# EE 660

# MACHINE LEARNING FROM SIGNALS: FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

**Lecture 22**

**Lecture 22**                    **EE 660**                    **Nov 5, 2020**

**Lecture 22 announcements**

-   Homework 8 is due tomorrow.

-   Homework 9 will be posted

**Lecture 22 outline**

-   Semi-supervised learning (SSL)  (part 2)

    • Self-training models

    • Mixture models and parametric classification (SL)

    • Mixture models and parametric classification (SSL)

# Self-Training Models

Use their prediction $\hat{f}(\underline{x}^{(u)})$ for additional training.

[ Self- tr. $\dot{\xi}$ prop INN algorithms]
[ Small ex.]

Variant of prop. INN alg.: prop. kNN — use a kNN for the S.L. module.

[ Prop. INN applied to 100 aliens data]

This alg. tends to work well when data forms $C$ dense, well separated clusters, or $C$ long chains of data that are separated.

# SSL Self-Training Models [Zhu and Goldberg text]

**Algorithm 2.4. Self-training.** *(wrapper)*

*Input: labeled data $\{(\mathbf{x}_i, y_i)\}_{i=1}^{l}$, unlabeled data $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$.*

*1. Initially, let $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^{l}$ and $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$.*

*2. Repeat:*

*3.      Train $\hat{f}$ from L using supervised learning.*

*4.      Apply $\hat{f}$ to the unlabeled instances in U.*

*5.      Remove a subset S from U; add $\{(\mathbf{x}, \hat{f}(\mathbf{x})) | \mathbf{x} \in S\}$ to L.*

S: data points with
highest confidence
of $\hat{f}$ prediction.

Assumption: data pts. with
highest confidence of
$\hat{f}(\underline{x})$ tend to be
correct.

*Specific example:*

**Algorithm 2.7. Propagating 1-Nearest-Neighbor.**

*Input: labeled data $\{(\mathbf{x}_i, y_i)\}_{i=1}^{l}$, unlabeled data $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$, distance function $d()$.*

*1. Initially, let $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^{l}$ and $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$.*

*2. Repeat until U is empty:*

pt. in
S

*3.      Select $\mathbf{x} = \operatorname{argmin}_{\mathbf{x} \in U} \min_{\mathbf{x}' \in L} d(\mathbf{x}, \mathbf{x}')$.* ← find closest 2 pts; one from L and one from U.
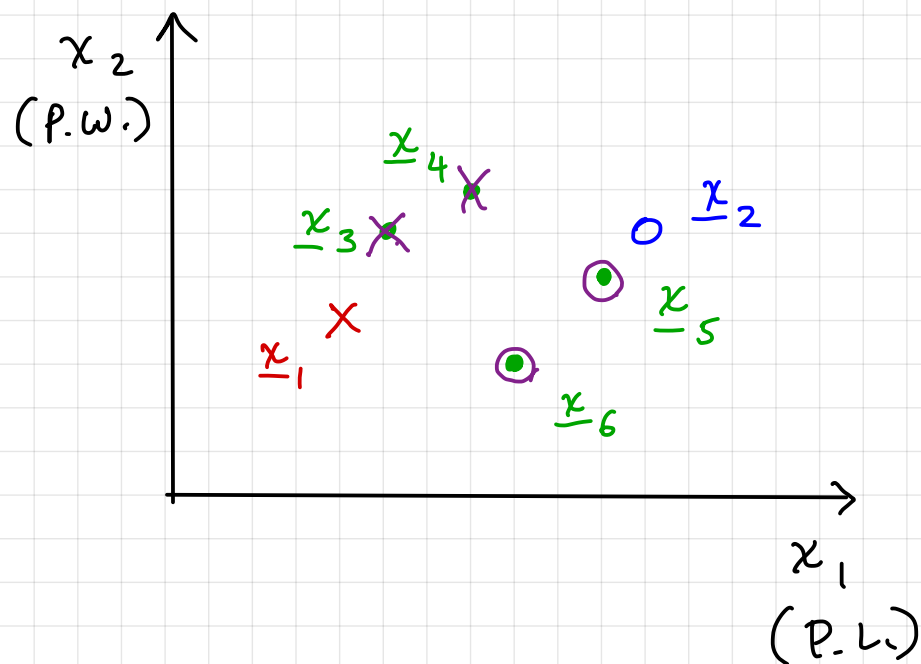
*4.      Set $\hat{f}(\mathbf{x})$ to the label of $\mathbf{x}$'s nearest instance in L. Break ties randomly.*

*5.      Remove $\mathbf{x}$ from U; add $(\mathbf{x}, f(\mathbf{x}))$ to L.*

Excerpts are from Xiaojin Zhu and Andrew B. Goldberg, *Introduction to Semi-Supervised Learning* (Morgan and Claypool, 2009)

# Propagating 1-NN example:



$x_2$ (P.W.)

$x_4$  $x_3$  $x_2$  $x_5$  $x_1$  $x_6$

$x_1$ (P.L.)

Legend:
- X — $y = +1$ (setosa)
- O — $y = -1$ (virginica)
- ● — unlabeled

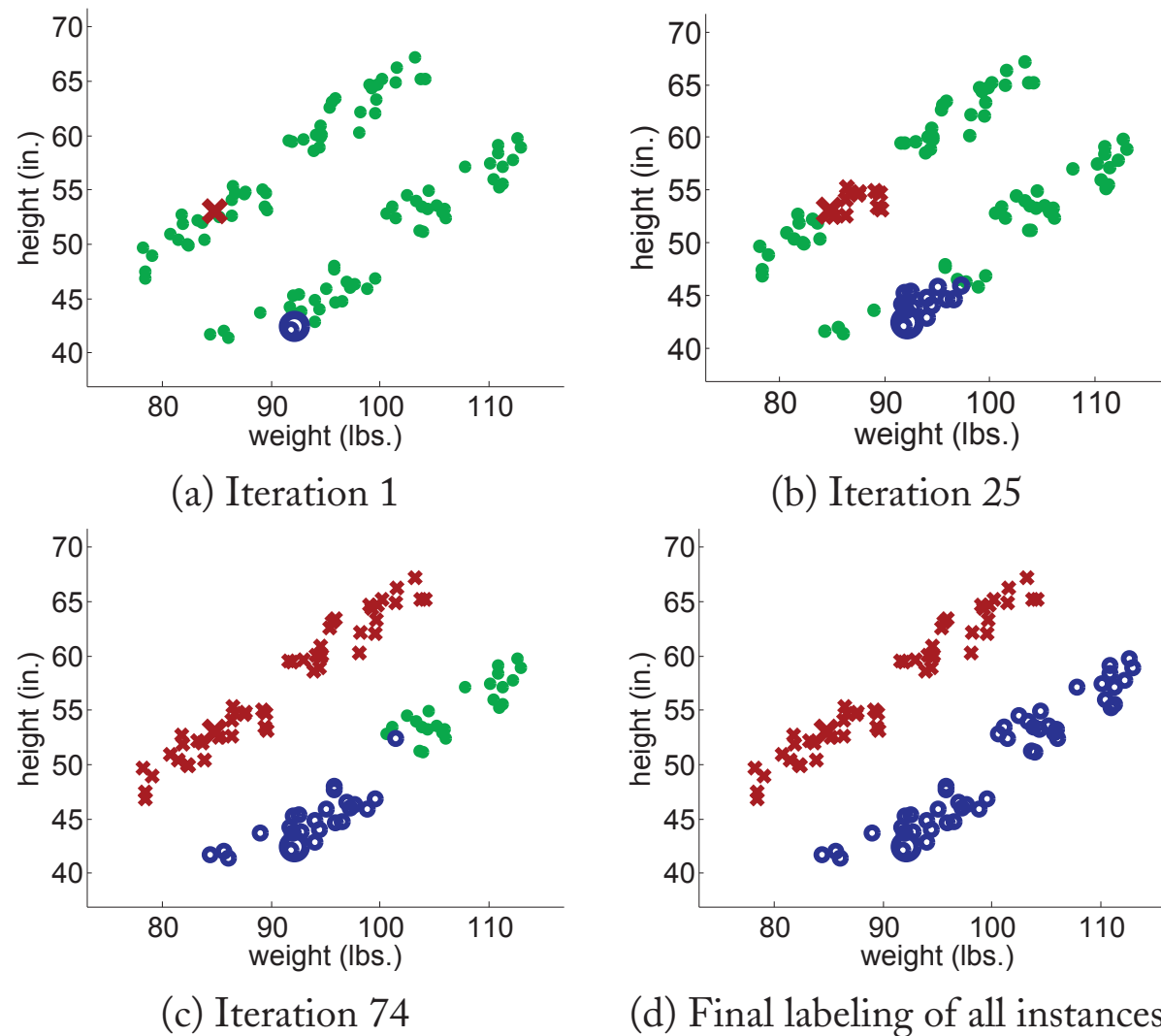| iteration # | closest unlabeled pt. $\underline{x}$ | $\hat{f}(\underline{x})$ | L | U |
|---|---|---|---|---|
| 1 | $\underline{x}_5$ | $-1$ | $\underline{x}_1, \underline{x}_2, (\underline{x}_5, -1)$ | $\underline{x}_3, \underline{x}_4, \underline{x}_6$ |
| 2 | $\underline{x}_3$ | $+1$ | $\underline{x}_1, \underline{x}_2, (\underline{x}_5, -1), (\underline{x}_3, +1)$ | $\underline{x}_4, \underline{x}_6$ |
| 3 | $\underline{x}_4$ | $+1$ | $\underline{x}_1, \underline{x}_2, (\underline{x}_5, -1), (\underline{x}_3, +1), (\underline{x}_4, +1)$ | $\underline{x}_6$ |
| 4 | $\underline{x}_6$ | $-1$ | all | ∅ |

(a) Iteration 1

(b) Iteration 25

(c) Iteration 74

(d) Final labeling of all instances

**Figure 2.3:** Propagating 1-nearest-neighbor applied to the 100-little-green-alien data.

Excerpts are from Xiaojin Zhu and Andrew B. Goldberg, *Introduction to Semi-Supervised Learning* (Morgan and Claypool, 2009)

# Mixture Models and Parametric Classification — Supervised Learning (SL)

Suppose we model $p(\underline{x}|y)$ as a pdf with some unknown parameters, e.g.:

$$p(\underline{x}|y, \underline{\theta}) = N(\underline{x}|\underline{\mu}_y, \underline{\underline{\Sigma}}_y)$$

$y = $ class index

$\underline{\mu}_y$ and $\underline{\underline{\Sigma}}_y$, $y = 1, 2, \cdots, C$, are unknown. $\Rightarrow$ parameters $\underline{\theta}$.

Posterior predictive $p(y|\underline{x}) = ?$

$$p(y|\underline{x}) = \frac{p(\underline{x}|y)\, p(y)}{p(\underline{x})}$$

prior

$$p(\underline{x}) = \sum_{y'} p(\underline{x}|y')\, p(y')$$

$= $ a mixture density.

To find $p(\underline{x}|y)$, how can estimate $\underline{\theta}$? (in SL)

One way: MLE

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

$$\hat{\theta}_{MLE} = \underset{\theta}{argmax} \sum_{i=1}^{\ell} \ln p(\underline{x}_i, y_i | \underline{\theta})$$

$$= \underset{\theta}{argmax} \sum_{i=1}^{\ell} \ln \left[ p(\underline{x}_i | y_i, \underline{\theta}) \, p(y_i | \underline{\theta}) \right]$$

$\underbrace{\phantom{p(y_i | \underline{\theta})}}$ prior on $y = \pi_y$, $y = $ class index.

$\rightarrow$ For a given fcn. $p(\underline{x} | y, \underline{\theta})$, solution gives $\hat{\underline{\theta}}_{MLE}$.

Given our estimates $\hat{\underline{\theta}}$, we have:

$$p(\underline{x}) = \sum_{y'} p(\underline{x} | y', \hat{\underline{\theta}}) \, p(y' | \hat{\underline{\theta}})$$

Let $\pi_{y'} = p(y') = p(y' | \hat{\underline{\theta}})$

$$\therefore \quad p(\underline{x}) = \sum_{y'} \pi_{y'} \, p(\underline{x} | y', \hat{\underline{\theta}}) = \text{a mixture density.}$$

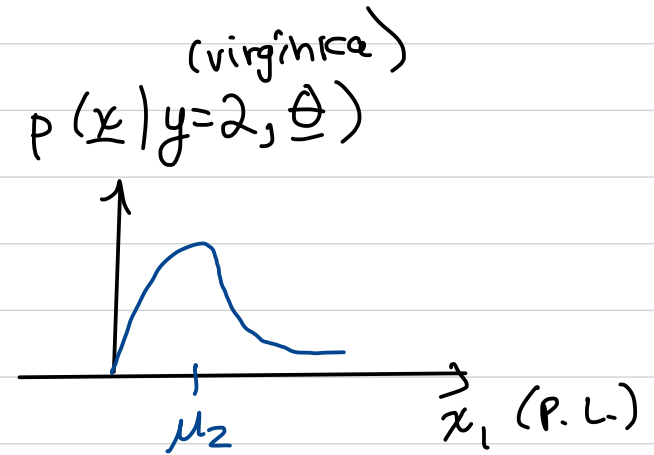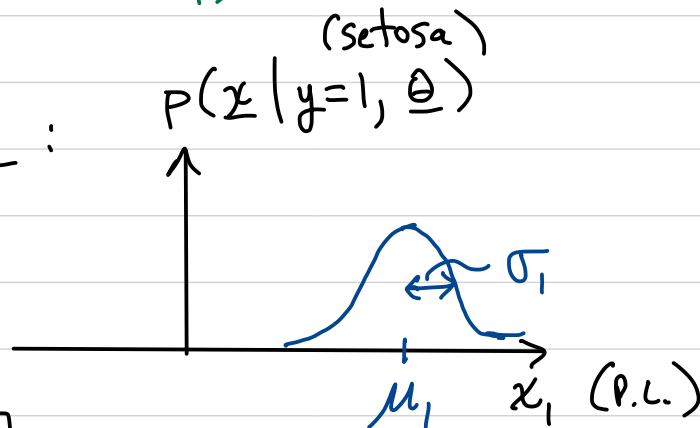# Mixture Models for SSL

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

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

(1) $\quad p(\underline{x}, y \mid \underline{\theta}) = p(\underline{x} \mid y, \underline{\theta}) \, p(y \mid \underline{\theta}) \qquad$ (assuming $\underline{\theta}$ is not random)

$$= \underbrace{p(\underline{x} \mid y, \underline{\theta})}_{\substack{\text{class-conditional} \\ \text{density, conditioned on } \underline{\theta}.}} \, p(y)$$
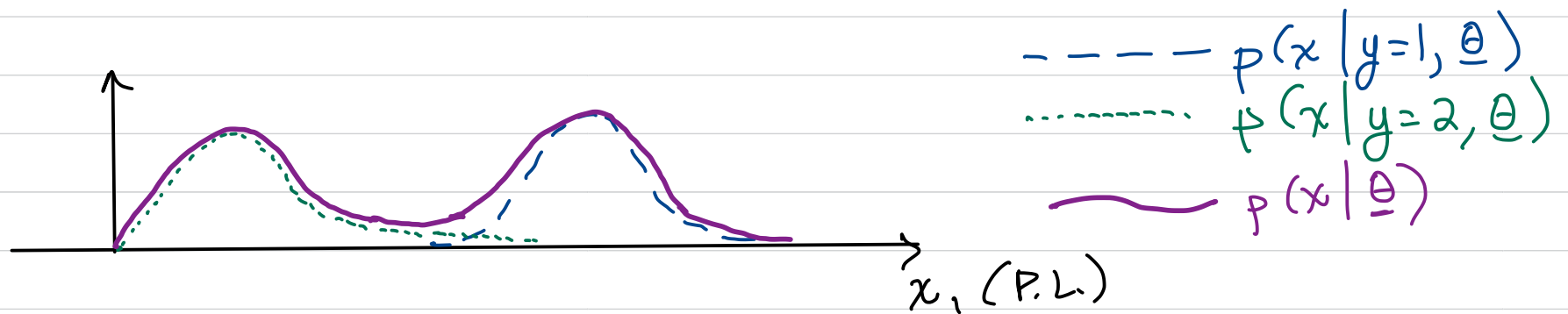
$$\curvearrowright \pi_y$$

Our model:
Labeled data $\mathcal{D}_L$:

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

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



$\sigma_1$

$\mu_1 \quad x_1 \text{ (P.L.)}$

$\mu_2 \quad x_1 \text{ (P.L.)}$

$$\underline{\theta} = \begin{bmatrix} \mu_1 \\ \sigma_1 \\ \mu_2 \\ \vdots \end{bmatrix}$$

Unlabeled data $\mathcal{D}_U$ :

(2) $\quad p(\underline{x}|\underline{\theta}) = \sum\limits_{y=1}^{c} \underbrace{p(\underline{x}|y,\underline{\theta})}_{\substack{\text{component}\\\text{density}}} \underbrace{\overbrace{p(y|\underline{\theta})}^{\pi_y}}_{\substack{\text{mixing}\\\text{parameter}}} = a \text{ mixture density}$



------ $p(x|y=1,\underline{\theta})$

........ $p(x|y=2,\underline{\theta})$

～～ $p(x|\underline{\theta})$

$x_1$ (P.L.)

$\Rightarrow$ How do we use both models (for $\mathcal{D}_L$ and $\mathcal{D}_U$) together to estimate $\underline{\theta}$?

Find $\underline{\theta}$ from data using MLE

(3)  $\hat{\underline{\theta}}_{MLE} = \underset{\underline{\theta}}{argmax}\ p(\mathcal{S}|\underline{\theta}) = \underset{\underline{\theta}}{argmax}\ \ln p(\mathcal{S}|\underline{\theta})$

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

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

(4) $\ln p(\mathcal{S}|\underline{\theta}) = \sum_{i>1}^{\ell} \left[\ln p(\underline{x}_i|y_i, \underline{\theta}) + \ln p(y_i|\underline{\theta})\right]$  $\Big\}\, \mathcal{S}_L$

$\qquad\qquad + \sum_{i=\ell+1}^{u} \ln\left[\sum_{y=1}^{C} p(\underline{x}_i|y, \underline{\theta})\underbrace{p(y|\underline{\theta})}_{\pi_y}\right]$  $\Big\}\, \mathcal{S}_U$

Let $\mathcal{S} \overset{\Delta}{=} \{\mathcal{S}_L, \mathcal{S}_U\}$;

and treat the unknown $y_i$ in $\mathcal{S}_U$ as "hidden variables", denoted $\mathcal{H}$.

$\longrightarrow$ If we knew $\mathcal{H}$, we could: $\underline{\theta}_{MLE} = \underset{\underline{\theta}}{argmax}\{\ln p(\mathcal{S}, \mathcal{H}|\underline{\theta})\}$
(as in SL). But, don't know $\mathcal{H}$.

$\longrightarrow$ Use Expectation maximization (EM) to estimate $\mathcal{H}$ and $\underline{\theta}$.