# 👍 Linear classifiers:
## 👎 Parameter learning

# Learn a probabilistic classification model

*"The sushi & everything else were awesome!"*

*"The sushi was good, the service was OK"*

Definite **+1**

Not sure

$P(y=+1|\mathbf{x}=$ *"The sushi & everything else were awesome!"* $)$
= 0.99

$P(y=+1|\mathbf{x}=$ *"The sushi was good, the service was OK"* $)$
= 0.55

Many classifiers provide a degree of certainty:

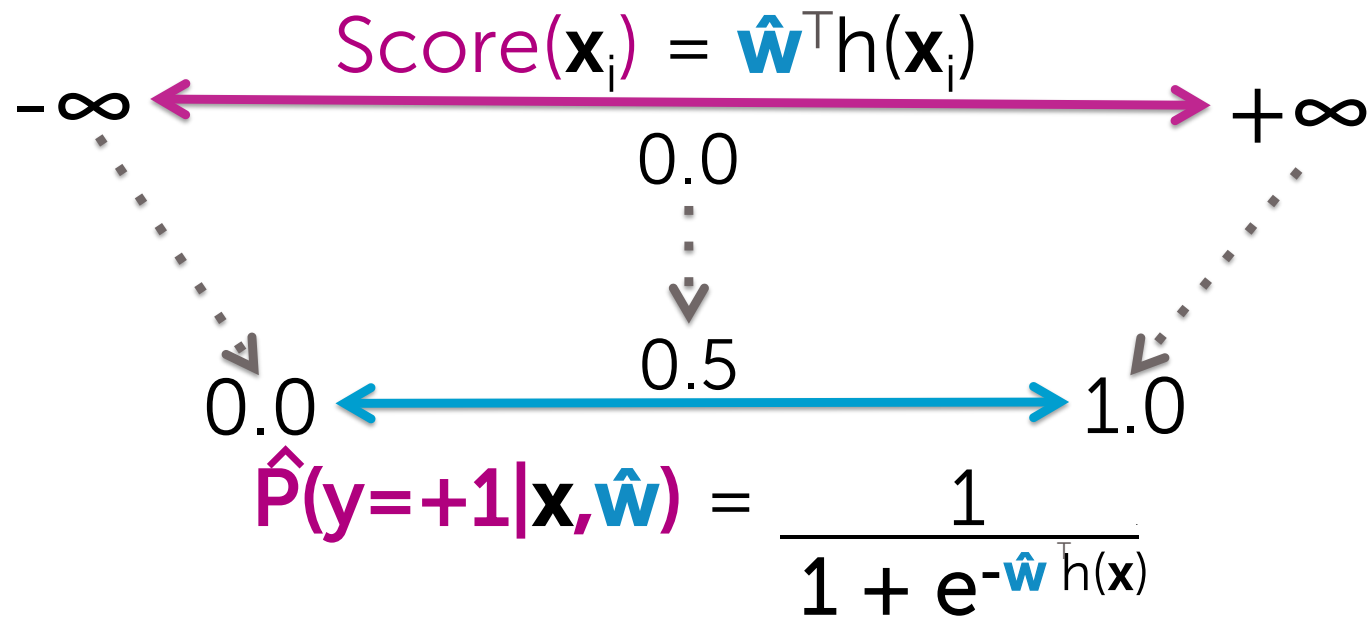Output label

Input sentence

$P(y|\mathbf{x})$

Extremely useful in practice

# A (linear) classifier

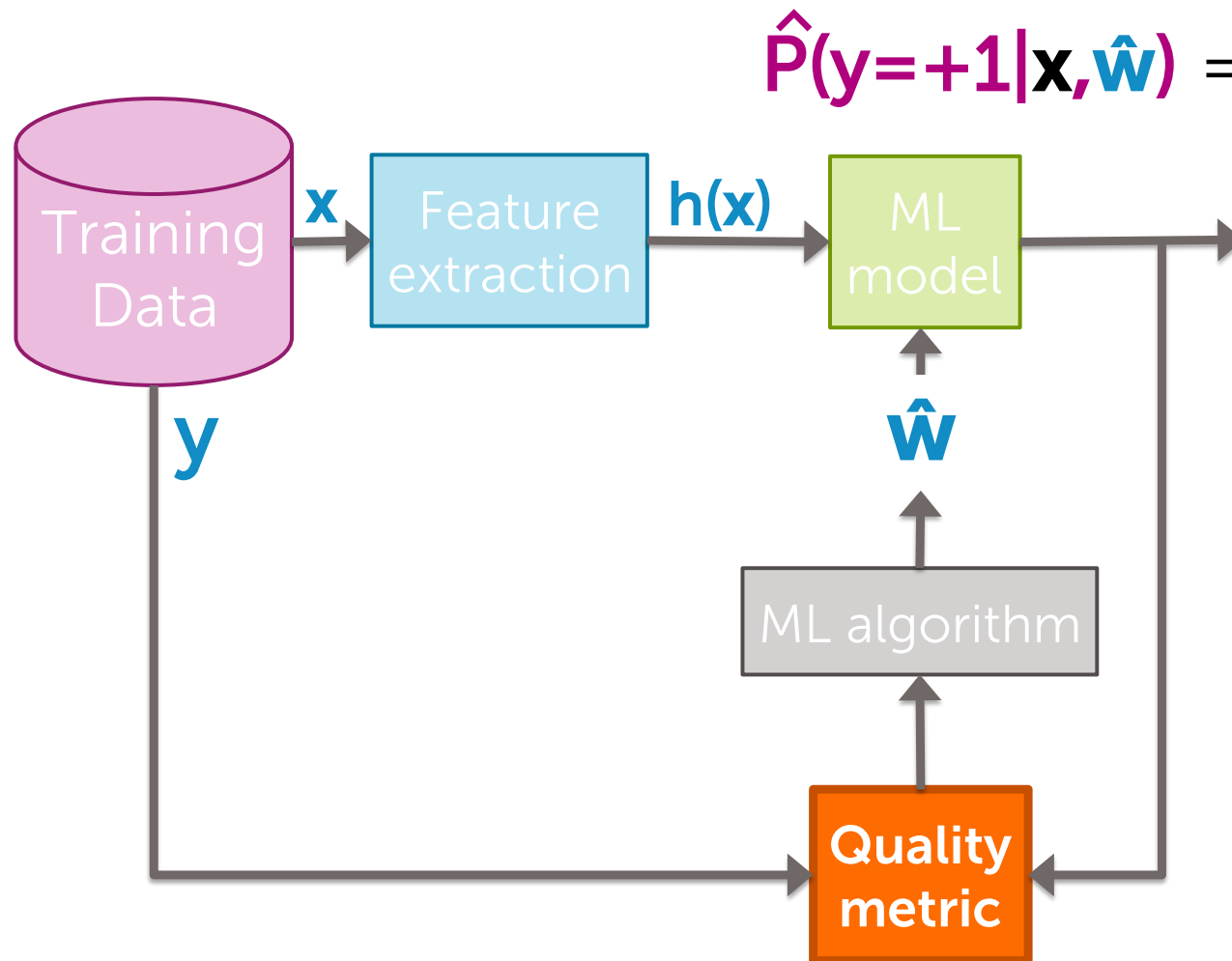- Will use training data to learn a weight or coefficient for each word

| Word | Coefficient | Value |
|---|---|---|
| | $\hat{w}_0$ | -2.0 |
| good | $\hat{w}_1$ | 1.0 |
| great | $\hat{w}_2$ | 1.5 |
| awesome | $\hat{w}_3$ | 2.7 |
| bad | $\hat{w}_4$ | -1.0 |
| terrible | $\hat{w}_5$ | -2.1 |
| awful | $\hat{w}_6$ | -3.3 |
| restaurant, the, we, … | $\hat{w}_7, \hat{w}_8, \hat{w}_{9,…}$ | 0.0 |
| … | | … |

Machine Learning Specialization

# Logistic regression model

$$\text{Score}(\mathbf{x}_i) = \hat{\mathbf{w}}^\top h(\mathbf{x}_i)$$

$-\infty$ ⟵━━━━━━━━━━━━⟶ $+\infty$

0.0

0.5

0.0 ⟵━━━━━━━━⟶ 1.0

$$\hat{P}(y=+1|\mathbf{x},\hat{\mathbf{w}}) = \frac{1}{1 + e^{-\hat{\mathbf{w}}^\top h(\mathbf{x})}}$$

# Quality metric for logistic regression: Maximum likelihood estimation

$$\hat{P}(y=+1|\mathbf{x},\hat{\mathbf{w}}) = \frac{1}{1 + e^{-\hat{\mathbf{w}}\, h(\mathbf{x})}}$$

Machine Learning Specialization

# Learning problem

Training data:
$N$ observations $(\mathbf{x}_i, y_i)$

| **x**[1] = #awesome | **x**[2] = #awful | y = sentiment |
|:---:|:---:|:---:|
| 2 | 1 | +1 |
| 0 | 2 | -1 |
| 3 | 3 | -1 |
| 4 | 1 | +1 |
| 1 | 1 | +1 |
| 2 | 4 | -1 |
| 0 | 3 | -1 |
| 0 | 1 | -1 |
| 2 | 1 | +1 |

Optimize **quality metric** on training data ➡ $\hat{\mathbf{w}}$

# Finding best coefficients

| x[1] = #awesome | x[2] = #awful | y = sentiment |
|:---:|:---:|:---:|
| 0 | 2 | -1 |
| 3 | 3 | -1 |
| 2 | 4 | -1 |
| 0 | 3 | -1 |
| 0 | 1 | -1 |

| x[1] = #awesome | x[2] = #awful | y = sentiment |
|:---:|:---:|:---:|
| 2 | 1 | +1 |
| 4 | 1 | +1 |
| 1 | 1 | +1 |
| 2 | 1 | +1 |

$P(y=+1|x_i, w) = 0.0$

$P(y=+1|x_i, w) = 1.0$

Pick $\hat{w}$ that makes

# Quality metric = Likelihood function

Negative data points

Positive data points

$P(y=+1|x_i, w) = 0.0$               $P(y=+1|x_i, w) = 1.0$

No $\hat{w}$ achieves perfect predictions (usually)

**Likelihood** $\ell(w)$: Measures quality of fit for model with coefficients $w$

# Find "best" classifier

Maximize likelihood over all possible $w_0, w_1, w_2$

$\ell(w_0=0, w_1=1, w_2=-1.5) = 10^{-6}$

$\ell(w_0=1, w_1=1, w_2=-1.5) = 10^{-5}$

#awful

*Best model:*
Highest likelihood $\ell(\mathbf{w})$
$\hat{\mathbf{w}} = (w_0=1, w_1=0.5, w_2=-1.5)$

$\ell(w_0=1, w_1=0.5, w_2=-1.5) = 10^{-4}$

gradient ascent to find $\hat{w}$

#awesome

Machine Learning Specialization

# Data likelihood

# Quality metric: probability of data

| **x**[1] = #awesome | **x**[2] = #awful | y = sentiment |
|:---:|:---:|:---:|
| 2 | 1 | +1 |

$x_1$                                              $y_1 :$

**If model good, should predict:**

$\hat{y}_1 = +1$

**Pick w to maximize:**

$P(y=+1 | x_1, w) = P(y=+1 | x[1]=2, x[2]=1, w)$

| **x**[1] = #awesome | **x**[2] = #awful | y = sentiment |
|:---:|:---:|:---:|
| 0 | 2 | -1 |

$x_2 =$                                              $y_2 =$

**If model good, should predict:**

$\hat{y}_2 = -1$

**Pick w to maximize:**

$P(y=-1 | x_2, w)$

15

# Maximizing likelihood (probability of data)

| Data point | x[1] | x[2] | y | Choose w to maximize |
|---|---|---|---|---|
| $x_1, y_1$ | 2 | 1 | +1 | $P(y=+1 \mid x_1, w) = P(y=+1 \mid x[1]=2, x[2]=1, w)$ |
| $x_2, y_2$ | 0 | 2 | -1 | $P(y=-1 \mid x_2, w)$ |
| $x_3, y_3$ | 3 | 3 | -1 | $P(y=-1 \mid x_3, w)$ |
| $x_4, y_4$ | 4 | 1 | +1 | $P(y=+1 \mid x_4, w)$ |
| $x_5, y_5$ | 1 | 1 | +1 | |
| $x_6, y_6$ | 2 | 4 | -1 | |
| $x_7, y_7$ | 0 | 3 | -1 | |
| $x_8, y_8$ | 0 | 1 | -1 | |
| $x_9, y_9$ | 2 | 1 | +1 | |

Must combine into single measure of quality ?

Multiply Probabilities

$P(y=+1 \mid x_1, w) \, P(y=-1 \mid x_2, w) \, P(y=-1 \mid x_3, w) \ldots$

Machine Learning Specialization

# Learn logistic regression model with maximum likelihood estimation (MLE)

| Data point | x[1] | x[2] | y | Choose w to maximize |
|------------|------|------|---|----------------------|
| $\mathbf{x}_1, y_1$ | 2 | 1 | $y:$ +1 | $P(y=+1 \mid \mathbf{x}[1]=2, \mathbf{x}[2]=1, \mathbf{w})$ |
| $\mathbf{x}_2, y_2$ | 0 | 2 | -1 | $P(y=-1 \mid \mathbf{x}[1]=0, \mathbf{x}[2]=2, \mathbf{w})$ |
| $\mathbf{x}_3, y_3$ | 3 | 3 | -1 | $P(y=-1 \mid \mathbf{x}[1]=3, \mathbf{x}[2]=3, \mathbf{w})$ |
| $\mathbf{x}_4, y_4$ | 4 | 1 | +1 | $P(y=+1 \mid \mathbf{x}[1]=4, \mathbf{x}[2]=1, \mathbf{w})$ |

$\ell(\mathbf{w}) =$

$$P(y_1 \mid \mathbf{x}_1, \mathbf{w}) \qquad P(y_2 \mid \mathbf{x}_2, \mathbf{w}) \qquad P(y_3 \mid \mathbf{x}_3, \mathbf{w}) \qquad P(y_4 \mid \mathbf{x}_4, \mathbf{w})$$

Num. of data points $\longrightarrow$

$$\ell(w) = \prod_{i=1}^{N} P(y_i \mid \mathbf{x}_i, \mathbf{w})$$

← pick w to make this fn. as large as possible

# Finding best linear classifier
# with gradient ascent

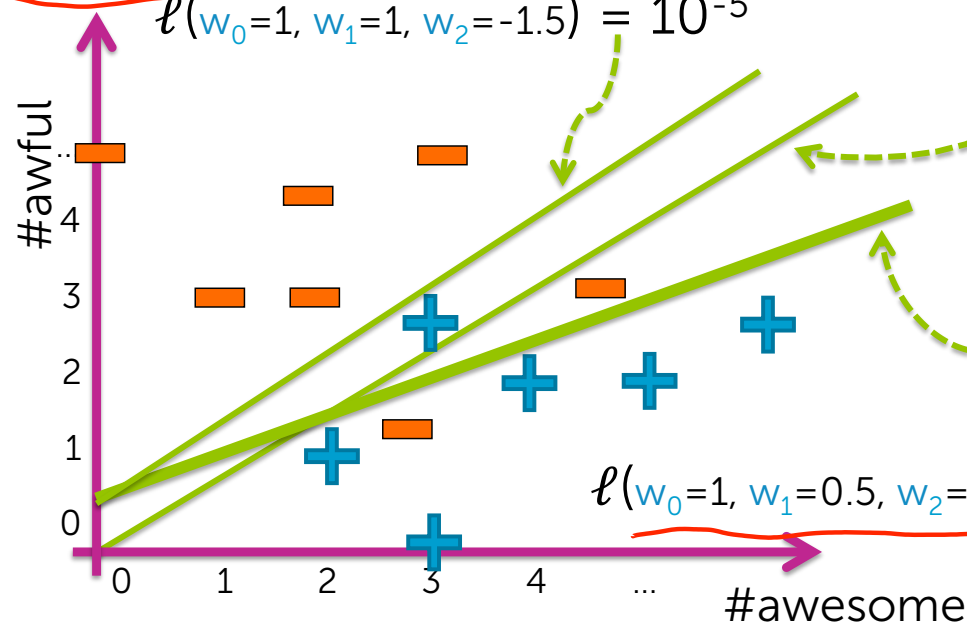$$\hat{P}(y=+1|\mathbf{x},\hat{\mathbf{w}}) = \frac{1}{1 + e^{-\hat{\mathbf{w}}\,h(\mathbf{x})}}$$

Training Data

**x** → Feature extraction → **h(x)** → ML model

**y**

$\hat{\mathbf{w}}$

ML algorithm

Quality metric

Machine Learning Specialization

# Find "best" classifier

Maximize likelihood over all possible $w_0, w_1, w_2$

$$\ell(\mathbf{w}) = \prod_{i=1}^{N} P(y_i \mid \mathbf{x}_i, \mathbf{w})$$

$\ell(w_0=0, w_1=1, w_2=-1.5) = 10^{-6}$

$\ell(w_0=1, w_1=1, w_2=-1.5) = 10^{-5}$

$\ell(w_0=1, w_1=0.5, w_2=-1.5) = 10^{-4}$

*Best model:*
Highest likelihood $\ell(\mathbf{w})$
$\hat{\mathbf{w}} = (w_0=1, w_1=0.5, w_2=-1.5)$

optimize with gradient ascent

#awful

4
3
2
1
0

0   1   2   3   4   ...

#awesome

Machine Learning Specialization

# Maximizing likelihood



Maximize function over all possible $w_0, w_1, w_2$

$$\max_{w_0, w_1, w_2} \prod_{i=1}^{N} P(y_i \mid \mathbf{x}_i, \mathbf{w})$$

$\ell(w_0, w_1, w_2)$ is a function of 3 variables

No closed-form solution ➔ use gradient ascent

# Review of gradient ascent

# Finding the max via hill climbing



Algorithm:

**while** not converged

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \left.\frac{d\ell}{dw}\right|_{w^{(t)}}$$
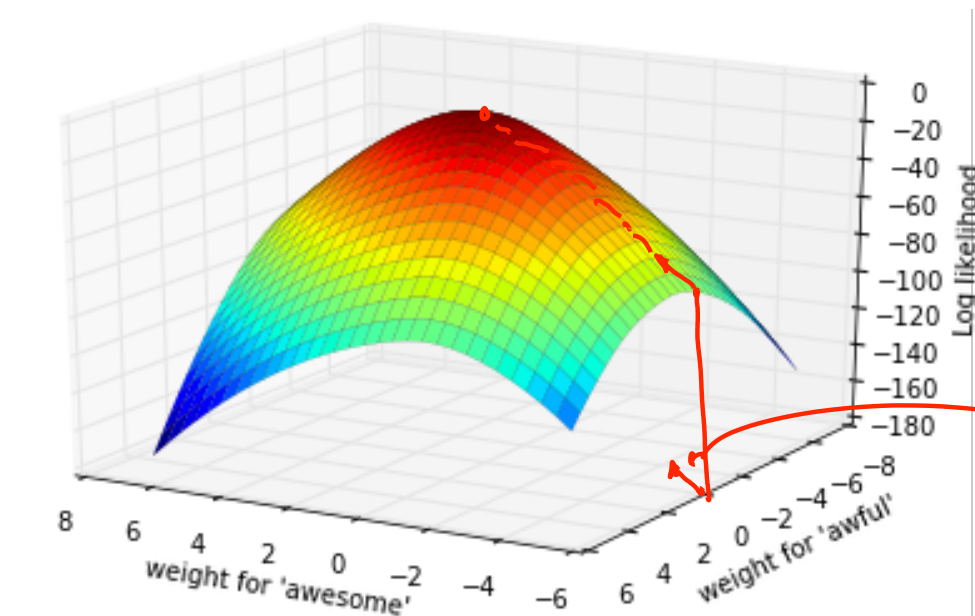
Step size

Machine Learning Specialization

# Convergence criteria

For convex functions, optimum occurs when

$$\frac{d\ell}{dw} = 0$$

$w^*$

In practice, stop when

$$\left.\frac{d\ell}{dw}\right|_{w^{(t)}} < \varepsilon$$

↑ tolerance

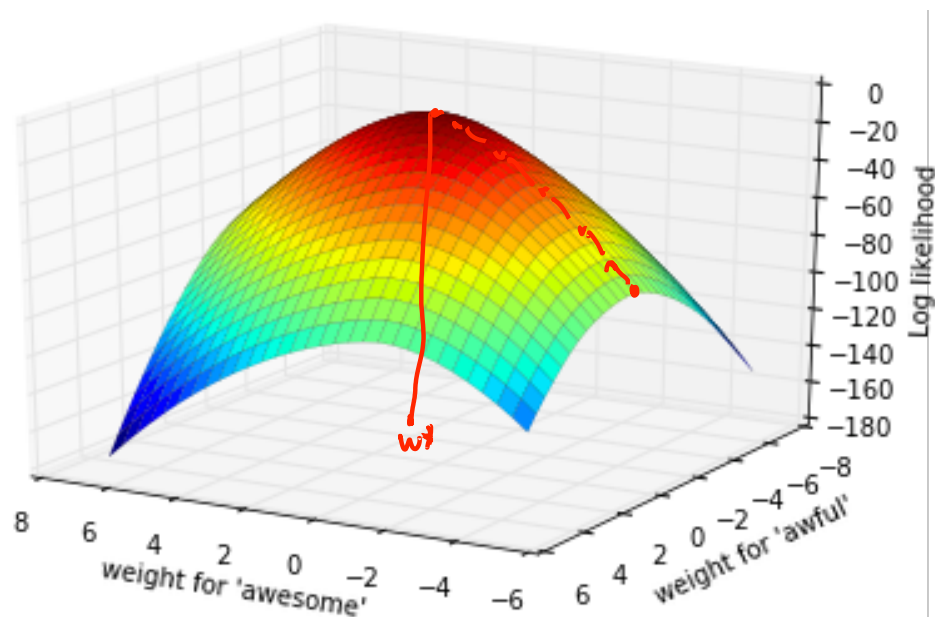<span style="color:magenta">Algorithm:</span>

**while** not converged

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \left.\frac{d\ell}{dw}\right|_{w^{(t)}}$$

Machine Learning Specialization
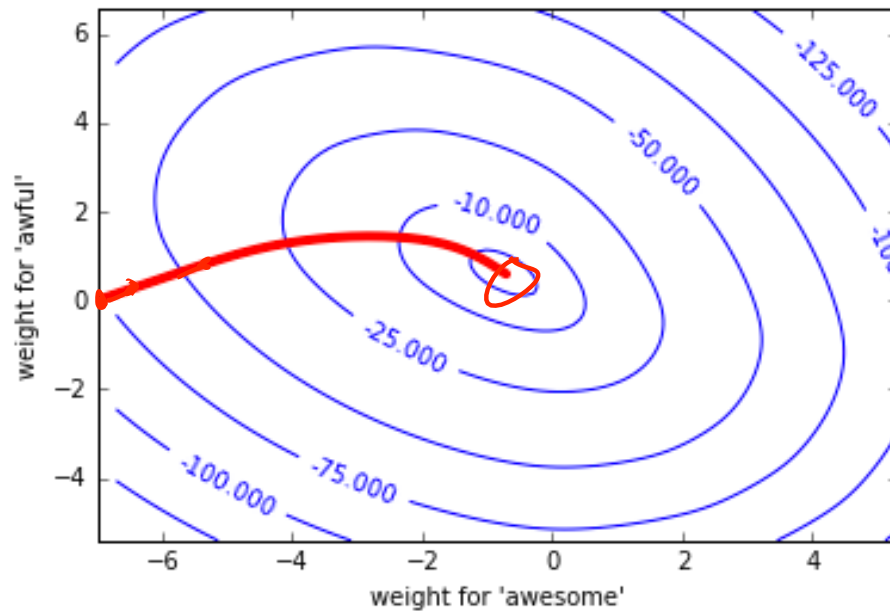
# Moving to multiple dimensions: Gradients



$$\nabla \ell(\mathbf{w}) = \begin{bmatrix} \frac{\partial \ell}{\partial w_0} \\ \frac{\partial \ell}{\partial w_1} \\ \vdots \\ \frac{\partial \ell}{\partial w_D} \end{bmatrix} \leftarrow D+1 \text{ dim vector}$$

# Contour plots

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Gradient ascent



Algorithm:

$w^{(0)} = 0$, random, or something smart.

**while** not converged
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \nabla \ell(\mathbf{w}^{(t)})$$

step size

Machine Learning Specialization

# Learning algorithm for logistic regression

# Derivative of (log-)likelihood

Sum over
data points

Feature
value

Difference between truth and prediction

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^{N} h_j(\mathbf{x}_i) \Big( \mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w}) \Big)$$

predict $x_i$ is positive

Indicator function:

$$\mathbb{1}[y_i = +1] = \begin{cases} 1 & \text{if } y_i = +1 \\ 0 & \text{if } y_i = -1 \end{cases}$$

# Computing derivative

$$\frac{\partial \ell(\mathbf{w}^{(t)})}{\partial \mathbf{w}_j} = \sum_{i=1}^{N} h_j(\mathbf{x}_i)\Big(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w}^{(t)})\Big)$$

$w^{(t)}$:

| $w_0^{(t)}$ | 0 |
|---|---|
| $w_1^{(t)}$ | 1 |
| $w_2^{(t)}$ | -2 |

$\frac{\partial \ell}{\partial w_1}$

$h_1(t) = H$ awesome

| x[1] | x[2] | y | P(y=+1|$\mathbf{x}_i$,w) | Contribution to derivative for $w_1$ |
|---|---|---|---|---|
| 2 | 1 | +1 | 0.5 | $2(1-0.5) = 1$ |
| 0 | 2 | -1 | 0.02 | $0(0-0.02) = 0$ |
| 3 | 3 | -1 | 0.05 | $3(0-0.05) = -0.15$ |
| 4 | 1 | +1 | 0.88 | $4(1-0.88) = 0.48$ |

Total derivative:

$\frac{\partial \ell(w^{(t)})}{\partial w_1} = 1 + 0 - 0.15 + 0.48 = 1.33$

$w_1^{(t+1)} = w_1^{(t)} + \eta \left. \frac{\partial \ell(w^{(t)})}{\partial w_1} \right| \eta = 0.1$

$= 1 + 0.1 \times 1.33 = 1.133$

Machine Learning Specialization

# Derivative of (log-)likelihood: Interpretation

Sum over
data points

Feature
value

Difference between truth and prediction

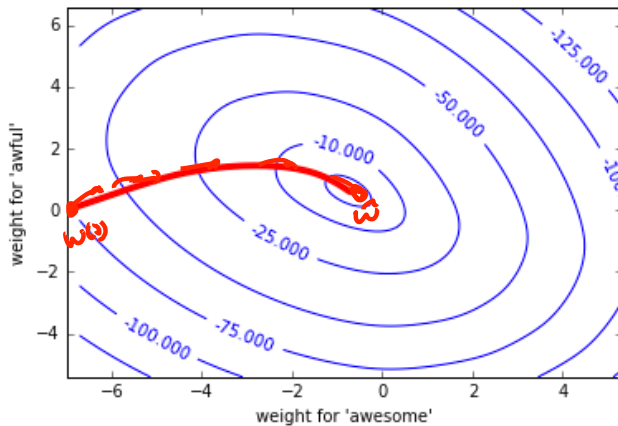$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^{N} h_j(\mathbf{x}_i)\Big(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w})\Big)$$

$\Delta_i$

| If $h_j(\mathbf{x}_i)=1$: | P(y=+1\|x_i,w) ≈ 1 | P(y=+1\|x_i,w) ≈ 0 |
|---|---|---|
| $y_i=+1$ | $\Delta_i = (1 - 1) \approx 0$ <br> ↳ don't change anything! | $\Delta_i \approx 1 \Rightarrow$ increase $w_j$ <br> $\Rightarrow$ Score($x_i$) becomes larger <br> $\Rightarrow P(y=+1\|x_i,w)$ increases |
| $y_i=-1$ | $\Delta_i = -1 \Rightarrow w_j$ to decrease <br> $\Rightarrow$ Score($x_i$) decreases $\Rightarrow P(y=+1\|x_i,w)$ decrease | $\Delta_i \approx 0$ <br> $\Rightarrow$ don't change anything |

# Summary of gradient ascent
# for logistic regression



init $\mathbf{w}^{(1)}=0$ (or randomly, or smartly), $t=1$

while $||\nabla\ell(\mathbf{w}^{(t)})|| > \varepsilon$ ← tolerance

    for $j=0,\ldots,D$ ← coefficient

$$\frac{1}{1+e^{-w^{(t)}\cdot h(x_i)}}$$

      partial[j] $= \displaystyle\sum_{i=1}^{N} h_j(\mathbf{x}_i)\left(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w}^{(t)})\right)$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta\ \text{partial[j]}$$

$t \leftarrow t + 1$

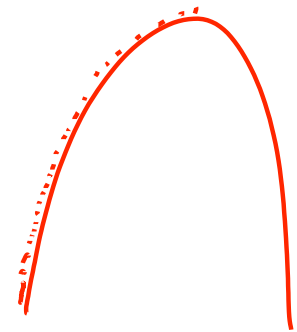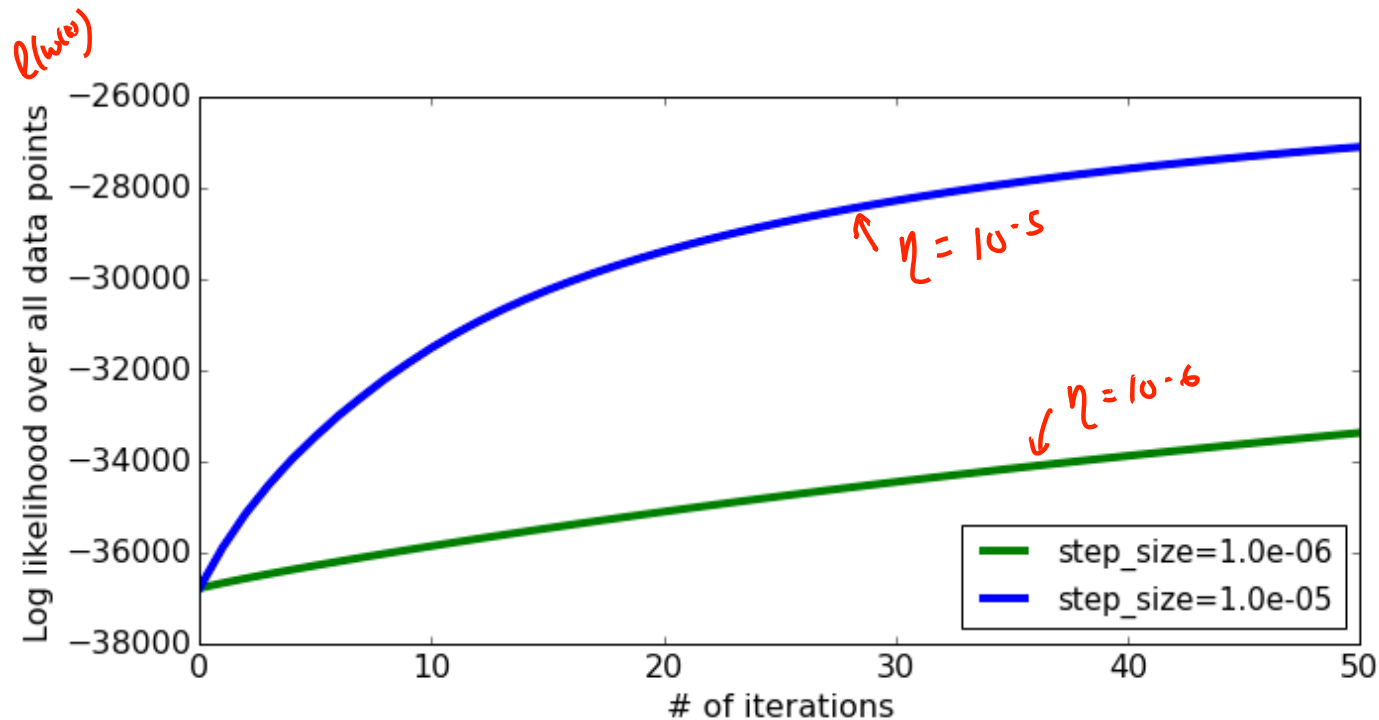step size     $\dfrac{\partial\ell(w^{(t)})}{\partial w_i}$

# Choosing the step size η

Machine Learning Specialization

# Learning curve:
## Plot quality (likelihood) over iterations



$$\ln \prod_{i=1}^{N} P(y_i \mid x_i, w^{(t)})$$

Log likelihood over all data points

−26000
−28000
−30000
−32000
−34000
−36000
−38000

0    10    20    30    40    50

# of iterations

$t$

$\eta = 10^{-5}$

— step_size=1.0e-05

→ need more than 50 iterations

# If step size is too small, can take a long time to converge

Machine Learning Specialization

# Compare converge with different step sizes



*higher is better*

$\eta = 10^{-5}$

$\eta = 1.5 \; 10^{-5}$

smooth faster progress

Early oscillation

Log likelihood over all data points

# of iterations

step_size=1.0e-05
step_size=1.5e-05

Machine Learning Specialization

# Careful with step sizes that are too large

Machine Learning Specialization

# Very large step sizes can even cause divergence or wild oscillations

Machine Learning Specialization

# Simple rule of thumb for picking step size η

- Unfortunately, picking step size requires a lot of trial and error ☹

- Try a several values, <u>exponentially spaced</u>
  - **Goal**: plot learning curves to
    - find one η <u>that is too small</u> (smooth but moving too slowly)
    - find one η <u>that is too large</u> (oscillation or divergence)

- Try values in between to find "best" η

  ↳ exponentially space, pick one that leads best training data likelihood

- *Advanced tip*: can also try step size that decreases with iterations, e.g.,

$$\eta_t = \frac{\eta_0}{t}$$

# Summary of logistic regression classifier

# What you can do now...

- Measure quality of a classifier using the likelihood function

- Interpret the likelihood function as the probability of the observed data

- Learn a logistic regression model with gradient descent