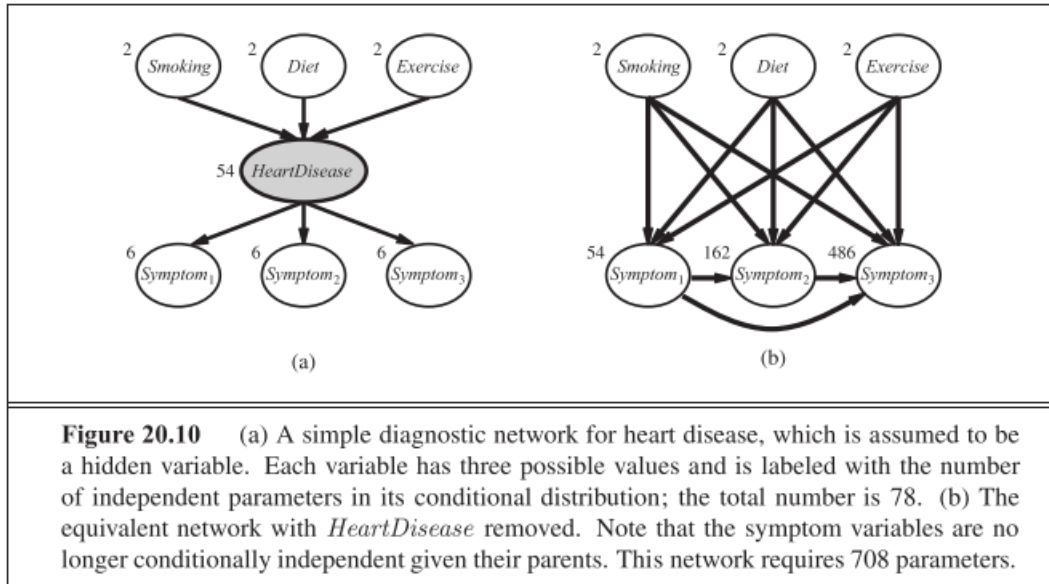
ters...?



3. Learning with hidden variables

3.1. Necessity

Hidden Variables can help us to eliminate the number of parameters



3.2. Problems

But we don't know the parameter of the hidden variables

3.3. EM algorithm

Intuition

learn better parameters given estimated variables

Basic Idea

$$\theta^{(k+1)} = \arg \max_{\theta} \sum_z P(\mathbf{Z} = \mathbf{z} \mid \mathbf{x}, \theta^{(k)}) L(\mathbf{x}, \mathbf{Z} = \mathbf{z} \mid \theta)$$

where:

- $\theta^{(k+1)}$ is the new parameter vector
- \mathbf{z} is the vector of values for latent variables \mathbf{Z}
- \mathbf{x} is the value of observed variables
- $P(\mathbf{Z} = \mathbf{z} \mid \mathbf{x}, \theta^{(k)})$ the 'estimation' of the latent variables given $\mathbf{x}, \theta^{(k)}$
- $L(\mathbf{x}, \mathbf{Z} = \mathbf{z} \mid \theta)$ the log likelihood:

$$L(\mathbf{x}, \mathbf{Z} = \mathbf{z} \mid \theta) = \log P(\mathbf{x}, \mathbf{Z} = \mathbf{z} \mid \theta)$$

Two Steps

1. **E-step**, where ‘E’ stands for *expectation*. Here the summation over \mathbf{z} is performed to compute the expectation. Note that, in order to accomplish this, it needs to compute, or *estimate*, the posterior $P(\mathbf{Z} = \mathbf{z} | \mathbf{x}, \theta^{(k)})$.
2. **M-step**. Which performs the maximization over parameters θ .

1. E-step:
estimate $p_{ij} = P(C = i | x_j)$ using current cluster parameters
2. M-step:
update cluster parameters using p_{ij}

4. Two examples of EM algorithm

4.1. EM for the Bayesian Network

- Due to hidden variable, directly optimizing log likelihood is hard:

$$\begin{aligned} L(\mathbf{x} | \theta) &= \log \Pr(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} \Pr(\mathbf{x}, \mathbf{Z} = \mathbf{z} | \theta) \\ &= \log \sum_{\mathbf{z}} \prod_{i=1}^N \Pr(x_i, Z_i = z_i | \theta) \\ &= \log \sum_{\mathbf{z}} \prod_{i=1}^N \Pr(z_i | \theta_B) \Pr(flavor_i | z_i, \theta_{Fz_i}) \Pr(wrapper_i | z_i, \theta_{Wz_i}) \Pr(hole_i | z_i, \theta_{Hz_i}) \end{aligned}$$

- With EM algorithm

$$\theta^{(k+1)} = \arg \max_{\theta} \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}, \theta^{(k)}) L(\mathbf{x}, \mathbf{z}|\theta)$$

which uses the *full* (or ‘*completed*’) *log-likelihood*:

$$\begin{aligned} L(\mathbf{x}, \mathbf{z}|\theta) &= \log \Pr(\mathbf{x}, \mathbf{z}|\theta) = \log \prod_{i=1}^N \Pr(x_i, z_i|\theta) \\ &= \sum_{i=1}^N \log \Pr(x_i, z_i|\theta) \\ &= \sum_{i=1}^N \log \Pr(z_i|\theta_B) + \sum_{i=1}^N \log \Pr(flavor_i|z_i, \theta_{Fz_i}) \\ &\quad + \sum_{i=1}^N \log \Pr(wrapper_i|z_i, \theta_{Wz_i}) + \sum_{i=1}^N \log \Pr(hole_i|z_i, \theta_{Hz_i}) \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^N \log \Pr(z_i|\theta_B) &= \sum_{i \text{ s.t. } z_i=1} \log \Pr(z_i|\theta_B) + \sum_{i \text{ s.t. } z_i=2} \log \Pr(z_i|\theta_B) \\ &= \sum_{i \text{ s.t. } z_i=1} \log \theta_B + \sum_{i \text{ s.t. } z_i=2} \log(1 - \theta_B) \\ &= N_1 \log \theta_B + N_2 \log(1 - \theta_B) \\ &= N_1 \log \theta_B + (N - N_1) \log(1 - \theta_B) \end{aligned}$$

Where $N_1 = N(bag = 1|\mathbf{z})$.

If \mathbf{z} was correct, this would then lead (by taking derivative and setting to 0) to

$$\theta_B \leftarrow \frac{N_1}{N}$$

$$\begin{aligned}
& \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}, \theta^{(k)}) [N_1 \log \theta_B + (N - N_1) \log(1 - \theta_B)] \\
&= \sum_n P(N_1 = n|\mathbf{x}, \theta^{(k)}) [n \log \theta_B + (N - n) \log(1 - \theta_B)] \\
&= \sum_n P(N_1 = n|\mathbf{x}, \theta^{(k)}) n \log \theta_B + \sum_n P(N_1 = n|\mathbf{x}, \theta^{(k)}) (N - n) \log(1 - \theta_B) \\
&= \log \theta_B \cdot \sum_n P(N_1 = n|\mathbf{x}, \theta^{(k)}) n + \log(1 - \theta_B) \cdot (N - \sum_n P(N_1 = n|\mathbf{x}, \theta^{(k)}) n) \\
&= \log \theta_B \cdot \hat{N}(\text{Bag} = 1) + \log(1 - \theta_B) \cdot (N - \hat{N}(\text{Bag} = 1))
\end{aligned}$$

with $\hat{N}(\text{Bag} = 1)$ is the *expected* counts for bag 1.

- Finally, this leads (by taking derivative and setting to 0) to

$$\theta_B^{(k+1)} \leftarrow \frac{\hat{N}(\text{Bag} = 1)}{N}$$

4.2. EM for HMMs