

Tuesday, 11/10/2020

EE 660

MACHINE LEARNING
FROM SIGNALS:
FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

Lecture 23

Lecture 23

EE 660**Nov 10, 2020**

Announcements

- Homework 9 will be posted later this week
- End-of-semester quiz
 - Tues., 11/24/2020, 5:30-7:00 PM
 - Will be short-answer questions, taken online
 - Will be 60-90 minutes in length
 - Materials allowed to be announced

Today's topics

- Semi-supervised learning (part 3)
 - Expectation Maximization (EM)
- Unsupervised learning (USL) (part 1)
 - Introduction
 - Mixture models and EM → *deferred*

EM Algorithm

Initialize $t = 0$ and $\underline{\theta}^{(0)}$



E step

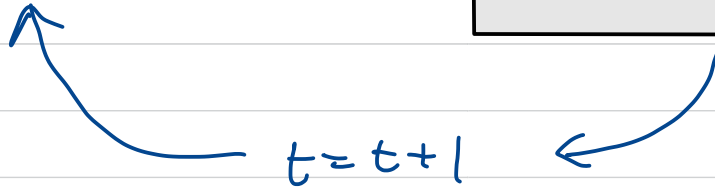
Compute best est.
of \mathcal{H} as $p(\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)})$



M step

Estimate parameters $\underline{\theta}^{(t+1)}$ by:

$$\underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)}} \{ \ln p(\mathcal{D}, \mathcal{H} | \underline{\theta}) \}$$



$t = t + 1$



Halt when $p(\mathcal{D} | \underline{\theta}^{(t)})$ converges

EM Properties

1. Can be shown that $p(\mathcal{D}|\underline{\theta})$ increases at every iteration.
2. Converges to local optimum.
3. Result depends on starting point $\underline{\theta}^{(0)}$.

Common choice is:

$$\underline{\theta}^{(0)} = \hat{\underline{\theta}}_{MLE} \text{ based on } \mathcal{D}_L.$$

How to use it: For E step

Let i index the data pts in \mathcal{D}_L ; h index the data pts. in \mathcal{D}_U .

$$(3) \quad p(\mathcal{H} | \mathcal{D}, \underline{\theta}) = \prod_{h=l+1}^{l+u} p(y_h | \underline{x}_u, \underline{y}_L, \underline{x}_L, \underline{\theta})$$

$$= \prod_h p(y_h | \underline{x}_h, \underline{\theta}) \quad (\text{if know } \underline{\theta}, \text{ don't need other } \underline{x}_n, \underline{x}_L, \text{ or } \underline{y}_L \text{ to predict } y_n).$$

$$(4) \quad p(y_h = c | \underline{x}_h, \underline{\theta}) = \frac{p(\underline{x}_h | y_h = c, \underline{\theta}) p(y_h = c | \underline{\theta})}{\sum_{y_u=1}^C p(y_u = c | \underline{\theta}) p(\underline{x}_h | y_u, \underline{\theta})}$$

π_y

π_{y_u}

Let $r_{hc} \triangleq p(y_h = c | \underline{x}_h, \underline{\theta})$ = responsibility of $y_h = c$ for data pt. \underline{x}_h .

data pt. index $\uparrow \uparrow$ class ass'n.

= "soft label" for \underline{x}_h .

$p(\mathcal{H} | \mathcal{D}, \underline{\theta})$ gives soft labels of all unlabeled data pts.

M step

$$\begin{aligned} \max_{\underline{\theta}} \mathbb{E}_{\mathcal{H}|\mathcal{D}, \underline{\theta}^{(t)}} \{ \ln p(\mathcal{D}, \mathcal{H} | \underline{\theta}) \} \\ = \max_{\underline{\theta}} \left\{ \sum_{\mathcal{H}} p(\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)}) \ln p(\mathcal{D}, \mathcal{H} | \underline{\theta}) \right\} \end{aligned}$$

$$p(\mathcal{D}, \mathcal{H} | \underline{\theta}) = \underbrace{p(\mathcal{H} | \mathcal{D}, \underline{\theta})}_{\substack{\text{given above - (3), (4)} \\ \text{from E step}}} \underbrace{p(\mathcal{D} | \underline{\theta})}_{\substack{\text{likelihood} \\ \text{of all known data}}}$$

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

$$\text{from (2): } p(\underline{x}_h | \underline{\theta}) = \sum_{y=1}^C p(\underline{x}_h | y, \underline{\theta}) \pi_y$$

(a mixture density)

$$p(\underline{x}_i, y_i | \underline{\theta}) = \underbrace{p(\underline{x}_i | y_i, \underline{\theta})}_{\substack{\text{model w/ labeled data} \\ \text{(class-conditional pdf)}}} \underbrace{p(y_i | \underline{\theta})}_{\substack{\text{class prior.}}} \rightarrow$$

Comments: (EM w/ mixture dens. for SSL)

1. Works well when model is \sim correct.
2. Otherwise might not work well.

[Fig. 3.2, 3.3 in Zhu SSL textbook]

Other topics in SSL (N.R.F.)

- Cluster-then-label methods (end of Ch.3)
- Co-training (Ch.4)
 - Each instance x_i has 2 views (feature sets): $x^{(1)}$, $x^{(2)}$.
 e.g.: words/letters in a phrase to classify: $x^{(1)}$
 context (nearby words): $x^{(2)}$
 - Use both for SSL
- Graph-based methods
- SVM based methods
- Bounds on E_{out} [Intro. in Ch.8]

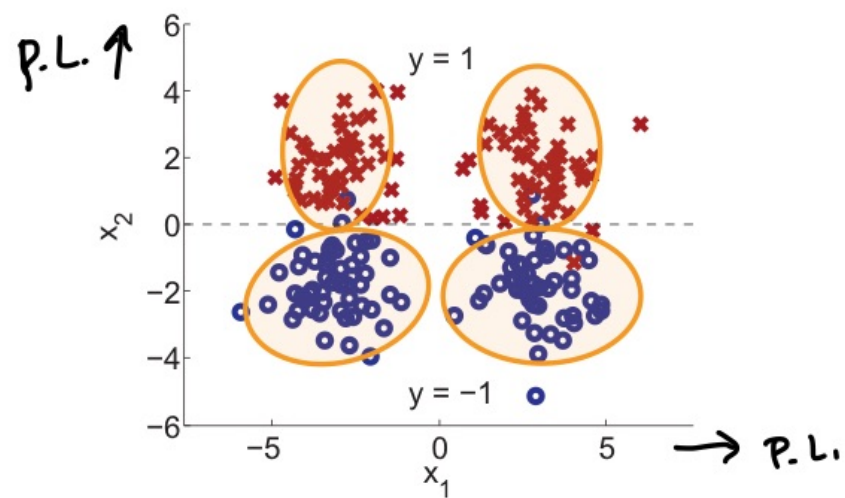


Figure 3.2: Two classes in four clusters (each a 2-dimensional Gaussian distribution).

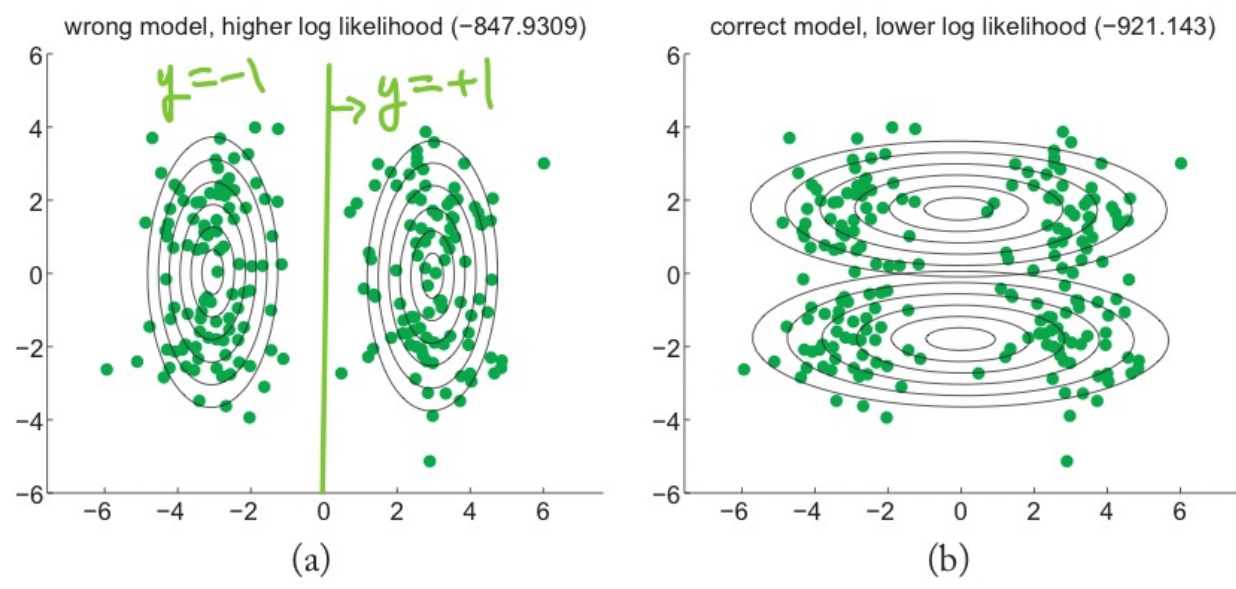


Figure 3.3: (a) Good fit under the wrong model assumption. The decision boundary is vertical, thus producing mass misclassification. (b) Worse fit under the wrong model assumption. However, the decision boundary is correct.

[From zhu et al., Intro.to Semi-Supervised Learning.]

Unsupervised Learning (USL)

Data pts. in \mathcal{D} have no class labels.

What can USL do for us?

- Learn structure in the data — significance of clusters.
- Feature selection or feature discovery
 - e.g.: cluster data
 - find centroid of each cluster μ_k
 - use distance $d(x_i, \mu_k)$ to each μ_k as a feature value.
 - bag-of-words representation
- Adapt a classifier to changes over time by revising cluster/class boundaries as new data is observed.