# EE 660

# MACHINE LEARNING FROM SIGNALS: FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

**Lecture 4**

# Lecture 4      EE 660      Sep 3, 2020

## Announcements

- Class projects and quizzes

  - There will be 1 end-of-semester quiz and no midterm quiz

  - There will be 2 projects - a (smaller) midterm project and a (larger) final project

  - The above were preferred by a large margin of students

- End-of-semester quiz

  - Tuesday, 11/24/2020,  5:30 - 7:00 PM

  - (will end earlier if the quiz is shorter than 90 min.)

  - This time has no remaining known conflicts

  - If you have an unmovable conflict, please email me asap.

- Final project

  - Will be due on Thur., 12/3/2020

- First homework will be posted this Friday

- From now on, most reading assignments will be given in each homework

## Today's Lecture

- Data notation

- Comment on definition of dataset D

  - y|x or y,x

- MLE Regression (part 2)

- Ridge regression

4

# Dataset Notation (1)

(2) $\quad \mathcal{D} = \{\underline{x}_i, y_i\}_{i=1}^{N}$. $\qquad$ Later: $\mathcal{D}_{Tr}$ with $N = N_{Tr}$ $\quad$ (training set)

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\mathcal{D}_{val}$ with $N = N_{val}$ $\quad$ (validation set)

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\mathcal{D}_{Test}$ with $N = N_{Test}$ $\quad$ (test set)

Design matrix: $\quad \underline{\underline{X}} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}_{N \times D}$ $\qquad$ Vector of data-point labels $\left.\right\}$ $\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

We use $\mathcal{D}$ for both set notation as in (2), and matrix-vector pair as in:

$$\mathcal{D} = \underline{\underline{X}}, \underline{y} \quad \text{or} \quad \underline{y}, \underline{\underline{X}}$$

$$\therefore \quad p(\mathcal{D} | \underline{\theta}) = p(\underline{y}, \underline{\underline{X}} | \underline{\theta})$$

but what we often want is $p(\underline{y} | \underline{\underline{X}}, \underline{\theta})$ or $p(\underline{\underline{X}} | \underline{y}, \underline{\theta})$.

**EE 660**  **Notes on** $p(\mathcal{D}|\underline{\theta})$, $p(\underline{y},\underline{\underline{X}}|\underline{\theta})$, $p(\underline{y}|\underline{\underline{X}},\underline{\theta})$  **Fall 2020**

Using probability relations we have:
$$p(\underline{y},\underline{x}|\underline{\theta}) = p(\underline{y}|\underline{x},\underline{\theta})\, p(\underline{x}|\underline{\theta}) = p(\underline{y}|\underline{x},\underline{\theta})\, p(\underline{x})$$

where for the last step we have dropped the last condition on $\underline{\theta}$ because it tells us nothing useful about $p(\underline{x})$. If we are interested in maximizing (or minimizing) the likelihood, we will take:

$$\arg\max_{\underline{\theta}} p(\mathcal{D}|\underline{\theta}) = \arg\max_{\underline{\theta}}\left\{\prod_{i=1}^{N} p(y_i,\underline{x}_i|\underline{\theta})\right\}$$

$$= \arg\max_{\underline{\theta}}\left\{\prod_{i=1}^{N} p(y_i|\underline{x}_i,\underline{\theta})\, p(\underline{x}_i)\right\} = \arg\max_{\underline{\theta}}\left\{\left(\prod_{i=1}^{N} p(\underline{x}_i)\right)\prod_{i=1}^{N} p(y_i|\underline{x}_i,\underline{\theta})\right\}$$

$$= \arg\max_{\underline{\theta}}\left\{\prod_{i=1}^{N} p(y_i|\underline{x}_i,\underline{\theta})\right\}$$

*Constant of $\underline{\theta}$*
*$\geq 0$*

and to obtain the last line, $\displaystyle\prod_{i=1}^{N} p(\underline{x}_i)$ was dropped because it is a positive multiplicative term that is a constant of $\underline{\theta}$. This can equivalently be seen by using the log likelihood instead:

$$\arg\max_{\underline{\theta}} p(\mathcal{D}|\underline{\theta}) = \arg\max_{\underline{\theta}}\left\{\ln\prod_{i=1}^{N} p(y_i,\underline{x}_i|\underline{\theta})\right\}$$

$$= \arg\max_{\underline{\theta}}\left\{\ln\prod_{i=1}^{N} p(y_i|\underline{x}_i,\underline{\theta})\, p(\underline{x}_i)\right\}$$

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

$$= \arg\max_{\underline{\theta}}\left\{\sum_{i=1}^{N}\left[\ln p(y_i|\underline{x}_i,\underline{\theta}) + \ln p(\underline{x}_i)\right]\right\}$$

$$= \arg\max_{\underline{\theta}}\left\{\sum_{i=1}^{N}\ln p(y_i|\underline{x}_i,\underline{\theta})\right\}$$

and to obtain the last line, the additive terms that don't depend on $\underline{\theta}$ have been dropped.

So, when the goal of using the likelihood is to find its argmax or argmin w.r.t. $\underline{\theta}$, we can replace $p(y_i,\underline{x}_i|\underline{\theta})$ directly with $p(y_i|\underline{x}_i,\underline{\theta})$.

# 2. Objective function

$$J_1(\underline{w}, \sigma) = ? = -\ln p(\mathcal{D}|\underline{w}) = NLL(\underline{w})$$
$$(\text{or } -p(\mathcal{D}|\underline{w}))$$

$$\ln p(\mathcal{D}|\underline{w}) = \sum_{i=1}^{N} \ln p(y_i|\underline{x}_i, \underline{w}_i) \qquad \text{(i.d.)}$$

$$= \sum_{i=1}^{N} \ln N(y_i|\underline{w}^T\underline{x}_i, \sigma^2) \qquad \text{(Gaussian model)}$$

Can re-write $J_1$ as: $\qquad [\text{M Eq. 7.5-7.9}]$

$$J_1(\underline{w}, \sigma) = \frac{1}{2\sigma^2} \sum_{i=1}^{N} (y_i - \underline{w}^T\underline{x}_i)^2 + \frac{N}{2} \ln(2\pi\sigma^2)$$
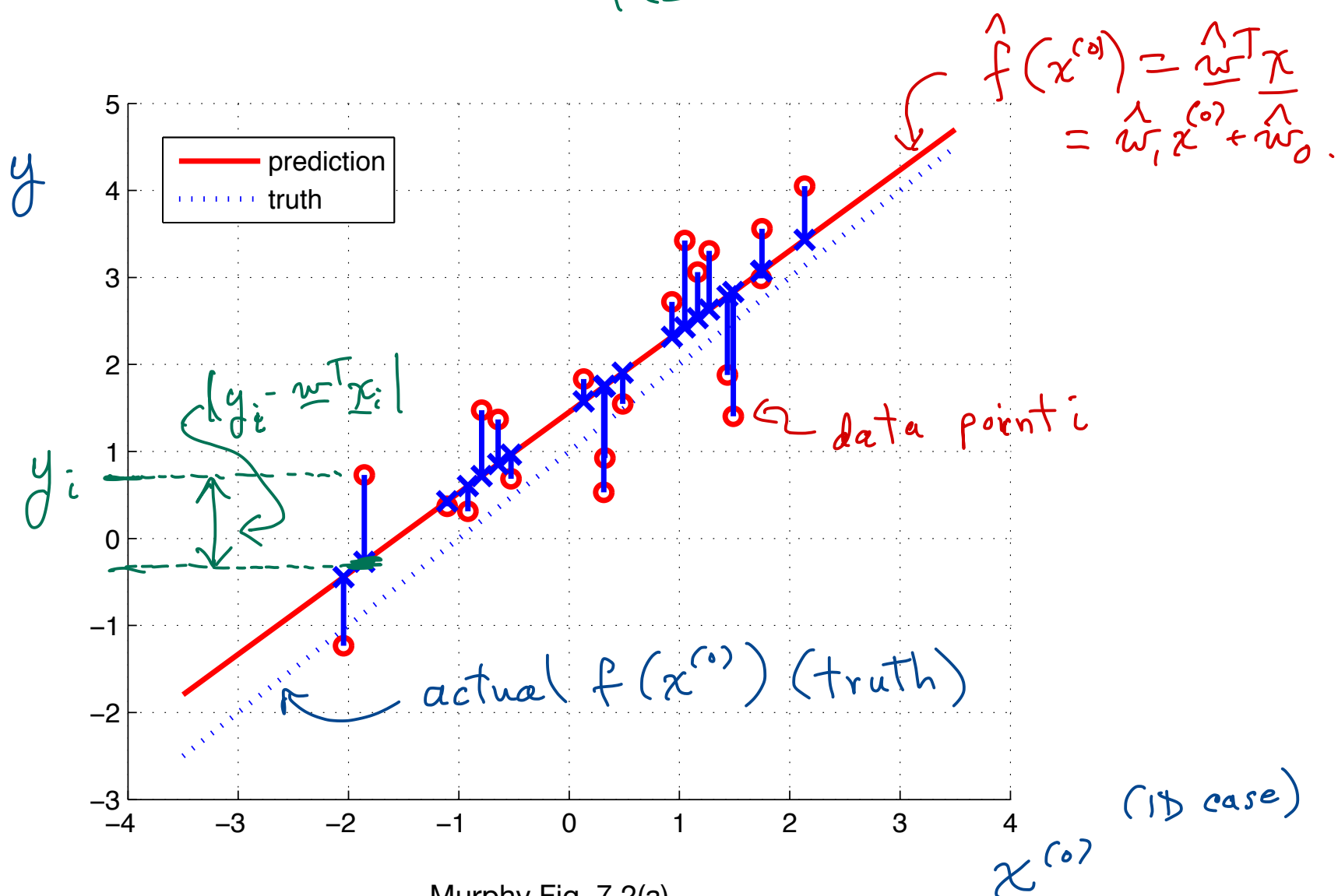
Simplify:
$\underbrace{\qquad}$
multiplicative
constant $> 0$

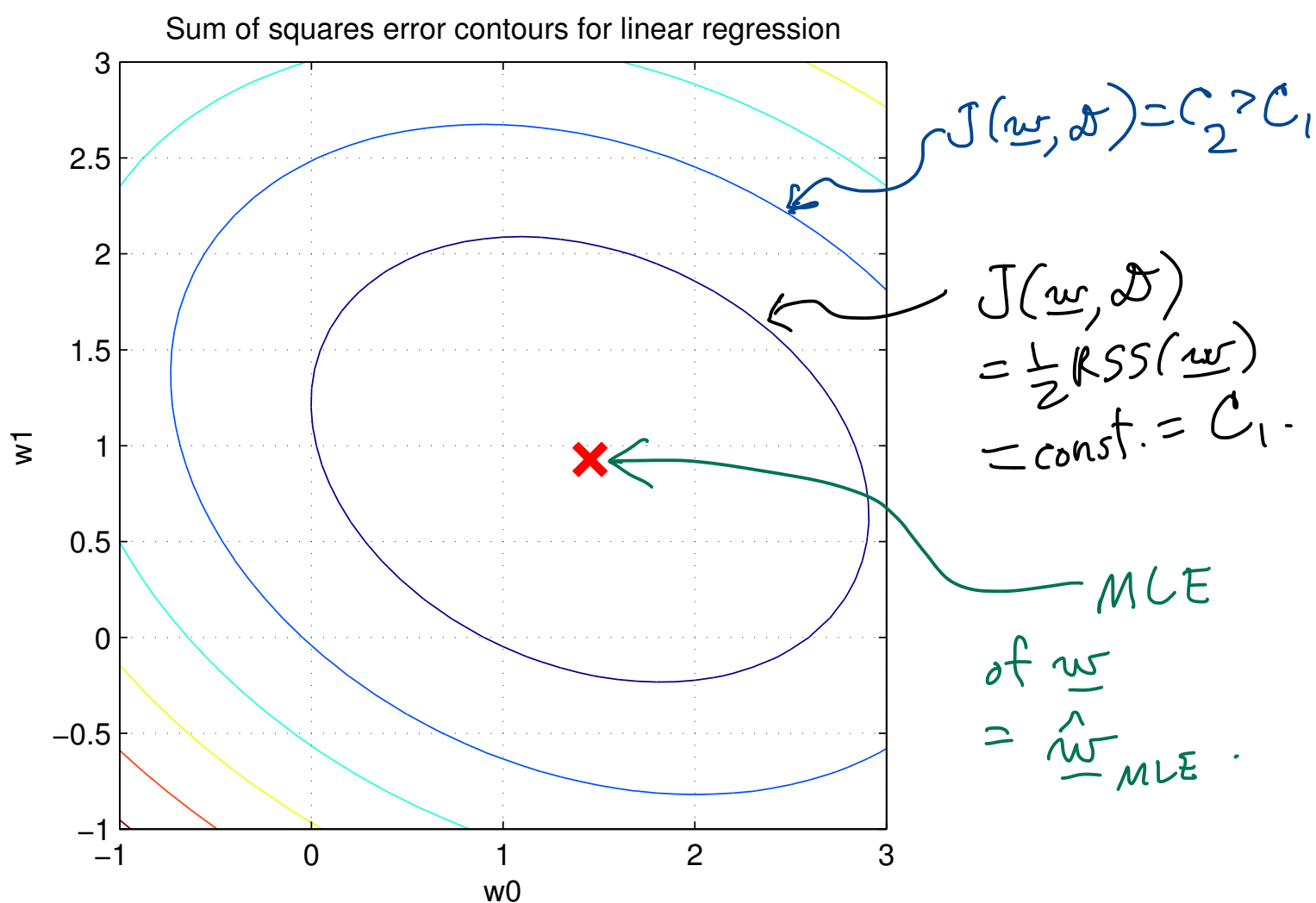$\underbrace{\qquad}$
constant of $\underline{w}$

$$\Rightarrow \boxed{\text{Let } J(\underline{w}, \sigma) = \frac{1}{2} RSS(\underline{w}) = \frac{1}{2} \sum_{i=1}^{N} (y_i - \underline{w}^T\underline{x}_i)^2 = \frac{1}{2}\|\underline{y} - \underline{X}\underline{w}\|_2^2}$$

$$\hat{y} \sim P(\hat{y} | \underline{x}, \underline{w}) = N(\hat{y} | \underbrace{\underline{w}^T \underline{x}}, \sigma^2)$$

$$\hat{f}(\underline{x}) = \hat{\underline{w}}^T \underline{x}$$

$$\hat{f}(\underline{x}^{(0)}) = \hat{\underline{w}}^T \underline{x}$$
$$= \hat{w}_1 x^{(0)} + \hat{w}_0.$$

$|y_i - \underline{w}^T \underline{x}_i|$

$y_i$

$y$

data point $i$

actual $f(x^{(0)})$ (truth)

$x^{(0)}$ (1D case)

Murphy Fig. 7.2(a)

Sum of squares error contours for linear regression

$J(\underline{w}, \mathscr{D}) = C_2 > C_1$

$J(\underline{w}, \mathscr{D})$
$= \frac{1}{2} RSS(\underline{w})$
$= const. = C_1.$

MLE of $\underline{w}$
$= \hat{\underline{w}}_{MLE}.$

w1

w0

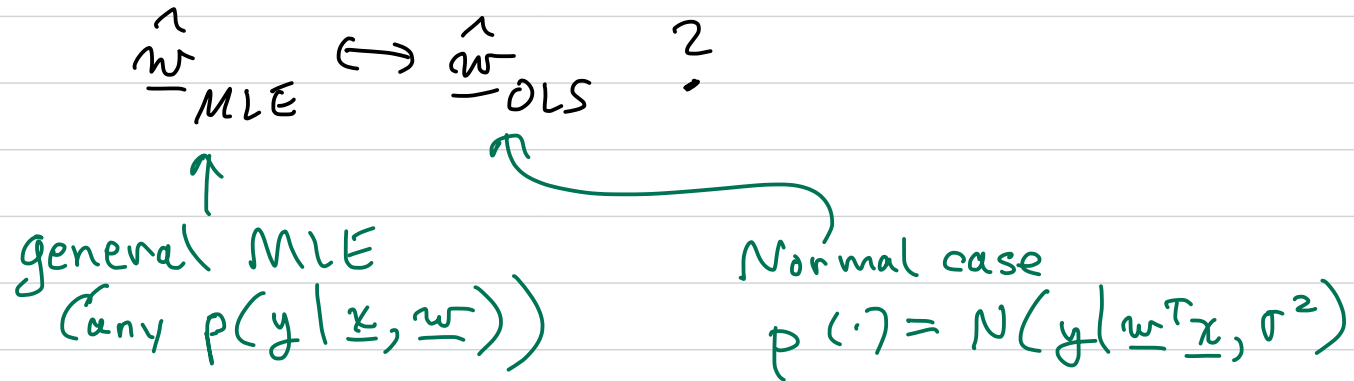Murphy Fig. 7.2(b)

# 3. Optimization method

Which method ?

- Gradient descent (stochastic or batch ...)
  → neural-network approach
  → very large dataset

- Solving $\nabla_{\underline{w}} J(\underline{w}, \mathcal{A}) = \underline{0}$    algebraically.

  → Pseudoinverse solution.
  → Non-neural approaches

  → Solving gives $\hat{\underline{w}}$

  $$\underline{\underline{X}}^T \underline{\underline{X}} \, \hat{\underline{w}} = \underline{\underline{X}}^T \underline{y}$$

  if $(\underline{\underline{X}}^T \underline{\underline{X}})$ is invertable, then:

  $$\boxed{\hat{\underline{w}}_{OLS} = \underline{\underline{X}}^{-} \underline{y} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{y}}$$

= ordinary least squares solution.
= pseudoinv. sol'n.

$$\hat{\underline{w}}_{MLE} \hookrightarrow \hat{\underline{w}}_{OLS} \quad ?$$

general MLE
(any $p(y \mid \underline{x}, \underline{w})$)

Normal case
$$p(\cdot) = N(y \mid \underline{w}^T \underline{x}, \sigma^2)$$

# Ridge Regression   [Murphy 7.5]

$(\underline{w} = \underline{w}^{(0)})$

$\rightarrow$ Use MAP (maximum a posteriori) estimation of $\underline{w}$.

Which is the MAP est. of $\underline{\theta}$?

(a)  $\hat{\underline{\theta}} = \underset{\theta}{argmax}\ p(\mathcal{D}|\underline{\theta})$

(b)  $\hat{\underline{\theta}} = \underset{\theta}{argmax}\ p(\underline{\theta}|\mathcal{D})$   $\leftarrow$

## Model is same as MLE regression:

$$
\begin{cases}
y \sim p(y|\underline{x}, \underline{\theta}) \\
\quad = N(y | \underline{w}^T \underline{x}, \sigma^2) \qquad -\ \text{Linear, Gaussian} \\
\text{or} \quad = N(y | \underline{w}^T \underline{\phi}(\underline{x}), \sigma^2) - \text{Nonlinear, Gaussian}
\end{cases}
$$

Hypothesis set (linear, Gaussian case with $\sigma^2$ given):

$\mathbb{Z}^{\geq 0}$ set of = positive integers

$$\mathcal{H} = \left\{ \hat{y}(\underline{x}) \sim N\left(\hat{y} | \underline{w}^T \underline{x}, \sigma^2\right) \ \middle|\ \underline{w} \in \mathbb{R}^{D+1}, \ D \in \mathbb{Z}^{\geq 0} \right\}$$

Is $\mathcal{H}$ the same as for MLE regression (linear, Gaussian case, $\sigma^2$ given)?

$\rightarrow$ Yes.

What's different here?    $\rightarrow$ Objective function

$$\hat{\underline{\theta}} = \underset{\theta}{\text{argmax}} \; p(\underline{\theta}\,|\,\mathcal{D}) = \underset{\theta}{\text{argmax}} \left\{ \frac{p(\mathcal{D}|\underline{\theta})\,p(\underline{\theta})}{p(\mathcal{D})} \right\}$$

$p(\mathcal{D})$ $\leftarrow$ const. of $\underline{\theta}$.

$$= \underset{\theta}{\text{argmax}} \left\{ p(\mathcal{D}|\underline{\theta})\,p(\underline{\theta}) \right\}$$

$\underbrace{\qquad\qquad}$ always $\geq 0$.

(1)  $$\hat{\underline{\theta}} = \underset{\theta}{\text{argmax}} \left\{ \ln p(\mathcal{D}|\underline{\theta}) + \ln p(\underline{\theta}) \right\}$$

$\underbrace{\qquad}$ likelihood of $\underline{\theta}$ (same as MLE)

$\underbrace{\qquad}$ prior for $\underline{\theta}$

(2)  $$\ln p(\mathcal{D}|\underline{\theta}) = \sum_{i=1}^{N} \ln N\left(y_i \,\middle|\, w_0 + \underline{w}^T \underline{x}_i,\; \sigma^2\right)$$