10-701 Recitation: Probabilistic and graphical models

Abulhair Saparov

- There are a lot of different naming conventions.
 - Seems to be the case a lot in machine learning
 - and outside machine learning as well...
- I will present some definitions using nomenclature that I think is intuitive, but also not far off from what other people use.

What is a **probabilistic model**?

A **probabilistic model** is a collection of random variables. The random variables can be divided into two categories:

- 1. the **observations** (data)
- 2. the hidden variables

Intuitively, a probabilistic model is a *description* of how your observations were generated.

It can be a hypothesis of the mechanism that underlies your data.

For more intuition, imagine we have a system of equations.

$$x + 2y = 4$$

If x is 0, what is y? If $x \sim \mathcal{N}(0, 1)$, what is y?

You can think of a probabilistic model as a system where the variables can be random.

A probabilistic model is a collection of random variables: $\{x, \theta\}$ The random variables can be divided into two categories:

- 1. **x:** the **observations** (data)
- 2. **0**: the hidden variables

the prior distribution is p(θ)
 the likelihood is p(x|θ)
the posterior distribution is p(θ|x)
 the joint distribution is p(x,θ)

So then, what's a graphical model?

A graphical model is a graphical representation of a probabilistic model.

There are different ways to represent probabilistic models graphically:

- Bayesian networks
- Factor graphs
- Markov random fields

When people say "graphical model", they usually mean

graph + probabilistic model.

Common convention: observations are depicted as shaded nodes, whereas hidden variables are unshaded.

Example

probabilistic model

w ~ Uniform(0,1)
x ~ Bernoulli(w)
y ~ Bernoulli(w)
z = x + y

You can think of **x** and **y** as two coin flips, and **w** represents the fairness of the coin. graphical model



Example

It is easy to factorize the joint distribution by looking at the graph structure:

 $p(\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}) =$ $p(\mathbf{w})p(\mathbf{x}|\mathbf{w})p(\mathbf{y}|\mathbf{w})p(\mathbf{z}|\mathbf{x}, \mathbf{y})$

graphical model



w ~ Uniform(0,1) x ~ Bernoulli(w) y ~ Bernoulli(w) z = x + y

Example

What is
$$p(x|w)$$
?
 $p\{x = 0|w\} = 1 - w \text{ and } p\{x = 1|w\} = w$
 $p\{y = 0|w\} = 1 - w \text{ and } p\{y = 1|w\} = w$

What is
$$p(\mathbf{z}|\mathbf{x},\mathbf{y},\mathbf{w}) = p(\mathbf{z}|\mathbf{x},\mathbf{y})$$
?
 $p(\mathbf{z}|\mathbf{x},\mathbf{y}) = \delta\{\mathbf{z} = \mathbf{x} + \mathbf{y}\}$

What if we want to get rid of **x** and **y**? We can **marginalize** out **x** and **y**.

Lets compute $p(\mathbf{z}|\mathbf{w})$.

w ~ Uniform(0,1)
x ~ Bernoulli(w)
y ~ Bernoulli(w)
z = x + y

$$\begin{split} p\{z=0|w\} &= \sum_{x\in\{0,1\}} p\{z=0|x,w\} p(x|w), \\ &= \sum_{x\in\{0,1\}} \sum_{y\in\{0,1\}} p\{z=0|x,y,w\} p(x|w) p(y|w), \\ &= \sum_{x\in\{0,1\}} \sum_{y\in\{0,1\}} \delta\{x=0,y=0\} p(x|w) p(y|w), \\ &= p\{x=0|w\} p\{y=0|w\}, \\ &= (1-w)^2. \end{split}$$

Example

Example

w ~ Uniform(0,1)
x ~ Bernoulli(w)
y ~ Bernoulli(w)
z = x + y

Repeating the process for the other possible values of **z**:

 $p\{z = 0 | w\} = (1 - w)^{2}$ $p\{z = 1 | w\} = 2w(1 - w)$ $p\{z = 2 | w\} = w^{2}$

This is a new model:

w ~ Uniform(0,1)

z ~ the above discrete distribution

The reasoning and intuition behind variable elimination and belief propagation is the same.





What is a good probabilistic model for this data?

The points look like they came from four normal distributions.

$$\mu_{1}, \dots, \mu_{k} \sim \mathcal{N}(0, 10I),$$

$$x_{1}^{(1)}, \dots, x_{n}^{(1)} \sim \mathcal{N}(\mu_{1}, I),$$

$$x_{1}^{(2)}, \dots, x_{n}^{(2)} \sim \mathcal{N}(\mu_{2}, I),$$
and so on . . . for k clusters



For simplicity, we fixed the covariances to the identity. If desired, you could make them unknown variables and use, for instance, an inverse-Wishart prior (see last week's recitation resources). We can re-write this model:

$$\mu_1, \dots, \mu_k \sim \mathcal{N}(0, 10I),$$

$$z_1, \dots, z_n = 1,$$

$$z_{n+1}, \dots, z_{2n} = 2,$$
etc...
$$x_i \sim \mathcal{N}(\mu_{z_i}, I) \text{ i.i.d. } i = 1, \dots, kn.$$





What is a good model for this data?

We no longer have knowledge of the class assignments z.

$$\mu_1, \dots, \mu_k \sim \mathcal{N}(0, 10I),$$

$$z_i = \text{Categorical}(\pi),$$

$$x_i \sim \mathcal{N}(\mu_{z_i}, I) \text{ i.i.d. } i = 1, \dots, kn.$$





$$\mu_1, \dots, \mu_k \sim \mathcal{N}(0, 10I),$$

$$z_i = \text{Categorical}(\pi),$$

$$x_i \sim \mathcal{N}(\mu_{z_i}, I) \text{ i.i.d. } i = 1, \dots, kn.$$

Let's try to do inference using MAP in this model. Write the log-posterior: $\log p(\mu_1, ..., \mu_k, z_1, ..., z_{kn} | x_1, ..., x_{kn})$.

$$\log p(\boldsymbol{\mu}, \boldsymbol{z} | \boldsymbol{x}) = \log p(\boldsymbol{x} | \boldsymbol{\mu}, \boldsymbol{z}) + \log p(\boldsymbol{\mu}) + \log p(\boldsymbol{z}) + C,$$

$$= \sum_{i=1}^{kn} \log p(x_i | \boldsymbol{\mu}, z_i) + \sum_{j=1}^{k} \log p(\mu_j) + \sum_{i=1}^{kn} \log p(z_i),$$

$$= -\frac{1}{2} \sum_{i=1}^{kn} (x_i - \mu_{z_i})^\top (x_i - \mu_{z_i}) - \frac{1}{20} \sum_{j=1}^{k} \mu_j^\top \mu_j + \sum_{i=1}^{kn} \log \pi_{z_i}.$$

Expectation maximization (EM)

- If we knew either μ or z, then MLE/MAP would be easier.
- EM is an inference algorithm for computing MLE or MAP.
- Given any probabilistic model with observations x and hidden variables θ, we first subdivide the hidden variables θ into two classes: z and μ.
- Start with a guess for **µ**.

E step. compute the expectation using our current estimate of μ :

$$q(\mathbf{\mu}) = \mathsf{E}_{\mathsf{p}(\mathsf{z}|\mathbf{x}, \boldsymbol{\mu}^*)}[\log \mathsf{p}(\mathbf{\mu}, \mathbf{z}|\mathbf{x})]$$

M step. update the estimate of μ^* by maximizing $q(\mu)$:

$$\mu^* = \arg \max q(\mu)$$

Repeat until convergence.

Expectation maximization

• We need to compute $p(z|x,\mu)$ to do the E step.

$$p(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{\mu})p(\boldsymbol{x},\boldsymbol{\mu}) = p(\boldsymbol{x},\boldsymbol{z},\boldsymbol{\mu}),$$

$$p(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{\mu}) = p(\boldsymbol{x}|\boldsymbol{z},\boldsymbol{\mu})p(\boldsymbol{\mu})p(\boldsymbol{z})/p(\boldsymbol{x},\boldsymbol{\mu}),$$

$$\log p(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{\mu}) = \log p(\boldsymbol{x}|\boldsymbol{z},\boldsymbol{\mu}) + \log p(\boldsymbol{z}) + C,$$

$$\log p\{z_i = j|\boldsymbol{x},\boldsymbol{\mu}\} = \log p(x_i|z_i = j,\mu_j) + \log p\{z_i = j\} + C,$$

$$= -\frac{1}{2}(x_i - \mu_j)^{\top}(x_i - \mu_j) + \log \pi_j + C,$$

$$p\{z_i = j|\boldsymbol{x},\boldsymbol{\mu}\} \propto \pi_j \exp\left\{-\frac{1}{2}(x_i - \mu_j)^{\top}(x_i - \mu_j)\right\}.$$



Expectation maximization

We can use $p(\mathbf{z}|\mathbf{x},\boldsymbol{\mu})$ to compute q. $q(\boldsymbol{\mu}) = \mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{\mu}^*)}[\log p(\boldsymbol{\mu},\boldsymbol{z}|\boldsymbol{x})],$ $= \mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{\mu}^*)}[\log p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{z})] + \log p(\boldsymbol{\mu}) + C,$ kn $= \sum \mathbb{E}_{p(z_i|\boldsymbol{x},\boldsymbol{\mu}^*)}[\log p(x_i|\boldsymbol{\mu}, z_i)] + \log p(\boldsymbol{\mu}) + C,$ i=1 $= -\frac{1}{2} \sum_{p(z_i|\boldsymbol{x},\boldsymbol{\mu}^*)} \mathbb{E}_{p(z_i|\boldsymbol{x},\boldsymbol{\mu}^*)} [(x_i - \mu_{z_i})^\top (x_i - \mu_{z_i})] + \log p(\boldsymbol{\mu}) + C,$ $= -\frac{1}{2} \sum_{k=1}^{k} \sum_{i=1}^{k} p\{z_i = j | \boldsymbol{x}, \boldsymbol{\mu}^*\} (x_i - \mu_j)^\top (x_i - \mu_j) + \log p(\boldsymbol{\mu}) + C,$ kn

$$q(\mu_j) = -\frac{1}{2} \sum_{i=1}^{\infty} p\{z_i = j | \boldsymbol{x}, \boldsymbol{\mu}^*\} (x_i - \mu_j)^\top (x_i - \mu_j) - \frac{1}{20} \mu_j^\top \mu_j + C.$$



Expectation maximization

• Now, for the M step, we just maximize q.



$$\mu_{j}^{*} = \arg \max_{\mu_{j}} q(\mu_{j}),$$

$$\frac{\partial q}{\partial \mu_{j}} = \sum_{i=1}^{kn} p\{z_{i} = j | \boldsymbol{x}, \boldsymbol{\mu}^{*}\}(x_{i} - \mu_{j}) - \frac{1}{10}\mu_{j},$$

$$0 = \sum_{i=1}^{kn} p\{z_{i} = j | \boldsymbol{x}, \boldsymbol{\mu}^{*}\}x_{i} - \mu_{j}^{*} \left(\frac{1}{10} + \sum_{i=1}^{kn} p\{z_{i} = j | \boldsymbol{x}, \boldsymbol{\mu}^{*}\}\right),$$

$$\mu_{j}^{*} = \left(\frac{1}{10} + \sum_{i=1}^{kn} p\{z_{i} = j | \boldsymbol{x}, \boldsymbol{\mu}^{*}\}\right)^{-1} \sum_{i=1}^{kn} p\{z_{i} = j | \boldsymbol{x}, \boldsymbol{\mu}^{*}\}x_{i}$$





























