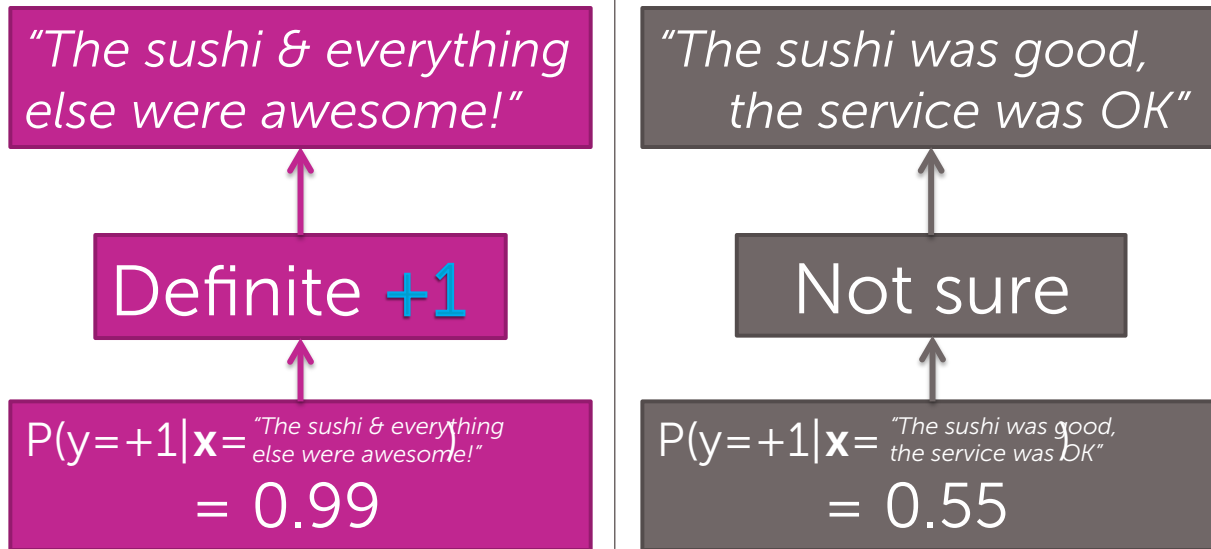# Linear classifiers:
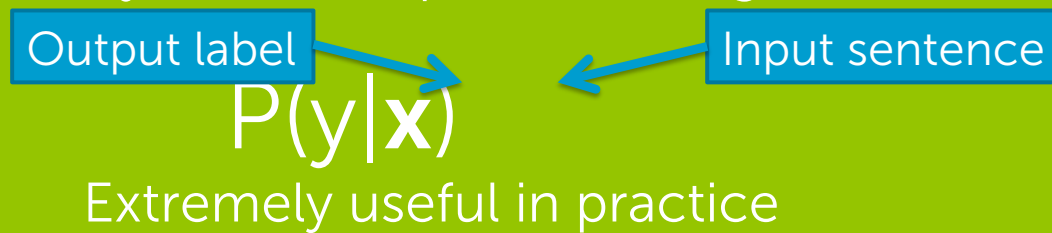# Parameter learning

Emily Fox & Carlos Guestrin

Machine Learning Specialization

University of Washington

# 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:

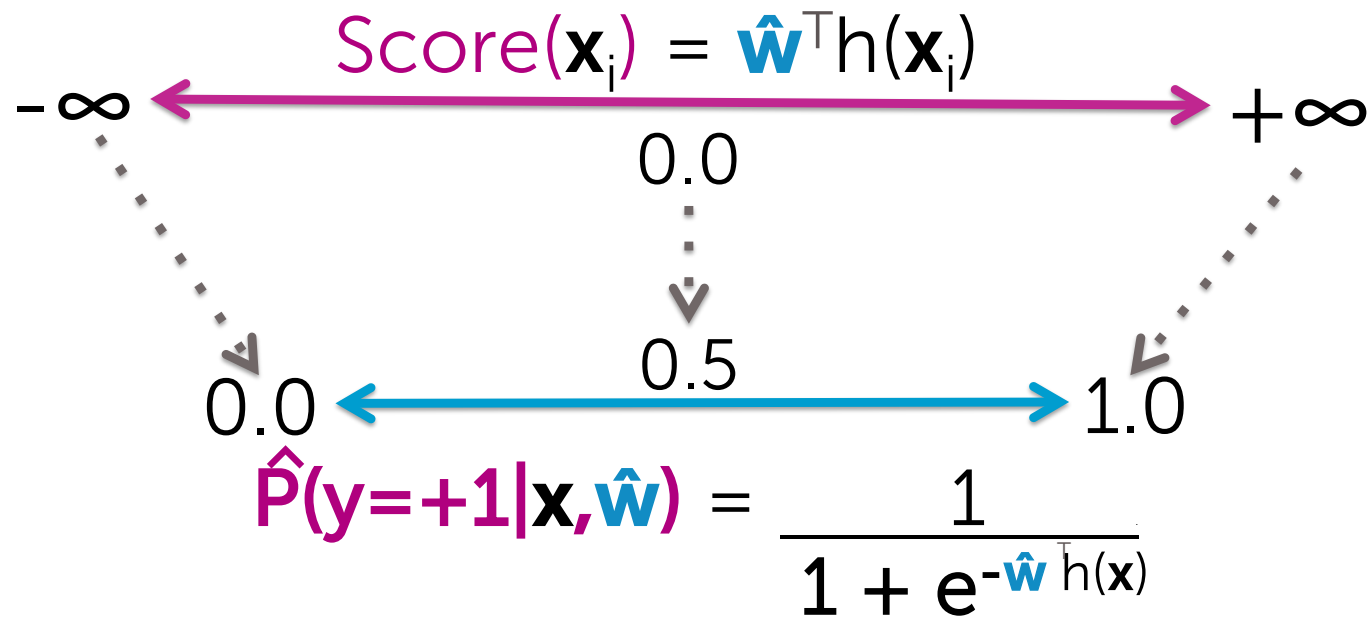Output label

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

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

Machine Learning Specialization

# Quality metric for logistic regression: Maximum likelihood estimation

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



Training Data

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

**y**

**ŵ**

ML algorithm

**Quality metric**

Machine Learning Specialization

# Learning problem

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

| $\mathbf{x}[1]$ = #awesome | $\mathbf{x}[2]$ = #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

$\hat{\mathbf{w}}$

Machine Learning Specialization

# MOVE TO HEAD SHOT

# Finding best coefficients

| 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 |

# 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 |
| 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 |
| 1 | 1 | +1 |
|  |  |  |
|  |  |  |
|  |  |  |
| 2 | 1 | +1 |

Machine Learning Specialization

# 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

Machine Learning Specialization

# 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

0   1   2   3   4   ...

Machine Learning Specialization

# Data likelihood

Machine Learning Specialization

# 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:

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

| **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:

$P(y=-1 \mid x_2, w)$

Machine Learning Specialization

# 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

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

The reason you multiply is that you assume that every row is independent of each other.

# 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 | $y: +1$ | $P(y=+1 \mid x[1]=2, x[2]=1, \mathbf{w})$ |
| $\mathbf{x}_2, y_2$ | 0 | 2 | -1 | $P(y=-1 \mid x[1]=0, x[2]=2, \mathbf{w})$ |
| $\mathbf{x}_3, y_3$ | 3 | 3 | -1 | $P(y=-1 \mid x[1]=3, x[2]=3, \mathbf{w})$ |
| $\mathbf{x}_4, y_4$ | 4 | 1 | +1 | $P(y=+1 \mid x[1]=4, x[2]=1, \mathbf{w})$ |

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

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

*num. of data points* $\longrightarrow$

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

$\Leftarrow$ *pick w to make this fn. as large as possible*

Machine Learning Specialization

# MOVE TO FULL BODY SHOT

# 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**

**ŵ**

ML algorithm

Quality metric

# MOVE TO HEAD SHOT

# Find "best" classifier

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

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

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

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

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

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

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

# MOVE TO FULL BODY SHOT

# Review of gradient ascent

Machine Learning Specialization

# MOVE TO HEAD SHOT

# 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

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

In practice, stop when

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

↑
tolerance

$w^*$

**Algorithm:**

**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



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

Machine Learning Specialization

# Contour plots

Machine Learning Specialization

# Gradient ascent



Algorithm:

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

**while** not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \nabla \ell(\mathbf{w}^{(t)})$$

step size

Machine Learning Specialization

# MOVE TO FULL BODY SHOT

# Learning algorithm for logistic regression

Machine Learning Specialization

# MOVE TO HEAD SHOT

# 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

j ranges from 0 ~ D, num of features
i ranges from 1~N, num of data points

- output = 1 if yi positive
- output = 0 if yi negative

Indicator function:

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

Machine Learning Specialization

# 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 |

$\dfrac{\partial \ell}{\partial w_1}$

$h_1(x) = $ # awesome   parameter w1 multiplies 1st feature h1(x)

| x[1] | x[2] | y | P(y=+1|$x_i$,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:

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

$w_1^{(t+1)} = w_1^{(t)} + \eta \dfrac{\partial \ell(w^{(t)})}{\partial w_1} \Big| \eta = 0.1$

$= 1 + 0.1 \times 1.33 = 1.133$

# Derivative of (log-)likelihood: Interpretation

Sum over data points

Feature value

Difference between truth and prediction

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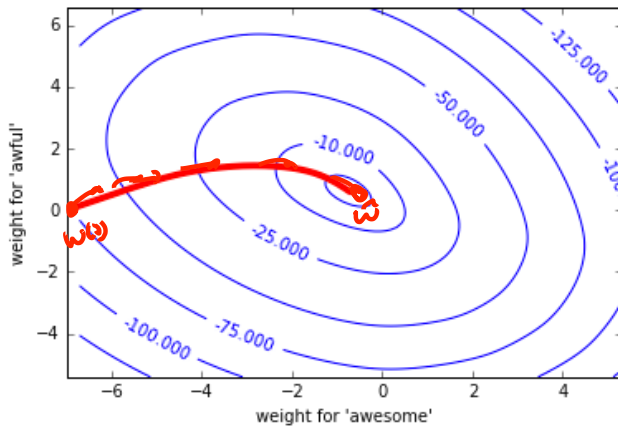
$\Delta_i$

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

increase parameter w
=> increase score
=> increase probability

Machine Learning Specialization

# 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

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

  for $j=0,\ldots,D$ ← coefficient

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

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

$t \leftarrow t+1$

step size

$\dfrac{\partial\ell(w^{(t)})}{\partial w_v}$

# MOVE TO FULL BODY SHOT

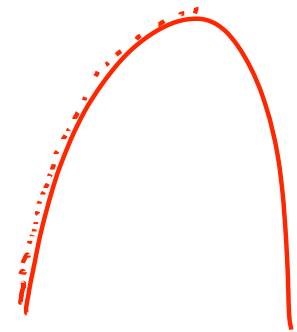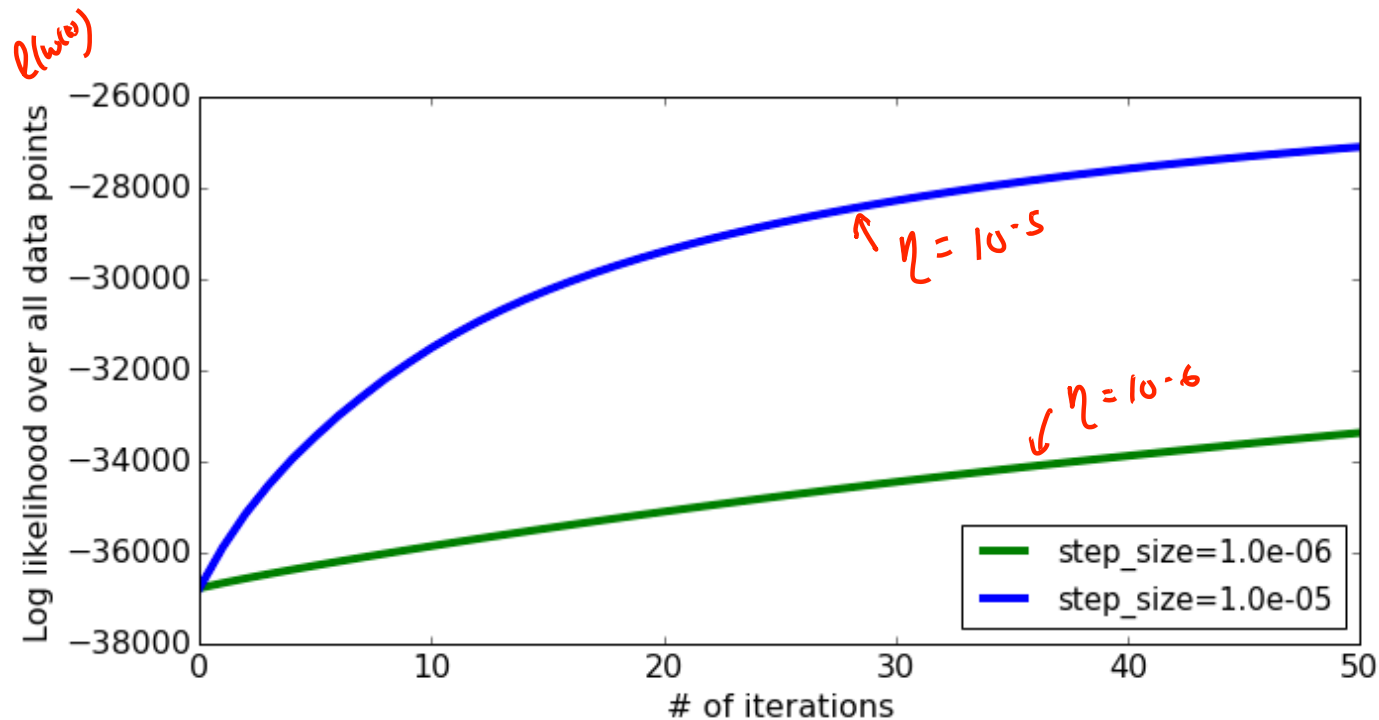# Choosing the step size η

Machine Learning Specialization

# MOVE TO HEAD SHOT

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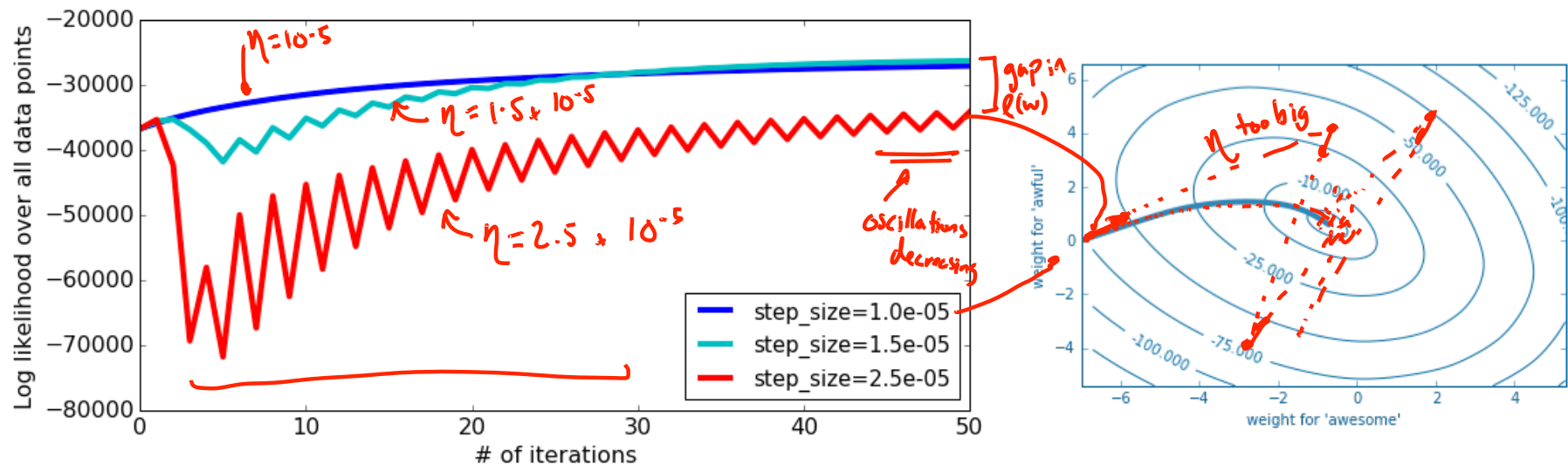


$\ln \prod_{i=1}^{N} P(y_i | x_i, w^{(t)})$ — Log likelihood over all data points (y-axis)

# of iterations — $t$ (x-axis)

$\eta = 10^{-5}$

step_size=1.0e-05

→ need more than 50 iterations

because the curve is still going up

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



$\ell(w^{(t)})$

Log likelihood over all data points

$\eta = 10^{-5}$

$\eta = 10^{-6}$

step_size=1.0e-06
step_size=1.0e-05

# of iterations

Machine Learning Specialization

# Compare converge with different step sizes



*Handwritten annotations:*
- higher is better
- $\eta = 10^{-5}$
- $\eta = 1.5 \cdot 10^{-5}$
- smooth faster progress
- early oscillation

*Plot axes:* Log likelihood over all data points (y-axis: −26000 to −42000), # of iterations (x-axis: 0 to 50)

*Legend:*
- 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 that is too small (smooth but moving too slowly)
    - find one η that is too large (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.,

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

# MOVE TO FULL BODY SHOT

# Deriving the gradient for logistic regression

VERY OPTIONAL

# MOVE TO HEAD SHOT

# Log-likelihood function

- Goal: choose coefficients **w** maximizing likelihood:

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

- Math simplified by using log-likelihood – taking (natural) log:

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

natural
log

Machine Learning Specialization

# The log trick, often used in ML...

- Products become sums:

$$\ln a \cdot b = \ln a + \ln b \qquad \ln \frac{a}{b} = \ln a - \ln b$$

- Doesn't change maximum!
  - If $\hat{\mathbf{w}}$ maximizes f($\mathbf{w}$):

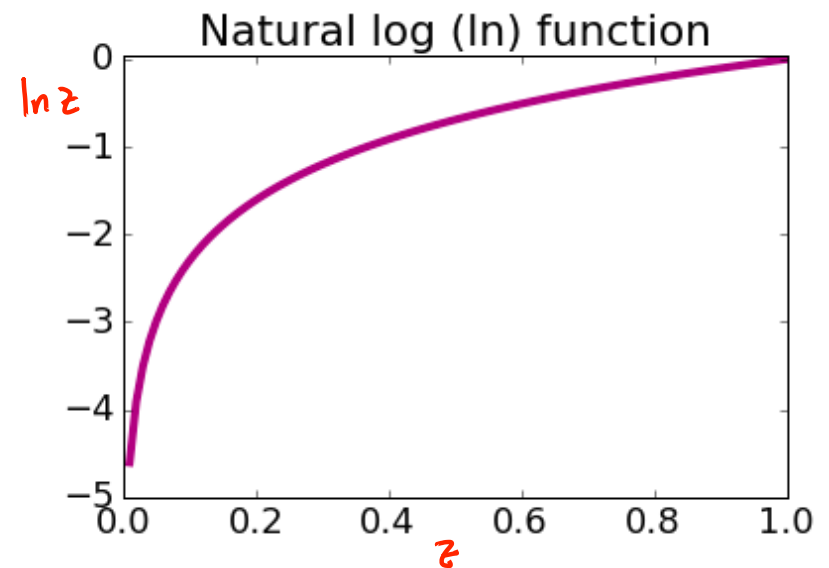$$\hat{w} = \underset{w}{\arg\max}\ f(w)$$

  the w that makes f(w) largest

  - Then $\hat{\mathbf{w}}_{ln}$ maximizes ln(f($\mathbf{w}$)):

$$\hat{w}_{ln} = \underset{w}{\arg\max}\ \ln\left(f(w)\right)$$

$$\hat{w} = \hat{w}_{ln}$$



Natural log (ln) function

$\ln z$

$z$

Machine Learning Specialization

Insert next title slide before Slide 52, around 4:55 in PL7_DerivingtheGradient_1stEdit

53

# Expressing the log-likelihood

**VERY OPTIONAL**

Machine Learning Specialization

# Using log to turn products into sums

$$\ln \prod_{i=1}^{N} f_i = \sum_{i=1}^{N} \ln f_i$$

- The log of the product of likelihoods becomes the sum of the logs:

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

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

# Rewriting log-likelihood

- For simpler math, we'll rewrite likelihood with indicators:

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

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

if $y_i = +1$

if $y_i = -1$

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

Insert next title slide before Slide 54, around 7:33 in PL7_DerivingtheGradient_1stEdit

# Deriving probability that y=-1 given x

VERY OPTIONAL

# Logistic regression model: P(y=-1|x,**w**)

- Probability model predicts y=+1:

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

- Probability model predicts y=-1:

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

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

Machine Learning Specialization

Insert next title slide before Slide 55, around 9:15 in PL7_DerivingtheGradient_1stEdit

# Rewriting the log-likelihood

VERY OPTIONAL

Machine Learning Specialization

# Plugging in logistic function for 1 data point

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

$$\ell\ell(\mathbf{w}) = \mathbb{1}[y_i = +1] \ln P(y = +1 \mid \mathbf{x}_i, \mathbf{w}) + \mathbb{1}[y_i = -1] \ln P(y = -1 \mid \mathbf{x}_i, \mathbf{w})$$

$$= \mathbb{1}[y_i = +1] \ln \frac{1}{1 + e^{-w^\top h(x_i)}} + \left(1 - \mathbb{1}[y_i = +1]\right) \ln \frac{e^{-w^\top h(y_i)}}{1 + e^{-w^\top h(x_i)}}$$

$$= -\mathbb{1}[y_i = +1] \ln\left(1 + e^{-w^\top h(x_i)}\right) + \left(1 - \mathbb{1}[y_i = +1]\right)\left[- w^\top h(x_i) - \ln\left(1 + e^{-w^\top h(x_i)}\right)\right]$$

$$= -\left(1 - \mathbb{1}[y_i = +1]\right) w^\top h(x_i) - \ln\left(1 + e^{-w^\top h(x_i)}\right)$$

Simpler form

$$\ln e^a = a$$

$$\mathbb{1}[y_i = -1] = 1 - \mathbb{1}[y_i = +1]$$

$$\ln \frac{1}{1 + e^{-w^\top h(x_i)}} = -\ln\left(1 + e^{-w^\top h(x_i)}\right)$$

$$\ln \frac{e^{-w^\top h(x_i)}}{1 + e^{-w^\top h(x_i)}} =$$

$$\ln e^{-w^\top h(x_i)} - \ln\left(1 + e^{-w^\top h(x_i)}\right)$$

$$-w^\top h(x_i)$$

Machine Learning Specialization

Insert next title slide before Slide 56, around 16:56 in PL7_DerivingtheGradient_1stEdit

# Deriving gradient of log-likelihood

VERY OPTIONAL

Machine Learning Specialization

# Gradient for 1 data point

$$\ell\ell(\mathbf{w}) = -(1 - \mathbb{1}[y_i = +1])\mathbf{w}^\top h(\mathbf{x}_i) - \ln\left(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}\right)$$

$$\frac{\partial \ell\ell}{\partial w_j} = -(1 - \mathbb{1}[y_i = +1]) \frac{\partial}{\partial w_j} \mathbf{w}^\top h(\mathbf{x}_i) - \frac{\partial}{\partial w_j} \ln\left(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}\right)$$

$$= -(1 - \mathbb{1}[y_i = +1]) h_j(\mathbf{x}_i) + h_j(\mathbf{x}_i) P(y = -1 \mid \mathbf{x}_i, \mathbf{w})$$

$$= h_j(\mathbf{x}_i)\left[\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w})\right]$$

$$\frac{\partial}{\partial w_j} \mathbf{w}^\top h(\mathbf{x}_i) = h_j(\mathbf{x}_i)$$

$$\frac{\partial}{\partial w_j} \ln\left(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}\right)$$

$$= -h_j(\mathbf{x}_i) \underbrace{\frac{e^{-\mathbf{w}^\top h(\mathbf{x}_i)}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}}}_{P(y = -1 \mid \mathbf{x}_i, \mathbf{w})}$$

Machine Learning Specialization

# Finally, gradient for all data points

- Gradient for one data point:

$$h_j(\mathbf{x}_i)\Big(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w})\Big)$$

- Adding over data points:

$$\frac{\partial \ell\ell}{\partial w_j} = \sum_{i=1}^{N} h_j(x_i)\Big(\mathbb{1}[y_i = +1] - P(y = +1 \mid x_i, w)\Big)\Bigg\}$$
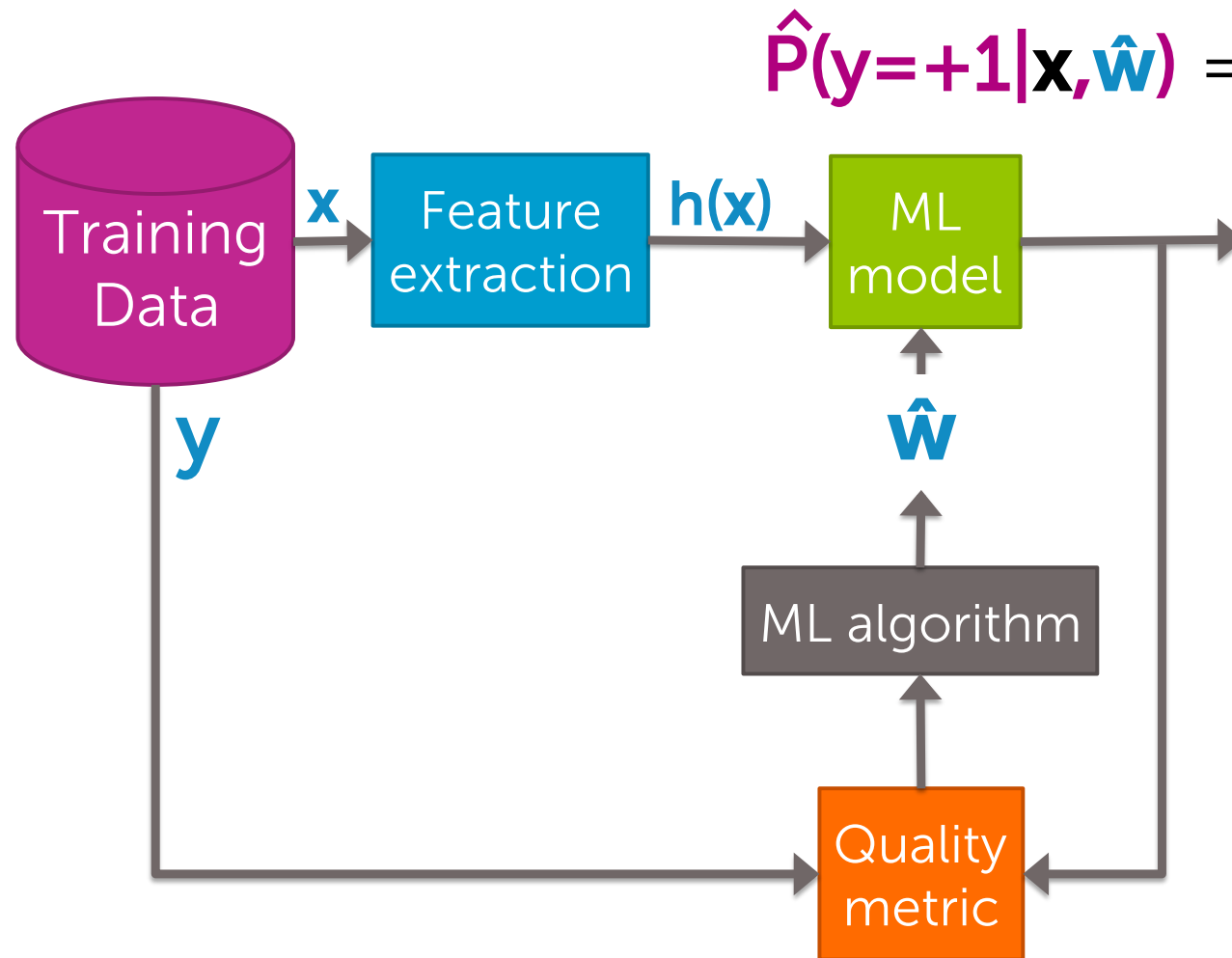
Machine Learning Specialization

# MOVE TO FULL BODY SHOT

# Summary of logistic regression classifier

# MOVE TO HEAD SHOT

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

Training Data

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

**y**

$\hat{\mathbf{w}}$

ML algorithm

Quality metric
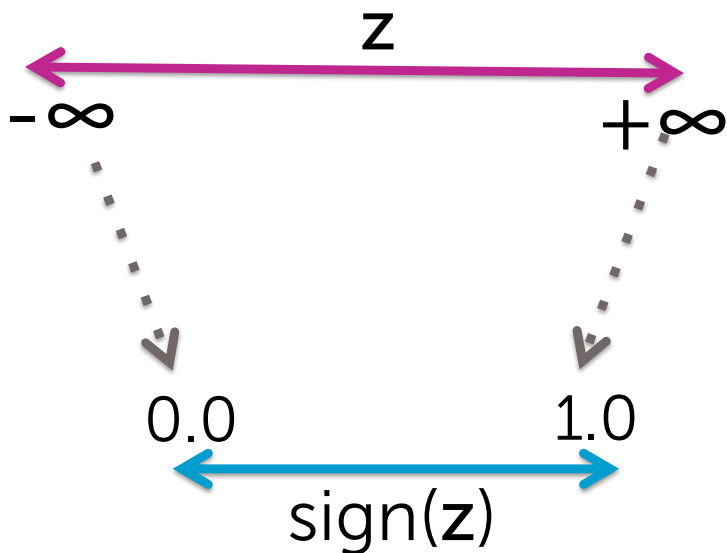
Machine Learning Specialization
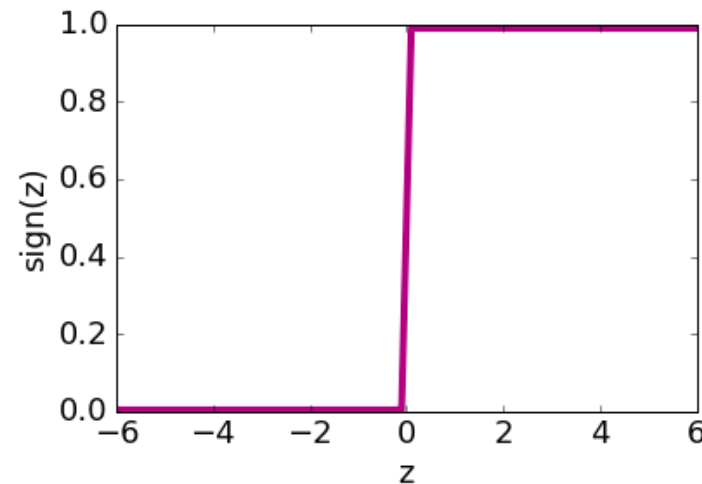
# MOVE TO FULL BODY SHOT

# 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

- (Optional) Derive the gradient descent update rule for logistic regression

Machine Learning Specialization

# Simplest link function: sign(**z**)



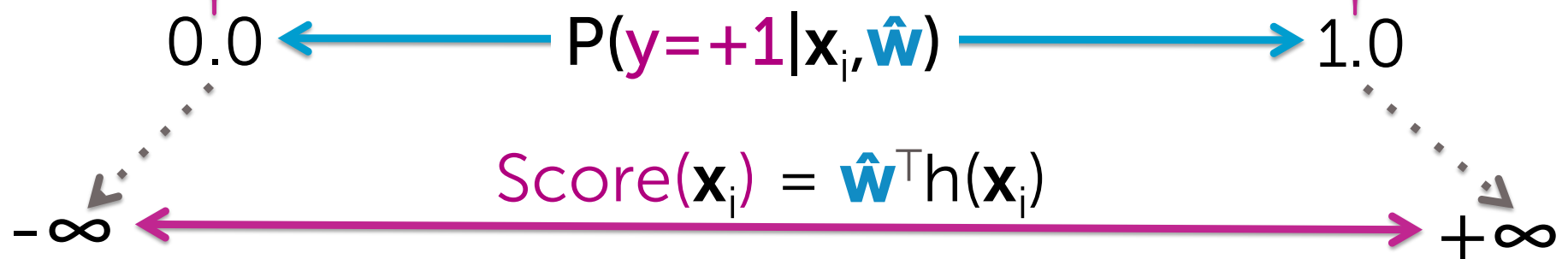$$sign(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

But, sign(z) only outputs –1 or +1, no probabilities in between

Machine Learning Specialization

# 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 |

$$0.0 \longleftarrow P(y=+1|x_i, \hat{w}) \longrightarrow 1.0$$

$$\text{Score}(x_i) = \hat{w}^\top h(x_i)$$

$$-\infty \longleftarrow \longrightarrow +\infty$$

# Quality metric: probability of data

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

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

If model good, should predict

Increase probability y=+1 when

Choose **w** to make

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

If model good, should predict

Increase probability y=-1 when

Choose **w** to make

Machine Learning Specialization

# Maximizing likelihood (probability of data)

| Data point | x[1] | x[2] | y | Choose **w** to maximize |
|------------|------|------|-----|--------------------------|
| **x**₁,y₁ | 2 | 1 | +1 | |
| **x**₂,y₂ | 0 | 2 | -1 | |
| **x**₃,y₃ | 3 | 3 | -1 | |
| **x**₄,y₄ | 4 | 1 | +1 | |
| **x**₅,y₅ | 1 | 1 | +1 | |
| **x**₆,y₆ | 2 | 4 | -1 | |
| **x**₇,y₇ | 0 | 3 | -1 | |
| **x**₈,y₈ | 0 | 1 | -1 | |
| **x**₉,y₉ | 2 | 1 | +1 | |

Must combine into single measure of quality

Machine Learning Specialization

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

- Choose coefficients **w** that maximize likelihood:

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

- No closed-form solution ➜ use gradient ascent

Machine Learning Specialization