

Logistic Regression:

- Borrow two big ideas from Gaussian Naive Bayes:

① It is nice to model $P(y|x)$ using just the sigmoid function

$$P(y|x) = \frac{1}{1 + e^{-y(w^T x)}}, \quad y \in \{1, -1\}$$

② it is nice to view w as a parameter of $P(y|x)$, that we can fit it into data

- Discriminative algorithm (Naive Bayes counter part)

- Assumption

$$P(y|x) = \frac{1}{1 + e^{-y(w^T x)}}$$

w is the parameter that we need to fit/estimate

- Estimating w :

MLE estimate: $w_{MLE} = \underset{w}{\operatorname{argmax}} P(y|x; w)$; เลือก w ที่ทำให้ $P(y|x)$ สูงสุด

$$= \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y_i|x_i; w)$$

$$w_{MLE} = \underset{w}{\operatorname{argmax}} \prod_{i=1}^n \frac{1}{1 + e^{-y_i(w^T x_i)}}$$

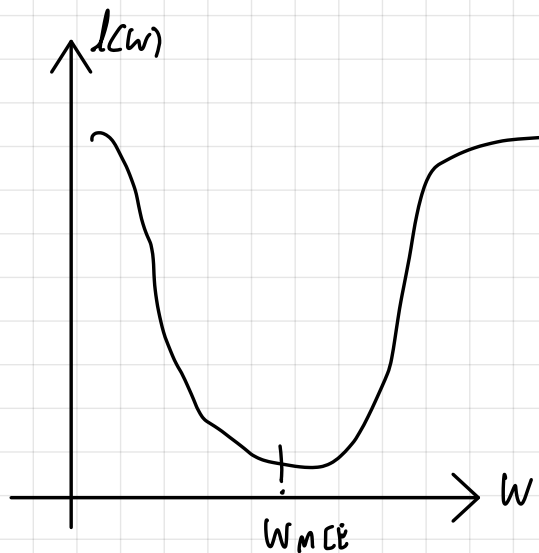
$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \log_e \left(\frac{1}{1 + e^{-y_i(w^T x_i)}} \right)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \log \frac{1}{1 + e^{-y_i(w^T x_i)}}$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log_e (1 + e^{-y_i(w^T x_i)}) \rightarrow \text{negative log likelihood } l(w)$$

Note - There is no closed-form solution to w_{MLE}
 - We will use "Gradient Descent" to find w_{MLE} ,
 that is, find w to minimize

Negative log likelihood $l(w) = \sum_{i=1}^n \log_e (1 + e^{-y_i(w^T x_i)})$

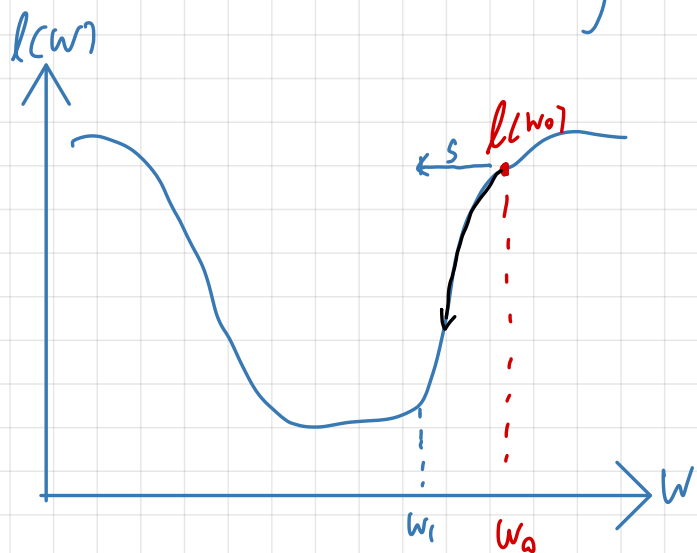


fact: $l(w)$ is a continuous, convex, and differentiable function
 ความต่อเนื่อง

การหาจุดต่ำสุดของ function ที่มี 3 ลักษณะ = convex, differentiable, continuous ***

Gradient Descent:

- Based on "Hill-climbing" scheme



Hill-climbing

- 1) start off by some point w_0
- 2) Repeat until convergence:

$$w_{t+1} = w_t + s$$

if $\|w_{t+1} - w_t\|_2 < \epsilon (0.1 \dots)$, convergence!

The problem is how to define s ?

- In Gradient Descent, we set

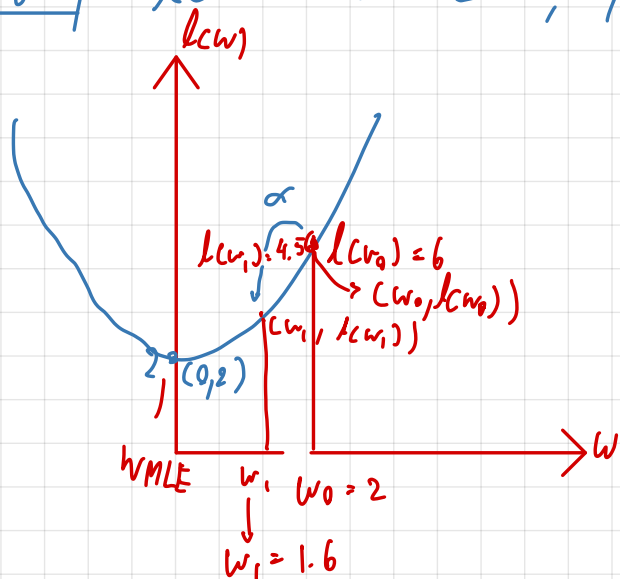
$$S = -\alpha \nabla l(w)$$

learning rate

$\nabla l(w) \approx \frac{dl}{dw}$ dimension of w (ตัวแปรตัว)
 gradient (ขนาด dimension)

Setting

let say $l(w) = w^2 + 2$; $y = x^2 + 2$



$$w_0 = 2 \Rightarrow l(2) = 6$$

$$S = -\alpha \frac{dl(w)}{dw} = -\alpha (2w) \\ = -2\alpha w \\ = -4\alpha$$

$$\text{set } \alpha = 0.1 \Rightarrow S = -0.4$$

update

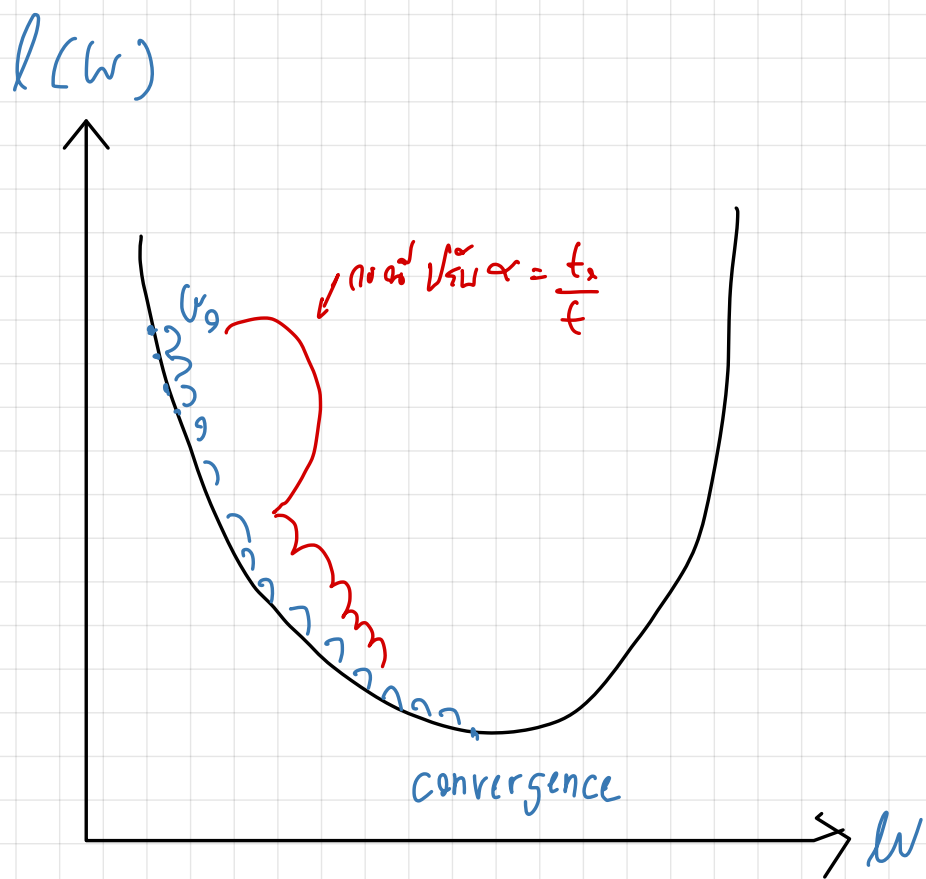
$$w_1 = w_0 + S = 2 + (-0.4) = 1.6$$

$$l(w_1) = (1.6)^2 + 2 = 4.56$$

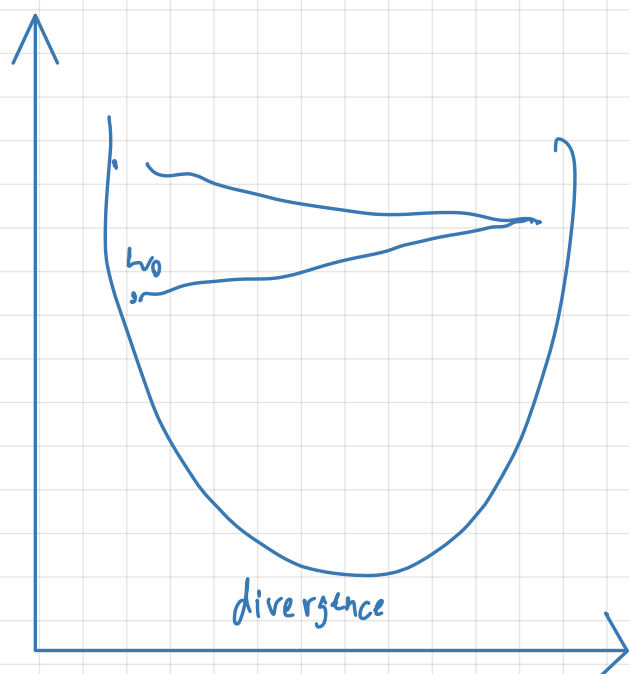
$\alpha > 0$ is a scalar called "learning rate"

$\nabla l(w) = \frac{dl}{dw}$ \rightarrow เปรียบเทียบค่าที่หาค่าต่ำสุดที่ w_{MLE} (ค่าจริงที่หาค่าต่ำสุด) และค่าที่หาค่าต่ำสุด w_{MLE}
 \rightarrow เปรียบเทียบค่าที่หาค่าต่ำสุด $l(w_1) \approx WMLE$

α : เปรียบเทียบ step หรือค่าที่หาค่าต่ำสุดที่หาค่าต่ำสุด step ถ้าให้ค่าที่หาค่าต่ำสุดที่หาค่าต่ำสุด w_{MLE} ได้
 ถ้าให้ค่าที่หาค่าต่ำสุดที่หาค่าต่ำสุด w_{MLE}



- If α is too small, the algorithm will find w_{MLE} very slowly.



- If α is too large, the algorithm will never find w_{MLE}

Best practice: A safer choice is to set $\alpha = \frac{t_0}{t}$, which guarantees that it will eventually become small enough to converge (for any $t_0 \geq 0$)
 \rightarrow เปรียบเทียบค่าที่หาค่าต่ำสุด w

Prove guarantess update S

Taylor's Expansion:

- if $\|s\|_2$ is small (meaning $w+s$ is close to w) then the following holds:

$$\boxed{l(w+s) \approx l(w) + \underbrace{\nabla l(w)^T}_{\text{vector} \times \text{vector} = \text{scalar}} s}$$
$$\nabla l(w) = \begin{bmatrix} \frac{\partial l}{\partial w_1} \\ \vdots \\ \frac{\partial l}{\partial w_d} \end{bmatrix} \quad w \in \mathbb{R}^d$$

- Gradient Descent's update $s = -\alpha \nabla l(w)$, where $\alpha > 0$

- Therefore, we can derive

$$\begin{aligned} l(w+s) &\approx l(w) + \nabla l(w)^T (-\alpha \nabla l(w)) \\ &\approx l(w) + (-\alpha \nabla l(w)^T \nabla l(w)) \\ &\approx l(w) - \underbrace{\alpha}_{>0} \underbrace{\|\nabla l(w)\|_2}_{>0} < l(w) \end{aligned}$$

\therefore เมื่อ $l(w) - s$ ยิ่งใกล้กับค่า $l(w)$ เนื่องจากค่า s มากกว่า 0 ดังนั้นเมื่อทำการ update โดยใช้ s จะเป็นการลดค่า $l(w)$ นั่นเอง

What's to pick up?

- logistic loss function $l(w) = \sum_{i=1}^n \log(1 + e^{-y_i w^T x_i})$ also measures how well the Classifier performs on the data (ideally, we wish for $l(w)$ to be smallest possible)
- The gradient Descent is an algorithm that finds an optimal point of several loss functions.

$P(y|x)$ $P(x|y)$
- logistic regression is a discriminative counterpart is Naive Bayes

Naive Bayes with Logistic Regression:

- with little data and if the modeling distribution is right, Naive Bayes tends to beat Logistic Regression.
- As data becomes larger, logistic Regression will outperform Naive Bayes which suffers from the fact that the modeling assumption may not be the right one
- เมื่อได้เปรียบ Naive Bayes คือการจำลอง distribution ขึ้นด้วย data ที่มี generate ได้ ถ้า data มากแล้ว อันนั้น distribution ในตัว ดังนั้น modeling distribution มันไม่จำเป็น