

Do not distribute course material

You may not and may not allow others to reproduce or distribute lecture notes and course materials publicly whether or not a fee is charged.

Topic 3

Linear Classification & Logistic Regression

PROF. LINDA SELLIE

- <http://cs229.stanford.edu/notes2020fall/notes2020fall/cs229-notes1.pdf>
- <https://eight2late.wordpress.com/2017/07/11/a-gentle-introduction-to-logistic-regression-and-lasso-regularisation-using-r/>

Outline

- ❑ Motivating example: How can we classify? ↗ How can we use a hyperplane for a classification problem?
- ❑ Estimating probabilities ↗ Can we predict not only which class an example belongs to - but a confidence score of that classification
- ↗ Maximum likelihood ↗ How can we find the most likely hyperplane? Could we write a function to describe how likely a hyperplane was to have generated the dataset?
 - Iterative approach - gradient ascent ↗ Maximizing the function
- ❑ Thinking about different types of error ↗ Some errors are more costly than other errors. Can we modify our predictions to decrease one type of error (and perhaps increase another type of error?)
- ❑ Transformation of the features ↗ Extending our algorithm to nonlinear decision boundaries
- ❑ Multiple classes ↗ What if we have more than two classes?

Outline

- ❑ Motivating example: How can we classify? ↗ How can we use a hyperplane for a classification problem?
- Which model
- Finding an objective function
 - Iterative Optimizer
 - gradient ascent
 - ↗ Maximizing the function
- How can we find the most likely hyperplane? Could we write a function to describe how likely a hyperplane was to have generated the dataset?
- ❑ Thinking about different types of error
 - ↗ Some errors are more costly than other errors. Can we modify our predictions to decrease one type of error (and perhaps increase another type of error?)
- ❑ Transformation of the features ↗ Extending our algorithm to nonlinear decision boundaries
- ❑ Multiple classes ↗ What if we have more than two classes?



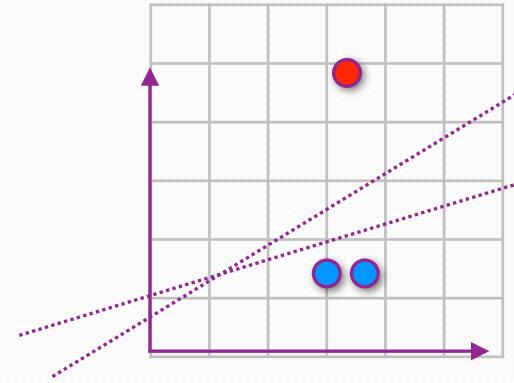
Given $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
how can we find the “best”
hyperplane, \mathbf{w} ?

Optimize \mathbf{w}

We first need to decide what makes
one hyperplane better than another
(i.e. an objective function) **Pair share**

MLE!

Likelihood of seeing the data



- Our model: $p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) & \text{for } y^{(i)} = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) & \text{for } y^{(i)} = 0 \end{cases}$

- Given the following data: $\mathbf{x}^{(1)\top} = [1, 3.2, 4.7] \quad y^{(1)} = 0$
 $\mathbf{x}^{(2)\top} = [1, 3.5, 1.4] \quad y^{(2)} = 1$
 $\mathbf{x}^{(3)\top} = [1, 3., 1.4] \quad y^{(3)} = 1$

- How likely were we to see the data if the line was:

$$\mathbf{w} = \begin{bmatrix} 1/2 \\ 2/3 \\ -1 \end{bmatrix} \quad \left(1 - \frac{1}{1 + e^{-(1/2+(2/3)3.2-4.7)}}\right)^{1-0.11} \left(\frac{1}{1 + e^{-(1/2+(2/3)3.5-1.4)}}\right)^{0.81} \left(\frac{1}{1 + e^{-(1/2+(2/3)3-1.4)}}\right)^{0.75} = 0.54$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ 1/3 \\ -1 \end{bmatrix} \quad \left(1 - \frac{1}{1 + e^{-(1+(1/3)3.2-4.7)}}\right)^{1-0.11} \left(\frac{1}{1 + e^{-(1+(1/3)3.5-1.4)}}\right)^{0.81} \left(\frac{1}{1 + e^{-(1+(1/3)3-1.4)}}\right)^{0.75} = 0.41$$

Pair share: Write the conditional likelihood function for these three examples

Our model:

$$p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) & \text{for } y^{(i)} = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) & \text{for } y^{(i)} = 0 \end{cases}$$

$$p(y = 1 | \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}}$$

$$p(y = 0 | \mathbf{x}; \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = 1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$D = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)})\}$$

setosa

$$\mathbf{x}^{(2)T} = [1, 3.5 \quad 1.4] \quad y^{(2)} = 1$$
$$\mathbf{x}^{(3)T} = [1, 3. \quad 1.4] \quad y^{(3)} = 1$$



https://en.wikipedia.org/wiki/Iris_flower_data_set#/media/

versicolor

$$\mathbf{x}^{(1)T} = [1, 3.2 \quad 4.7] \quad y^{(1)} = 0$$



<https://commons.wikimedia.org/>

The conditional likelihood function

versicolor

$$\mathbf{x}^{(1)\top} = [1 \ 3.2 \ 4.7] \quad y^{(1)} = 0$$

setosa

$$\begin{aligned}\mathbf{x}^{(2)\top} &= [1 \ 3.5 \ 1.4] & y^{(2)} &= 1 \\ \mathbf{x}^{(3)\top} &= [1 \ 3.0 \ 1.4] & y^{(3)} &= 1\end{aligned}$$

$$p(y = 1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$p(y = 0 \mid \mathbf{x}; \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = 1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$L(\mathbf{w}) = \left(1 - \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x}^{(1)})}}\right) \left(\frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x}^{(2)})}}\right) \left(\frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x}^{(3)})}}\right)$$

$$L(\mathbf{w}) = \left(1 - \frac{1}{1 + e^{-(w_0 + w_1 3.2 + w_2 4.7)}}\right) \left(\frac{1}{1 + e^{-(w_0 + w_1 3.5 + w_2 1.4)}}\right) \left(\frac{1}{1 + e^{-(w_0 + w_1 3 + w_2 1.4)}}\right)$$

$$L(\mathbf{w}) = (1 - p(y = 1 \mid \mathbf{x}^{(1)}; \mathbf{w})) p(y = 1 \mid \mathbf{x}^{(2)}; \mathbf{w}) p(y = 1 \mid \mathbf{x}^{(3)}; \mathbf{w}) = \prod_{i=1}^N p(y^{(i)} \text{ correctly predicted} \mid \mathbf{x}^{(i)}; \mathbf{w})$$

$$L(\mathbf{w}) = \prod_{i:y^{(i)}=1} p(y^{(i)} = 1 \mid \mathbf{x}^{(i)}; \mathbf{w}) \prod_{i:y^{(i)}=0} (1 - p(y^{(i)} = 1 \mid \mathbf{x}^{(i)}; \mathbf{w}))$$

The conditional likelihood function

Conditional likelihood function (conditioned on \mathbf{x}). Larger value means more likely

$$L(\mathbf{w}) = \prod_{i:y^{(i)}=1} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}) \prod_{i:y^{(i)}=0} (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}))$$

Here we assume all the examples are independent

$$\prod_{i:y^{(i)}=1} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}} (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}))^{1-y^{(i)}} \prod_{i:y^{(i)}=0} (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}))^{1-y^{(i)}} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}}$$

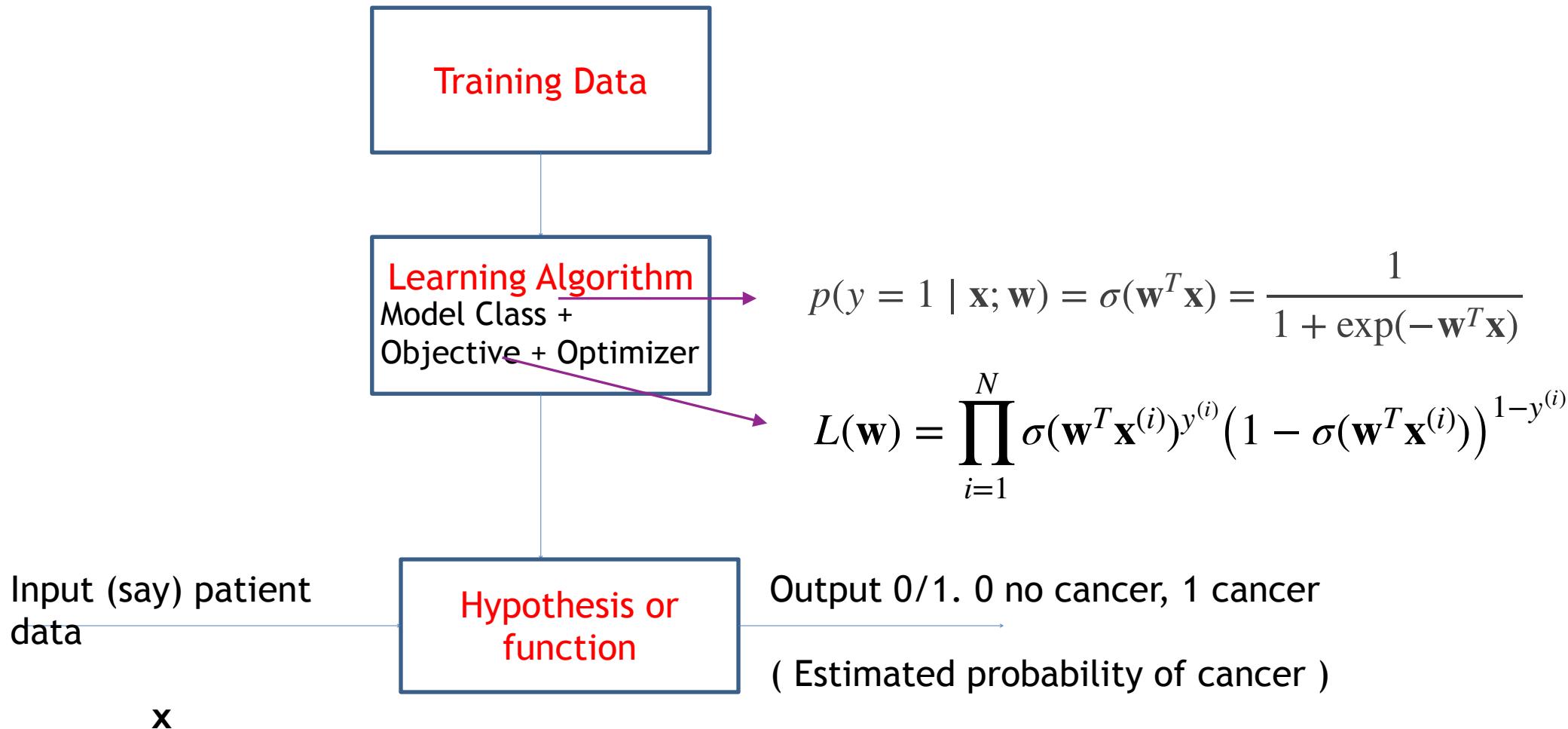
$$L(\mathbf{w}) = \prod_{i=1}^N p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}} (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}))^{1-y^{(i)}}$$

Define: $p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}^{(i)})$

$$= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))^{1-y^{(i)}}$$

$$= \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}^{(i)}}}$$

Supervised Learning



Pair share: how do we find the \mathbf{w} that maximizes this function

$$\text{Maximize } L(\mathbf{w}) = \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))^{1-y^{(i)}}$$

The log-likelihood function

Define: $p(y^{(i)} = 1 \mid \mathbf{x}^{(i)}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$
 $= \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$

□ We wanted to maximize $L(\mathbf{w}) = \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))^{1-y^{(i)}}$

□ This is the same as maximizing $\ell(\mathbf{w}) = \ln(L(\mathbf{w}))$

$$= \ln \left[\prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))^{1-y^{(i)}} \right]$$

$$= \sum_{i=1}^N \ln \left[\underbrace{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})^{y^{(i)}}}_{a^c} \underbrace{(1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))^{1-y^{(i)}}}_{b^d} \right]$$

$$\log a^c b^d = c \log a + d \log b$$

$$= \sum_{i=1}^N \left[\underbrace{y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}_{c \log a} + \underbrace{(1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))}_{d \log b} \right]$$

□ How do we maximize the conditional likelihood?

Equivalent objective function choices

$$\ell(\mathbf{w}) = \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

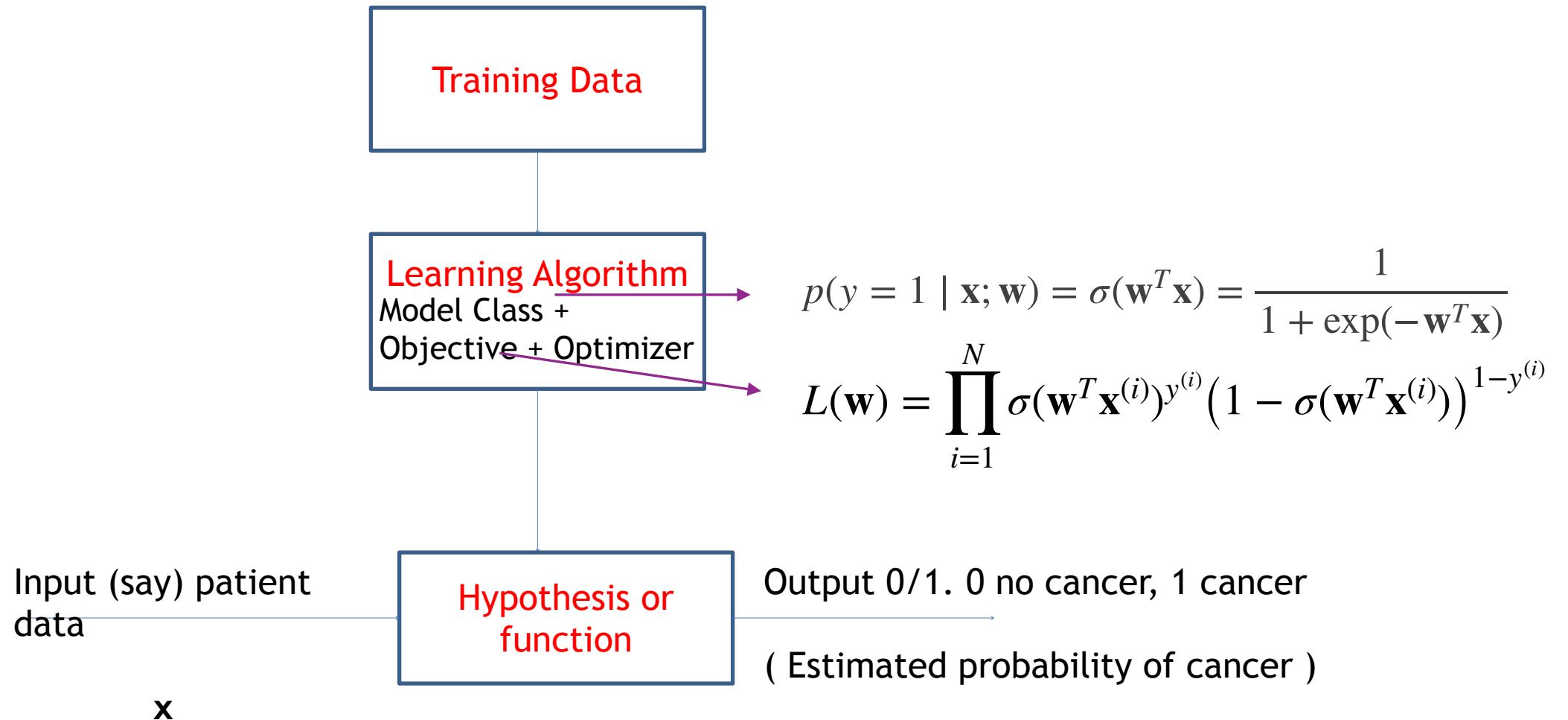
Finds same \mathbf{w}

$$\frac{1}{N} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

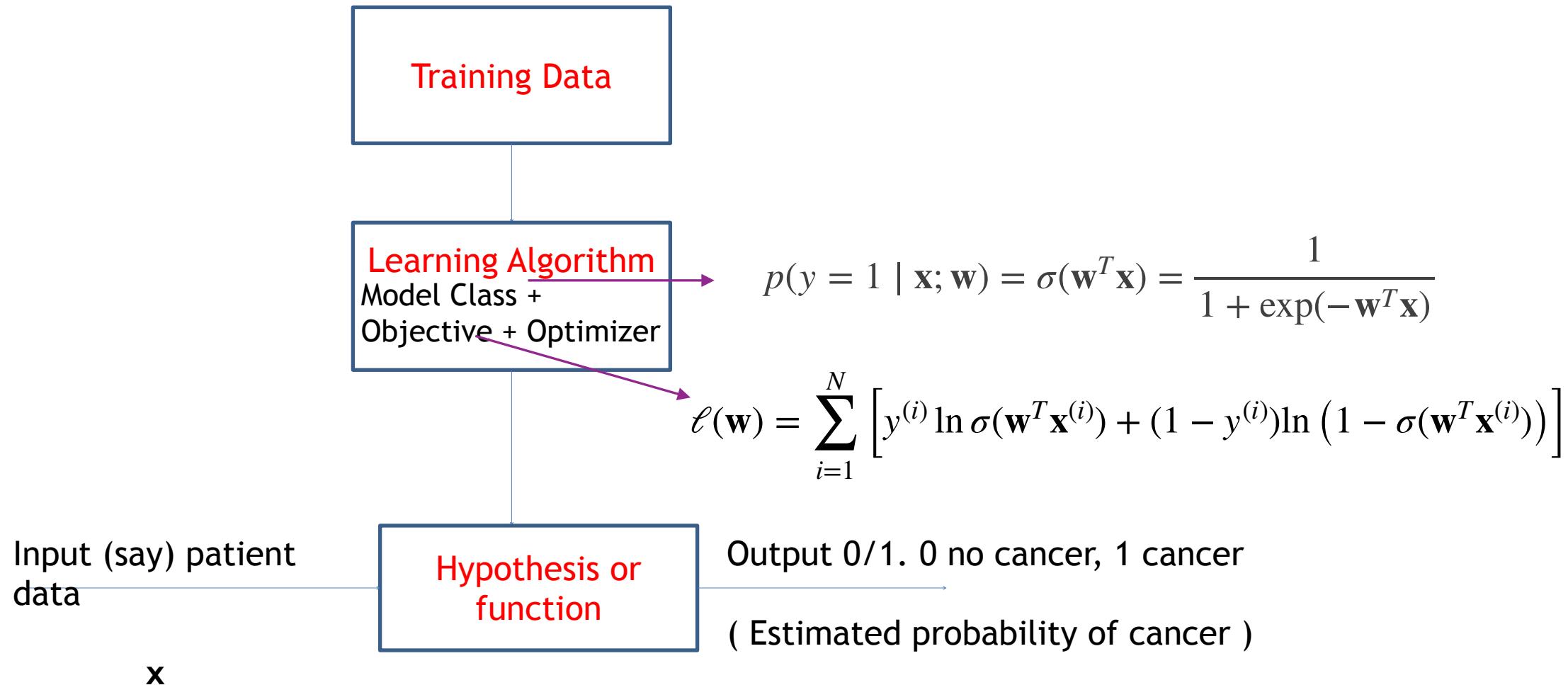
$$-\ell(\mathbf{w}) = \sum_{i=1}^N - \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

Negative log likelihood is the
Cross entropy error
(Find which parameters minimize
the function)

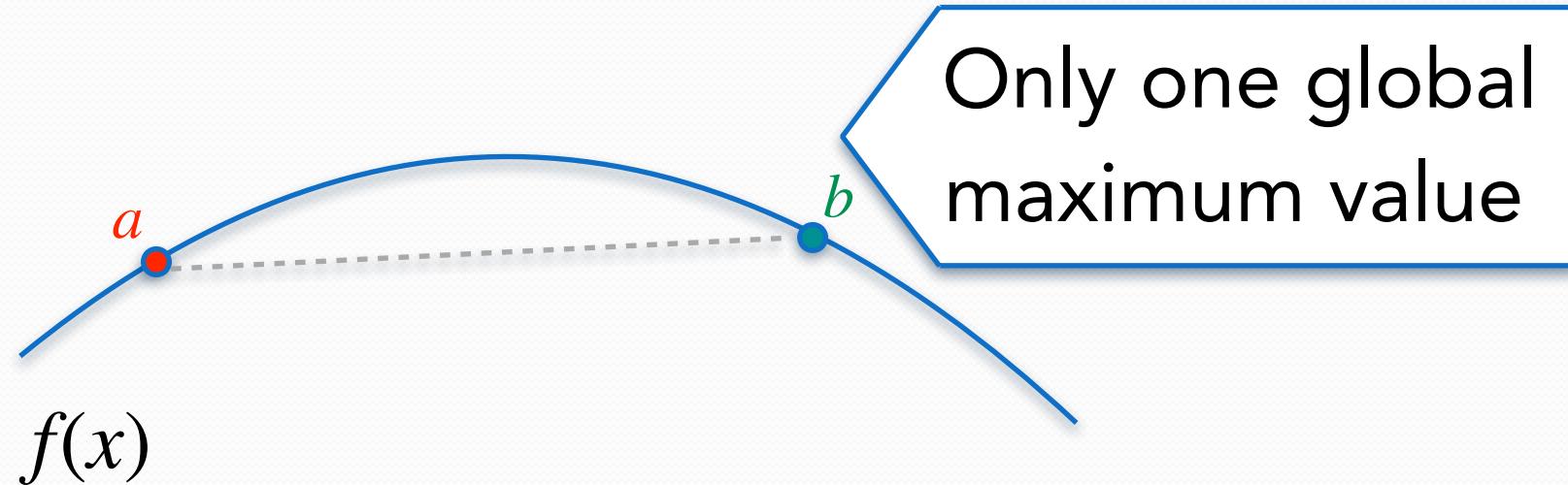
Supervised Learning



Supervised Learning

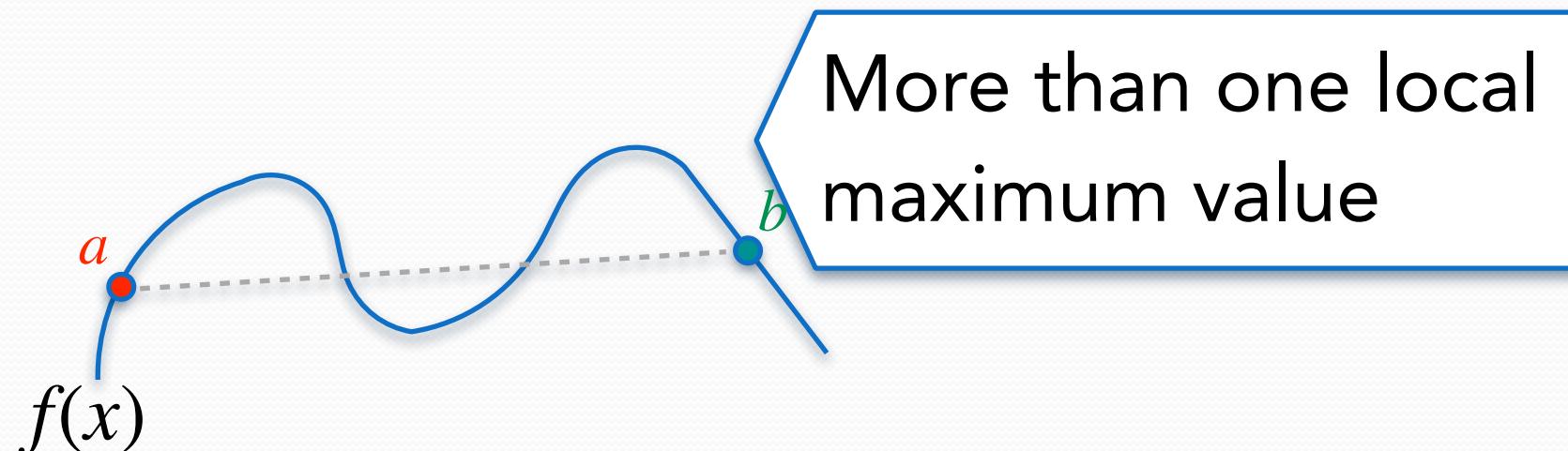


Concave function



Only one global maximum value

Not a concave function



More than one local maximum value

Finding \mathbf{w}

$$\ell(\mathbf{w}) = \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

- There is no closed form solution (i.e. we don't have a global optimization technique)
- We can move toward the optimal value using gradient *ascent*
- To find which way to move, we take the gradient of $\ell(\mathbf{w})$.

Outline

- ❑ Motivating example: How can we classify? ↗ How can we use a hyperplane for a classification problem?
- ❑ Estimating probabilities ↗ Can we predict not only which class an example belongs to - but a confidence score of that classification
- ❑ Maximum likelihood ↗ How can we find the most likely hyperplane? Could we write a function to describe how likely a hyperplane was to have generated the dataset?

- ❑ Iterative approach - gradient ascent ↗ Maximizing the function
- ❑ Thinking about different types of error ↗ Some errors are more costly than other errors. Can we modify our predictions to decrease one type of error (and perhaps increase another type of error?)
- ❑ Transformation of the features ↗ Extending our algorithm to nonlinear decision boundaries
- ❑ Multiple classes ↗ What if we have more than two classes?



Maximize a function by repeatedly moving toward the maximum

1. For $i = 1$ to num_iters :
if $f'(w) > 0$ then f is increasing,
move w a little to the __

- if $f'(w) < 0$ then f is decreasing,
move w a little to the __

right, right

left, left

right, left

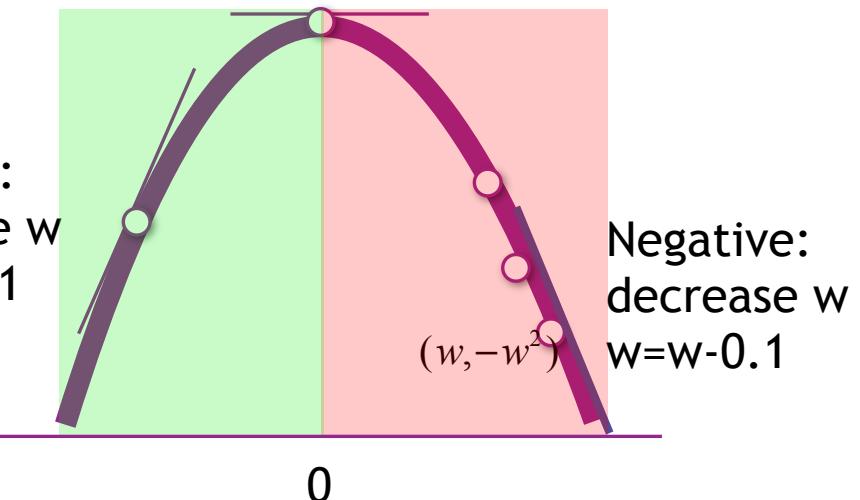
left, right

none of these options

$$f(w) = -w^2$$

$$\frac{df}{dw} = -2w$$

Zero at
maximum value



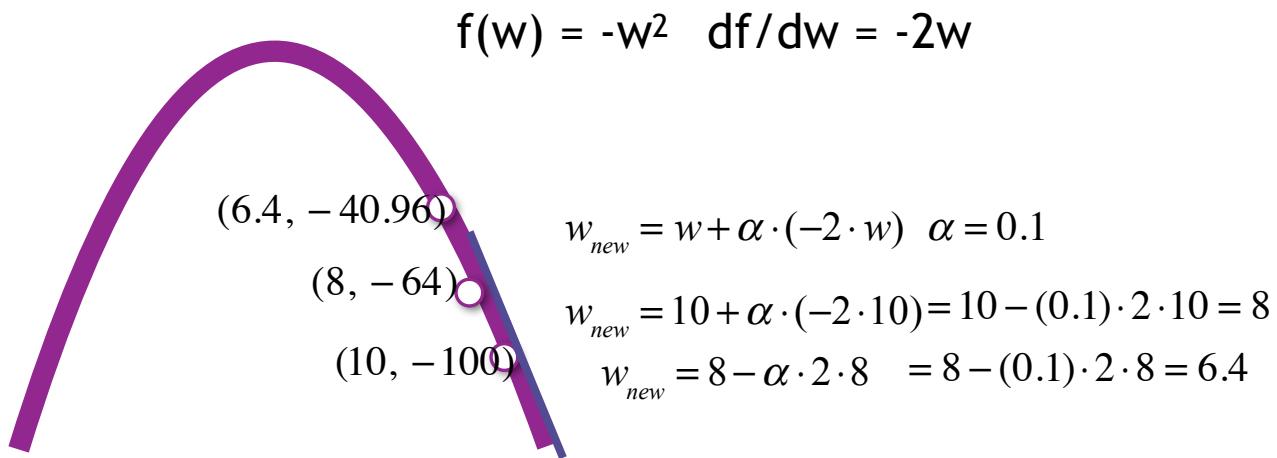
Maximize a function by repeatedly moving toward the maximum

Algorithm:

For i = 1 to num_iters:

if $f'(w) > 0$ then f is increasing,
move w a little to the right

if $f'(w) < 0$ then f is decreasing,
move w a little to the left



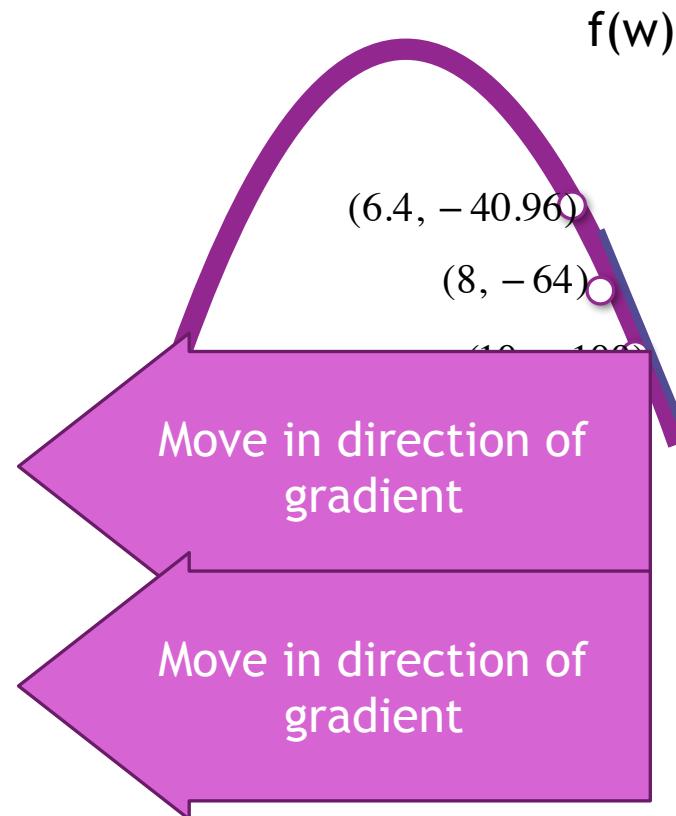
Maximize a function by repeatedly moving toward the maximum

Algorithm:

For i = 1 to num_iters:

 if $f'(w) > 0$ then f is increasing,
 move w a little to the right

 if $f'(w) < 0$ then f is decreasing,
 move w a little to the left



$$f(w) = -w^2 \quad df/dw = -2w$$

$$w_{new} = w + \alpha \cdot (-2 \cdot w) \quad \alpha = 0.1$$

$$w_{new} = 10 + \alpha \cdot (-2 \cdot 10) = 10 - (0.1) \cdot 2 \cdot 10 = 8$$

$$w_{new} = 8 - \alpha \cdot 2 \cdot 8 = 8 - (0.1) \cdot 2 \cdot 8 = 6.4$$

The Gradient

□ We wanted to maximize $\frac{1}{N} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N [y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))]$

□ We would like to find where the derivative is zero (ie the most likely values for \mathbf{w} .)
We do this by taking steps towards the maximum value

$$w_j = w_j + \frac{\alpha}{N} \frac{\partial}{\partial w_j} \ell(\mathbf{w})$$

for $k = 1$ to `num_iter`

$$\text{temp0} = w_0 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_0}$$

$$\text{temp1} = w_1 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_1}$$

$$\text{temp2} = w_2 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_2}$$

$$w_0^{new} = \text{temp0}$$

$$w_1^{new} = \text{temp1}$$

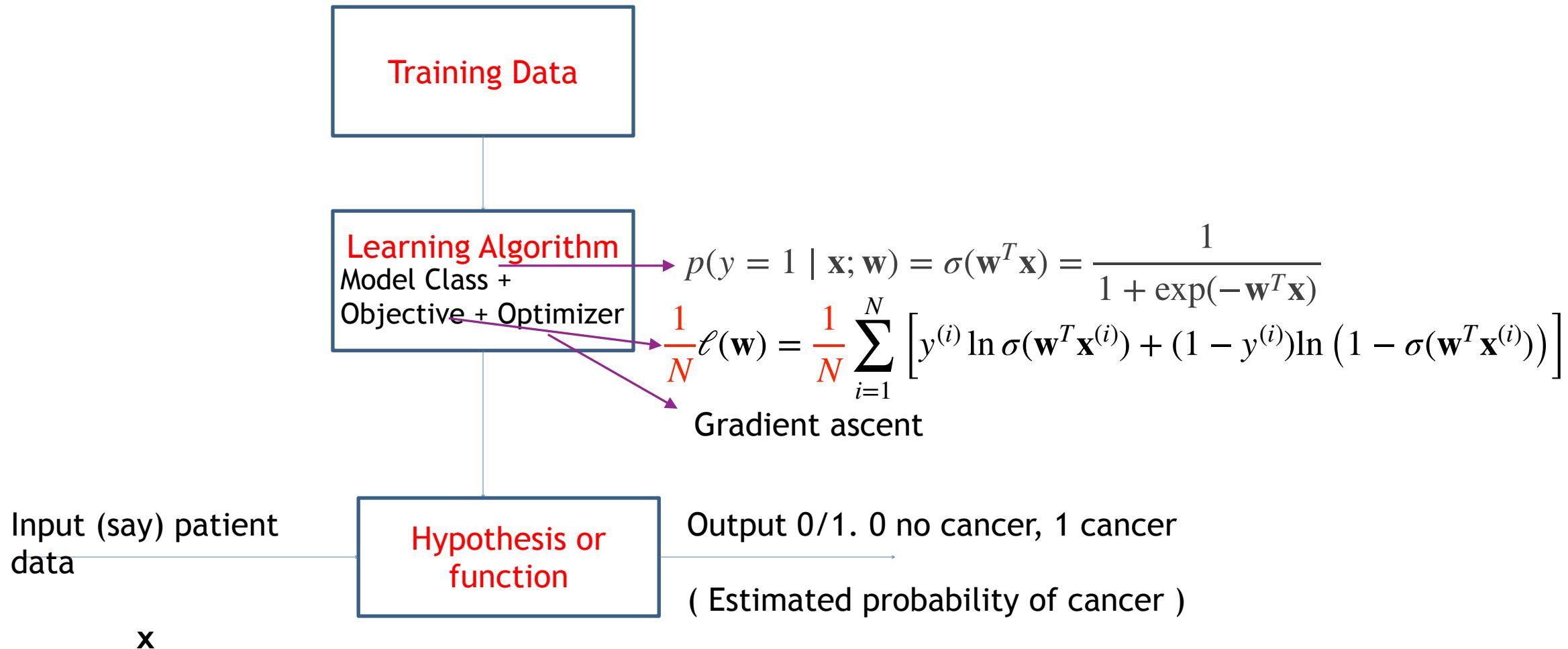
$$w_2^{new} = \text{temp2}$$

[1 3.2 4.7]	[1 3.5 1.4]
[1 3.2 4.5]	[1 3. 1.4]
[1 3.1 4.9]	[1 3.2 1.3]
[1 2.3 4.]	[1 3.1 1.5]
[1 2.8 4.6]	[1 3.6 1.4]
[1 2.8 4.5]	[1 3.9 1.7]
[1 3.3 4.7]	[1 3.4 1.4]
[1 2.4 3.3]	[1 3.4 1.5]
[1 2.9 4.6]	[1 2.9 1.4]
[1 2.7 3.9]	[1 3.1 1.5]

Label 0 examples

23
Label 1 examples

Supervised Learning



Its just math...

z z z

$$h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}}$$

Our update rule:

- We wanted to maximize $\frac{1}{N} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$
- We would like to find where the gradient is zero (ie the most likely values for \mathbf{w} .) We do this by taking steps towards the maximum value $w_j = w_j + \frac{\alpha}{N} \frac{\partial}{\partial w_j} \ell(\mathbf{w})$ for each j

The derivative of the sigmoid function:

$$\begin{aligned}\frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) &= \frac{1}{(1+e^{-z})^2} e^{-z} \\&= \frac{1}{(1+e^{-z})^2} (e^{-z} + 1 - 1) = \frac{1+e^{-z}}{(1+e^{-z})^2} - \frac{1}{(1+e^{-z})^2} \\&= \frac{1}{(1+e^{-z})} - \frac{1}{(1+e^{-z})^2} = \frac{1}{(1+e^{-z})} \left(1 - \frac{1}{(1+e^{-z})} \right)\end{aligned}$$

You will not be tested on this material with the pink background

Thus:

$$\sigma(z) = \left(\frac{1}{1+e^{-z}} \right)$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1-\sigma(z))$$

Useful to know before we derive $\frac{\partial}{\partial w_j} \ell(\mathbf{w})$

$$\frac{\partial \mathbf{w}^T \mathbf{x}}{\partial w_j} = \frac{\partial (w_0x_0 + w_1x_1 + \dots + w_jx_j + \dots + w_dx_d)}{\partial w_j} = x_j$$

$$\frac{\partial \sigma(\mathbf{w}^T \mathbf{x})}{\partial w_j} = \sigma(\mathbf{w}^T \mathbf{x})(1 - \sigma(\mathbf{w}^T \mathbf{x})) \frac{\partial \mathbf{w}^T \mathbf{x}}{\partial w_j} = \sigma(\mathbf{w}^T \mathbf{x})(1 - \sigma(\mathbf{w}^T \mathbf{x})) x_j$$

$$\boxed{\frac{\partial \sigma(\mathbf{w}^T \mathbf{x})}{\partial w_j} = \sigma(\mathbf{w}^T \mathbf{x})(1 - \sigma(\mathbf{w}^T \mathbf{x})) x_j}$$

The Gradient

$$\sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}}$$

$$\frac{\partial \sigma(\mathbf{w}^T \mathbf{x})}{\partial w_j} = \sigma(\mathbf{w}^T \mathbf{x})(1 - \sigma(\mathbf{w}^T \mathbf{x})) \frac{\partial \mathbf{w}^T \mathbf{x}}{w_j}$$

$$\frac{\partial \mathbf{w}^T \mathbf{x}}{\partial w_j} = x_j$$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} \ell(\mathbf{w}) &= \frac{\partial}{\partial w_j} \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right] = \sum_{i=1}^N \left[y^{(i)} \frac{\partial}{\partial w_j} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \frac{\partial}{\partial w_j} \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right] \\
 &= \sum_{i=1}^N \left[y^{(i)} \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} \frac{\partial \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}{w_j} - (1 - y^{(i)}) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})} \frac{\partial \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}{w_j} \right] = \sum_{i=1}^N \left[y^{(i)} \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})} \right] \frac{\partial \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}{w_j} \\
 &= \sum_{i=1}^N \left[y^{(i)} \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})} \right] \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \frac{\partial \mathbf{w}^T \mathbf{x}^{(i)}}{w_j} = \sum_{i=1}^N \left(y^{(i)} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) - (1 - y^{(i)}) \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) \right) x_j^{(i)} \\
 &= \sum_{i=1}^N \left(y^{(i)} \cdot 1 - y^{(i)} \cdot \cancel{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} - 1 \cdot \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + y^{(i)} \cancel{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} \right) x_j^{(i)} \\
 &= \sum_{i=1}^N \left(y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) \right) x_j^{(i)}
 \end{aligned}$$

Error of i^{th} training example

Interpretation:

- ❑ We wanted to maximize $\frac{1}{N} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$
- ❑ We would like to find where the partial derivative is zero (ie the most likely values for \mathbf{w}). We do this by taking steps towards the maximum value
- ❑ How does an example affect the update?
 - ★ if $y=1$ and $\sigma(\mathbf{w}^T \mathbf{x}) \approx 1$, then $(y - \sigma(\mathbf{w}^T \mathbf{x})) \approx 0$, almost no change!
 - ★ if $y=1$ and $\sigma(\mathbf{w}^T \mathbf{x}) \approx 0$, then $(y - \sigma(\mathbf{w}^T \mathbf{x})) \approx 1$, approx α/N times the j^{th} feature
 - ★ if $y=0$ and $\sigma(\mathbf{w}^T \mathbf{x}) \approx 0$, then $(y - \sigma(\mathbf{w}^T \mathbf{x})) \approx 0$, almost no change!
 - ★ if $y=0$ and $\sigma(\mathbf{w}^T \mathbf{x}) \approx 1$, then $(y - \sigma(\mathbf{w}^T \mathbf{x})) \approx -1$, approx $-\alpha/N$ times the j^{th} feature



$$\frac{\partial}{\partial w_j} \ell(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_j^{(i)}$$

Gradient Ascent Algorithm

$$\frac{1}{N} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left[y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

for k = 1 to num_iter

$$temp0 = w_0 + \alpha \frac{\partial \ell(\mathbf{w}) / N}{\partial w_0} = w_0 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_0^{(i)}$$

$$temp1 = w_1 + \alpha \frac{\partial \ell(\mathbf{w}) / N}{\partial w_1} = w_1 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_1^{(i)}$$

:

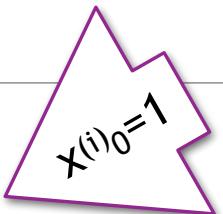
$$tempd = w_d + \alpha \frac{\partial \ell(\mathbf{w}) / N}{\partial w_d} = w_d + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_d^{(i)}$$

$$w_0 = temp0$$

$$w_1 = temp1$$

:

$$w_d = tempd$$



Simultaneous update

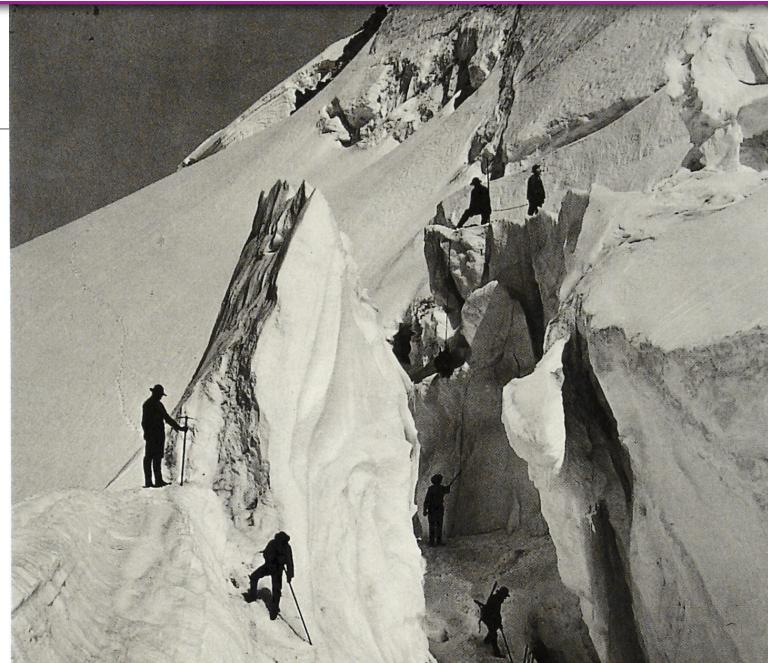
Running time:
 $O(Nd \cdot \# \text{ iter})$

1. Does this algorithm work for any choice of initial values for \mathbf{w} ?

Yes

No

It depends on the dataset



Gradient Ascent Step

$$X = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 3 & 5 \\ 1 & 4 & 1 \\ 1 & 3 & 1.5 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \rightarrow \quad X\mathbf{w} = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 3 & 5 \\ 1 & 4 & 1 \\ 1 & 3 & 1.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \rightarrow \quad \hat{\mathbf{y}} = \begin{bmatrix} 0.5 \\ 0.5 \\ 0.5 \\ 0.5 \end{bmatrix}$$

for k = 1 to num_iter

$$temp0 = w_0 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_0} = w_0 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_0^{(i)} = w_0 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) + (-0.5) + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) + (-0.5) + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) + (0.5) + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) + (0.5) \right]$$

$$temp1 = w_1 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_1} = w_1 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_1^{(i)} = w_1 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) x_1^{(1)} + (-0.5)3 + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) x_1^{(2)} + (-0.5)3 + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) x_1^{(3)} + (0.5)4 + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) x_1^{(4)} + (0.5)3 \right]$$

$$temp2 = w_2 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_2} = w_2 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_2^{(i)} = w_2 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) x_2^{(1)} + (-0.5)4 + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) x_2^{(2)} + (-0.5)5 + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) x_2^{(3)} + (0.5)1 + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) x_2^{(4)} + (0.5)(1.5) \right]$$

$$w_0^{new} = temp0 \quad 0$$

$$w_1^{new} = temp1 \quad 0.025$$

$$w_2^{new} = temp2 \quad -0.1625$$

Gradient Ascent Example

$$X = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 3 & 5 \\ 1 & 4 & 1 \\ 1 & 3 & 1.5 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 0 \\ 0 \\ 0.025 \\ -0.1625 \end{bmatrix} \quad X\mathbf{w} = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 3 & 5 \\ 1 & 4 & 1 \\ 1 & 3 & 1.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0.025 \\ -0.1625 \end{bmatrix} = \begin{bmatrix} -0.575 \\ -0.7375 \\ -0.0625 \\ -0.16875 \end{bmatrix} \quad \hat{\mathbf{y}} = \begin{bmatrix} 0.3601 \\ 0.3236 \\ 0.4844 \\ 0.4579 \end{bmatrix}$$

for k = 1 to num_iter

$$temp0 = w_0 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_0} = w_0 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_0^{(i)} = w_0 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) \right]$$

$$(0-0.3601) \quad (0-0.3236) \quad (1-0.4844) \quad (1-0.4579)$$

$$temp1 = w_1 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_1} = w_1 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_1^{(i)} = w_1 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) x_1^{(1)} + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) x_1^{(2)} + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) x_1^{(3)} + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) x_1^{(4)} \right]$$

$$(0-0.3601)*3 \quad (0-0.3236)*3 \quad (1-0.4844)*4 \quad (1-0.4579)*3$$

$$temp2 = w_2 + \frac{\alpha}{N} \frac{\partial \ell(\mathbf{w})}{\partial w_2} = w_2 + \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_2^{(i)} = w_2 + \frac{0.2}{4} \left[(0 - \sigma(\mathbf{w}^T \mathbf{x}^{(1)})) x_2^{(1)} + (0 - \sigma(\mathbf{w}^T \mathbf{x}^{(2)})) x_2^{(2)} + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(3)})) x_2^{(3)} + (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(4)})) x_2^{(4)} \right]$$

$$(0-0.3601)*4 \quad (0-0.3236)*5 \quad (1-0.4844)*1 \quad (1-0.4579)*(1.5)$$

$$w_0^{new} = temp0 \quad 0.019$$

$$w_1^{new} = temp1 \quad 0.107$$

$$w_2^{new} = temp2 \quad -0.249$$