## Lecture 24 announcements

- HW 12 - due Thursday

## Lecture 24 outline

- Semi-supervised learning (part 2)

  • Mixture models for SSL

  • Maximum likelihood estimate (MLE)

  • Expectation maximization (EM)

✳ CORRECTIONS ON p.4, 10.

# Mixture Models for SSL

We want to find $p(y \mid \underline{x})$

Let's model each class as a specified density with unknown parameters:

$$p(\underline{x}, y \mid \underline{\theta}) = p(\underline{x} \mid y, \underline{\theta})\, p(y \mid \underline{\theta})$$

$$(1) \qquad = \underbrace{p(\underline{x} \mid y, \underline{\theta})}\; \pi_y$$

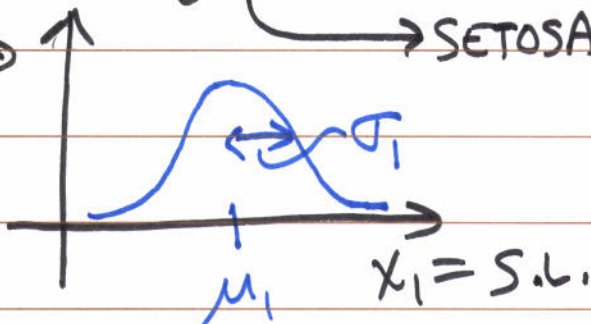Class-conditional density, conditioned on $\underline{\theta}$.

$\left( \pi_y = p(y) \right.$ = prior on $y$ $\left. \right).$

OUR MODEL:

LABELED DATA:

$p(\underline{x} \mid y=1, \underline{\theta}) \longrightarrow$ SETOSA

$p(\underline{x} \mid y=2, \underline{\theta}) \longrightarrow$ VIRGINICA

$\sigma_1$

$x_1 = S.L.$

$\mu_1$

$\mu_2$

$x_1 = S.L.$
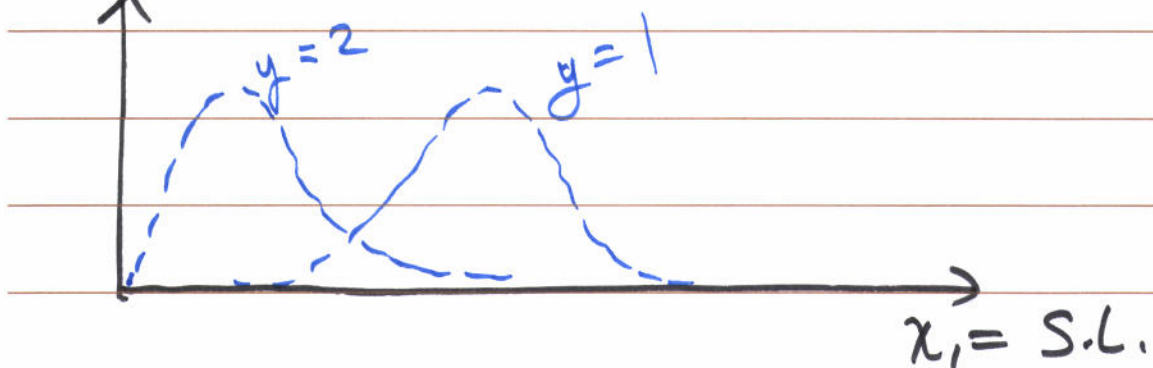
$$\underline{\theta} = \begin{bmatrix} \mu_1 \\ \sigma_1 \\ \mu_2 \\ \vdots \end{bmatrix}$$

UNLABELED DATA:

$\pi_y$

$$(2) \quad p(\underline{x} \mid \underline{\theta}) = \sum_{y=1}^{C} \underbrace{p(\underline{x} \mid y, \underline{\theta})}_{\text{COMPONENT DENS.}} \underbrace{p(y \mid \underline{\theta})}_{\substack{\text{MIXING} \\ \text{PARAMETER}}}$$

$= $ A MIXTURE DENSITY.

$p(\underline{x} \mid \underline{\theta})$

$y=2$ $\qquad$ $y=1$

$x_1 = S.L.$

How use both models (for $\mathcal{D}_L$ and $\mathcal{D}_U$) together?

FIND $\underline{\theta}$ FROM DATA USING MLE

$$\hat{\underline{\theta}}_{MLE} = \underset{\underline{\theta}}{\text{argmax}} \; p(\mathcal{D}|\underline{\theta}) = \underset{\underline{\theta}}{\text{argmax}} \; \ln p(\mathcal{D}|\underline{\theta})$$

$$p(\mathcal{D}|\underline{\theta}) = \prod_{i=1}^{\ell} p(\underline{x}_i, y_i | \underline{\theta}) \cdot \prod_{i=\ell+1}^{u+\ell} p(\underline{x}_i | \underline{\theta})$$

$$\ln p(\mathcal{D}|\underline{\theta}) = \sum_{i=1}^{\ell} \ln p(\underline{x}_i, y_i | \underline{\theta}) + \sum_{i=\ell+1}^{\ell+u} \ln p(\underline{x}_i | \underline{\theta})$$

$$\underbrace{\qquad\qquad}_{\mathcal{D}_L \text{ use } (1)} \qquad \underbrace{\qquad\qquad}_{\mathcal{D}_U \text{ use } (2)}$$

$$= \sum_{i=1}^{\ell} \ln \underbrace{p(\underline{x}_i | \underline{\theta}) + \ln \pi_{y_i}}_{y_i}$$

$$+ \sum_{i=\ell+1}^{u} \ln \left[ \sum_{y=1}^{c} p(\underline{x}_i | y, \underline{\theta}) \, \pi_{y} \right] \qquad \circledast$$

LET $\mathcal{D} \overset{\triangle}{=} \{ \mathcal{D}_L, \mathcal{D}_U \}$ $\underbrace{i = \ell+1, \cdots, u+\ell,}$

TREAT UNKNOWN LABELS $y_i$, ~~the x~~ AS "HIDDEN VARIABLES", DENOTED $\mathcal{H}$.

$\rightarrow$ USE EXPECTATION MAXIMIZATION (EM).

TO EST. $\mathcal{H}$ AND $\underline{\theta}$.

# EM Algorithm (General Formulation)
## [Follows SSL text; also
## Murphy 11.4]

COMMONLY USED FOR:

- WORKING WITH MISSING DATA.
- FINDING MLE IN DIFFICULT SITUATIONS.
- ESTIMATING QUANTITIES IN MIXTURE MODELS.

LET $\mathcal{D}$ BE ALL THE DATA:

$$\{ (\underline{x}_i, y_i), i=1, \cdots, \ell ; \underline{x}_h, h=\ell+1, \cdots, \ell+u \}$$

LET $\mathcal{H}$ BE THE HIDDEN LABELS

$$\{ y_h, h=\ell+1, \cdots, \ell+u \}$$

# EM ALGORITHM ($t$ = ITERATION INDEX)

INITIALIZE $t = 0$ AND $\underline{\theta}^{(0)}$

E STEP:

COMPUTE BEST EST. OF $\mathcal{H}$ AS:
$$p(\mathcal{H} \mid \mathcal{D}, \underline{\theta}^{(t)})$$

$\underline{\theta}$

M STEP:

EST. PARAMETERS $\underline{\theta}^{(t+1)}$ BY:

$$\underline{\theta}^{(t+1)} = \arg\max$$

$$E_{\mathcal{H} \mid \mathcal{D}, \underline{\theta}^{(t)}} \{ \ln p(\mathcal{D}, \mathcal{H} \mid \underline{\theta}) \}$$

(EST. HIDDEN LABELS $y_h$)

(USE $y_h$ EST.'s TO COMPUTE NEW MLE OF $\underline{\theta}$).

$t \leftarrow t+1$

HALT WHEN $p(\mathcal{D} \mid \underline{\theta}^{(t)})$ CONVERGES.

# EM PROPERTIES

1. CAN BE SHOWN THAT $p(\mathcal{D}|\underline{\theta})$ INCREASES AT EVERY ITERATION.

2. CONVERGES TO A <u>LOCAL OPTIMUM</u>.

    Is $-\ln p(\mathcal{D}|\underline{\theta})$ A CONVEX FCN. OF $\underline{\theta}$?

    $\ell$ $\rightarrow$ GENERRALLY, NO.

3. RESULT DEPENDS ON STARTING POINT $\underline{\theta}^{(0)}$.

---

COMMON CHOICE: $\underline{\theta}^{(0)} = \underline{\hat{\theta}}_{MLE}$ BASED ON $\mathcal{D}_L$.

How to USE IT:

### FOR E STEP

LET $i$ INDEX THE DATA PTS. IN $\mathcal{D}_L$;

$\quad h \quad " \quad " \quad " \quad " \quad$ IN $\mathcal{D}_U$.

$$p(\mathcal{H} \mid \mathcal{D}, \underline{\theta}) = \prod_{h=\ell+1}^{\ell+u} p\left(y_h \mid \underline{\underline{X}}=U, \underline{Y}=L, \underline{\underline{X}}=L\right)$$

(3)
$$= \prod_h p\left(y_h \mid x_h, \underline{\theta}\right)^{\underline{\theta}}$$

(KNOW $\underline{\theta}$ $\Rightarrow$) DON'T NEED OTHER DATA.)

(4) $p\left(y_h = c \mid x_h, \underline{\theta}\right) = \dfrac{p\left(x_h \mid y_h = c, \underline{\theta}\right) p\left(y_h \mid \underline{\theta}\right)}{\displaystyle\sum_{y_h=1}^{c} p\left(y_h \mid \underline{\theta}\right) p\left(x_h \mid y_h, \underline{\theta}\right)}$

LET $\gamma_{hc} \overset{\Delta}{=} p\left(y_h = c \mid x_h, \underline{\theta}\right)$

DATA PT. INDEX (UNLABELED)

CLASS ASSIGNMENT

$$\gamma_{hc} = \text{RESPONSIBILITY OF } y_h = c \text{ LABEL}$$

$$\text{FOR DATA PT. } \underline{x}_h.$$

$$= \text{``SOFT LABEL'' FOR DATA PT. } \underline{x}_h.$$

$$p(y_h = c \mid \underline{x}_h, \underline{\theta}).$$

FOR M STEP

$$\max_{\underline{\theta}} \; E_{\mathcal{H} \mid \mathcal{D}, \underline{\theta}^{(t)}} \left\{ \ln p(\mathcal{D}, \mathcal{H} \mid \underline{\theta}) \right\}$$

$$= \max_{\underline{\theta}} \left\{ \sum_{\mathcal{H}} p(\mathcal{H} \mid \mathcal{D}, \underline{\theta}^{(t)}) \ln p(\mathcal{D}, \mathcal{H} \mid \underline{\theta}) \right\}$$

$$p(\mathcal{D}, \mathcal{H} \mid \underline{\theta}) = \underbrace{p(\mathcal{H} \mid \mathcal{D}, \underline{\theta})}_{\substack{\text{GIVEN} \\ \text{ABOVE.} \\ \text{EQ. (3),(4).}}} \underbrace{p(\mathcal{D} \mid \underline{\theta})}_{\substack{\text{LIKELIHOOD OF} \\ \text{ALL KNOWN (GIVEN)} \\ \text{DATA.}}}$$

$$\cancel{p(\mathcal{D} \mid \underline{\theta})}$$

$$p(\mathscr{D} \mid \underline{\theta}) = \left[ \prod_{i=1}^{\ell} p(\underline{x}_i, y_i \mid \underline{\theta}) \right] \prod_{h=\ell+1}^{\ell+u} p(\underline{x}_h \mid \underline{\theta})$$

$$p(\underline{x}_i, y_i \mid \underline{\theta}) = \underbrace{p(\underline{x}_i \mid y_i, \underline{\theta})}_{\substack{\text{MODEL FOR} \\ \text{LABELED DATA.}}} \underbrace{p(y_i \mid \underline{\theta})}_{\pi_{y_i}}$$

MIXTURE DENSITY $\left( \text{EQ. (2)} \right)$:

$$p(\underline{x}_h \mid \underline{\theta}) = \sum_{y_h=1}^{c} p(\underline{x}_h \mid y_h, \underline{\theta}) \, \pi_{y_h} \, . \qquad \circledast$$