



Linear classifiers: Logistic regression



Predicting sentiment by topic: *An intelligent restaurant review system*

It's a big day & I want to book a table at
a nice Japanese restaurant

You are looking for
★★★★
sushi restaurants



What are people
saying about
the food?
the ambiance?...



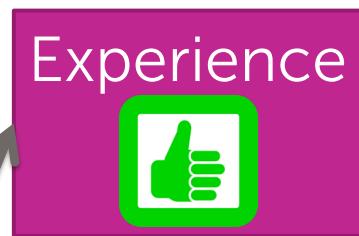


Positive reviews not positive about everything



Sample review:

Watching the chefs create incredible edible art made the experience very unique.



My wife tried their ramen and it was pretty forgettable.



All the sushi was delicious! Easily best sushi in city.



Classifying sentiment of review

Easily best sushi in city.

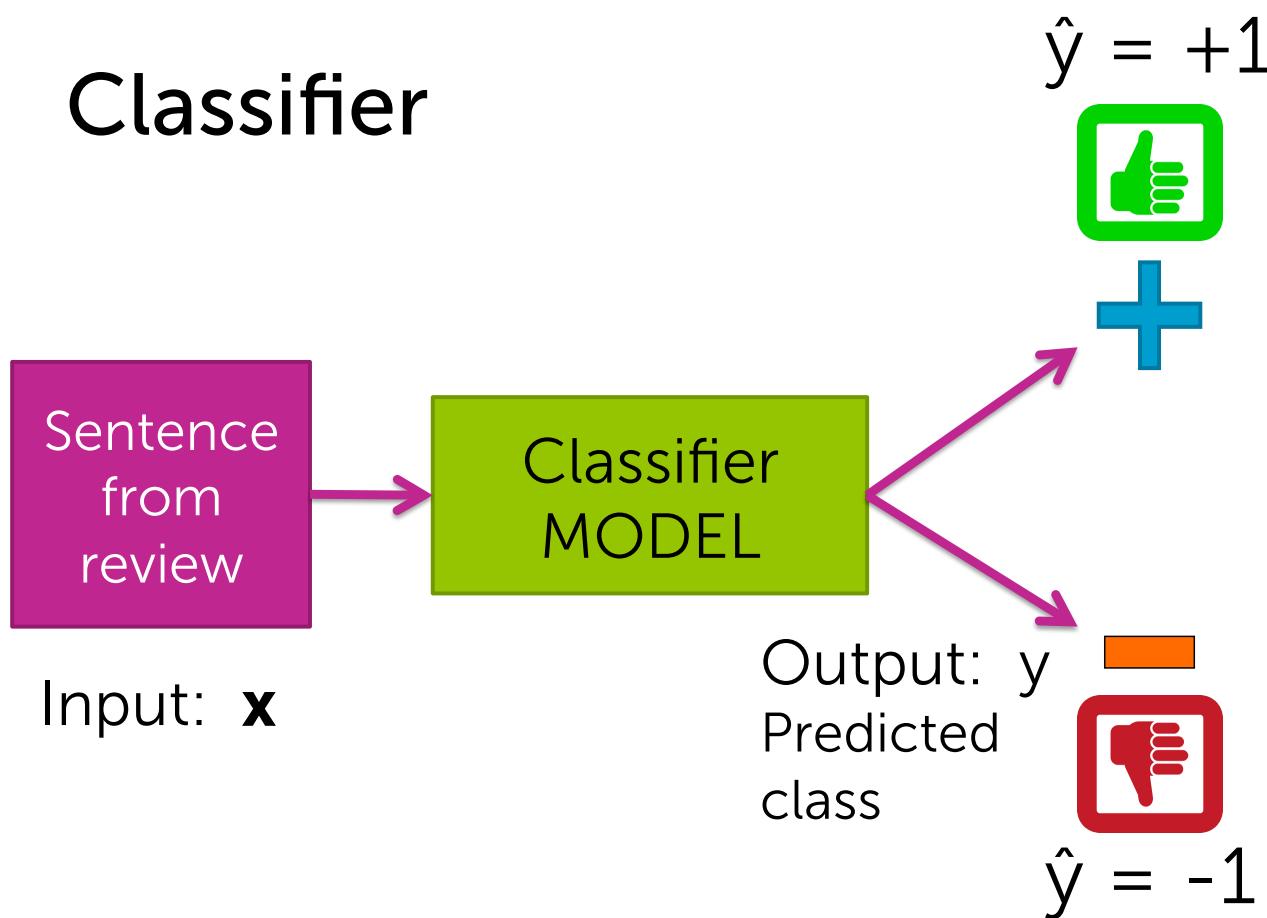


Sentence Sentiment
Classifier



Linear classifier: Intuition

Classifier



Note: we'll start talking about 2 classes, and address multiclass later

A (linear) classifier

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

Word	Coefficient
good	1.0
great	1.5
awesome	2.7
bad	-1.0
terrible	-2.1
awful	-3.3
restaurant, the, we, where, ...	0.0
...	...

Scoring a sentence

Word	Coefficient
good	1.0
great	1.2
awesome	1.7
bad	-1.0
terrible	-2.1
awful	-3.3
restaurant, the, we, where, ...	0.0
...	...

Input \mathbf{x}_i :

Sushi was great,
the food was awesome,
but the service was terrible.

Called a linear classifier, because output is weighted sum of input.

Sentence
from
review

Input: \mathbf{x}

Word	Coefficient
...	...



Simple linear classifier

$\text{Score}(\mathbf{x})$ = weighted count of
words in sentence

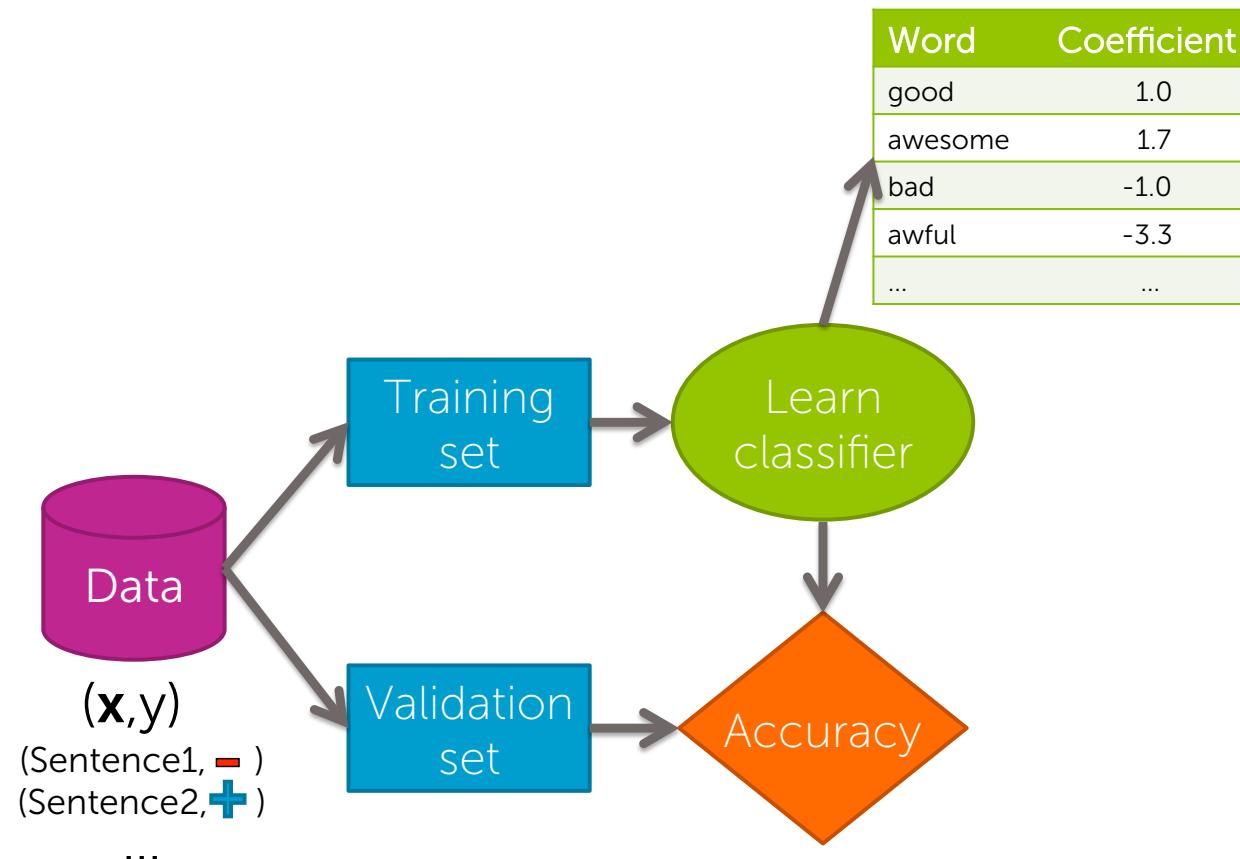
If $\text{Score}(\mathbf{x}) > 0$:

$\hat{y} = +1$

Else:

$\hat{y} = -1$

Training a classifier = Learning the coefficients





Decision boundaries



Suppose only two words had non-zero coefficient

Word	Coefficient
#awesome	1.0
#awful	-1.5

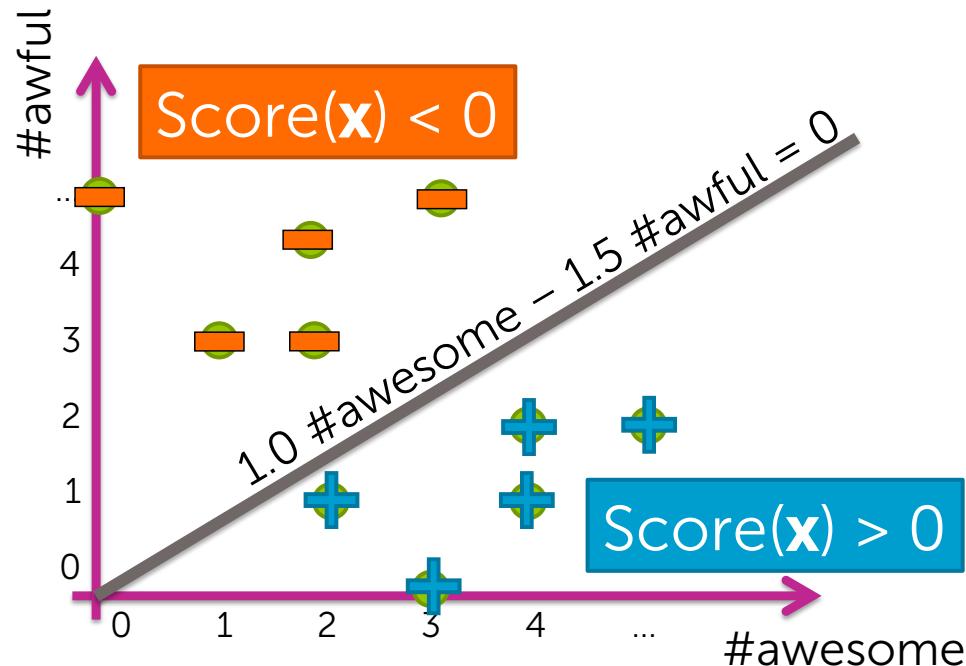
$$\text{Score}(x) = 1.0 \text{ #awesome} - 1.5 \text{ #awful}$$



Decision boundary example

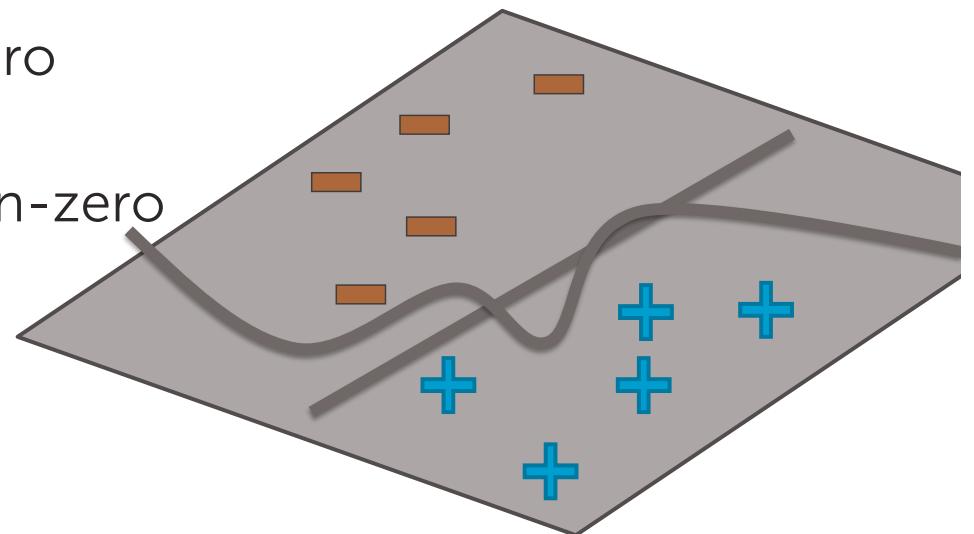
Word	Coefficient
#awesome	1.0
#awful	-1.5

$$\rightarrow \text{Score}(x) = 1.0 \text{ #awesome} - 1.5 \text{ #awful}$$

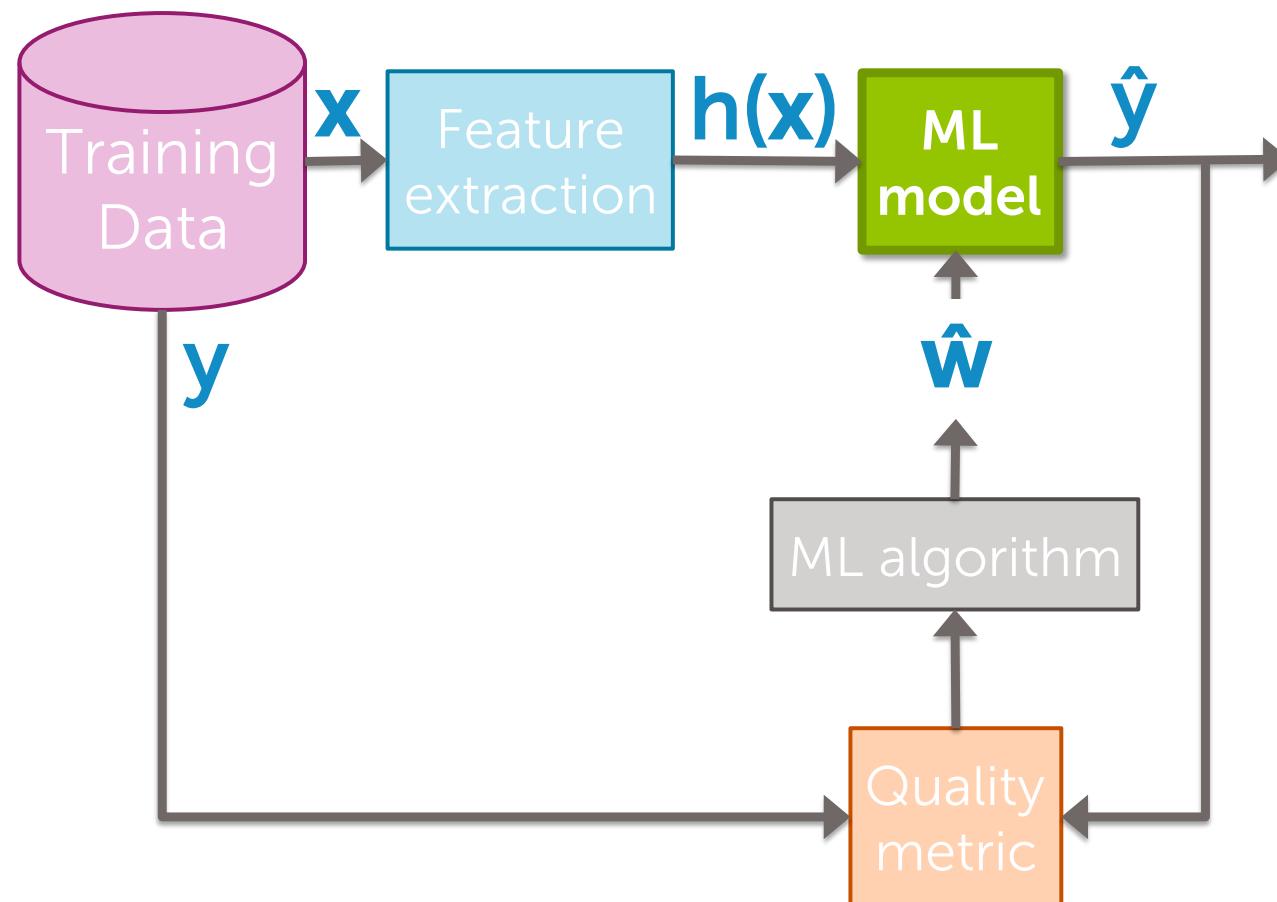


Decision boundary separates positive & negative predictions

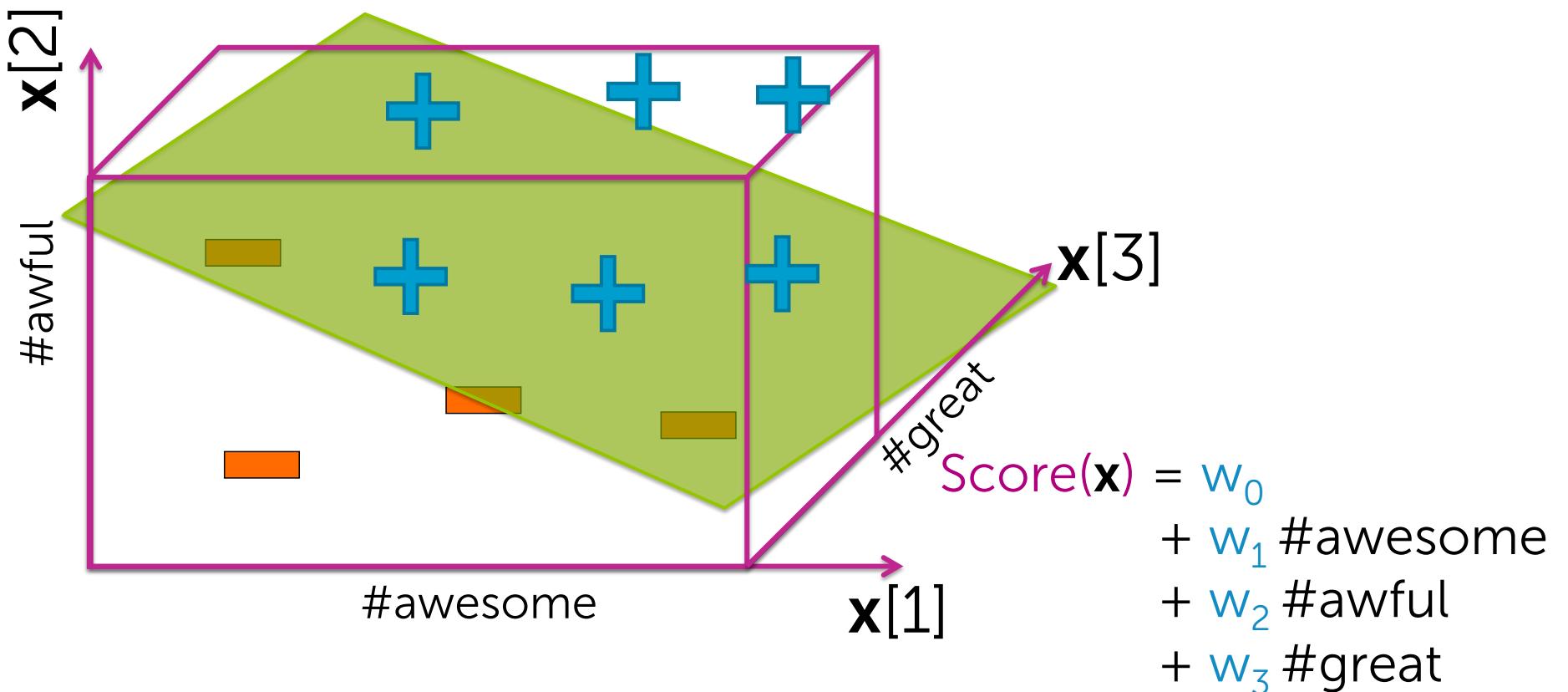
- For linear classifiers:
 - When 2 coefficients are non-zero
→ line
 - When 3 coefficients are non-zero
→ plane
 - When many coefficients are non-zero
→ hyperplane
- For more general classifiers
→ more complicated shapes



Linear classifier: Model



Coefficients of classifier



General notation

Output: $y \leftarrow \{-1, +1\}$

Inputs: $\mathbf{x} = (\mathbf{x}[1], \mathbf{x}[2], \dots, \mathbf{x}[d])$



d-dim vector

Notational conventions:

$\mathbf{x}[j] = j^{\text{th}}$ input (*scalar*)

$h_j(\mathbf{x}) = j^{\text{th}}$ feature (*scalar*)

$\mathbf{x}_i =$ input of i^{th} data point (*vector*)

$\mathbf{x}_i[j] = j^{\text{th}}$ input of i^{th} data point (*scalar*)

Simple hyperplane

Model: $\hat{y}_i = \text{sign}(\text{Score}(x_i))$

$$\text{Score}(x_i) = w_0 + w_1 x_i[1] + \dots + w_d x_i[d] = \mathbf{w}^\top \mathbf{x}_i$$

feature 1 = 1

feature 2 = $x_i[1]$... e.g., #awesome

feature 3 = $x_i[2]$... e.g., #awful

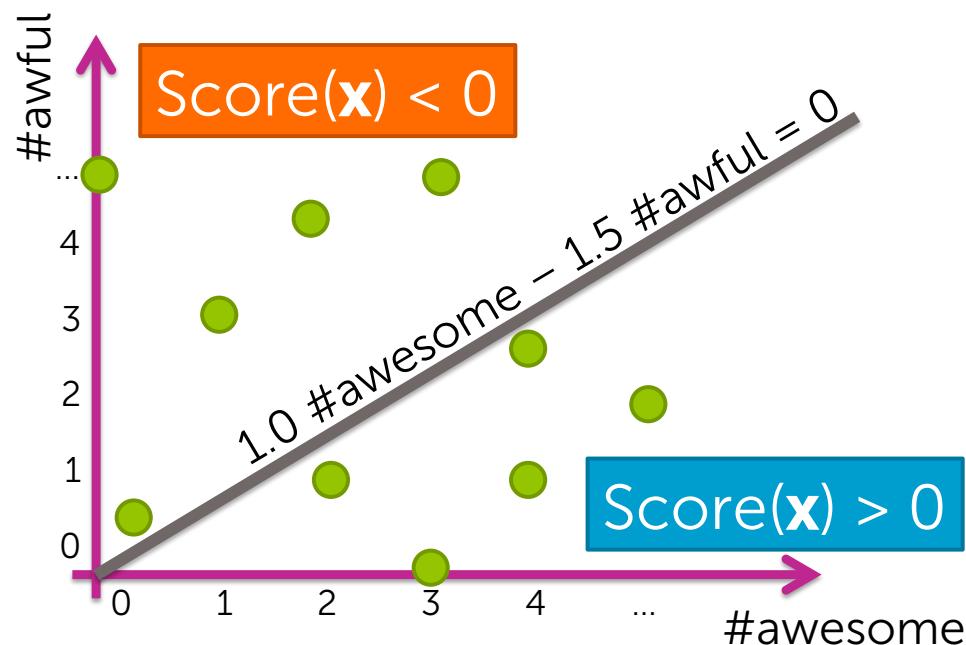
...

feature $d+1 = x_i[d]$... e.g., #ramen

Decision boundary: effect of changing coefficients

Input	Coefficient	Value
	w_0	0.0
#awesome	w_1	1.0
#awful	w_2	-1.5

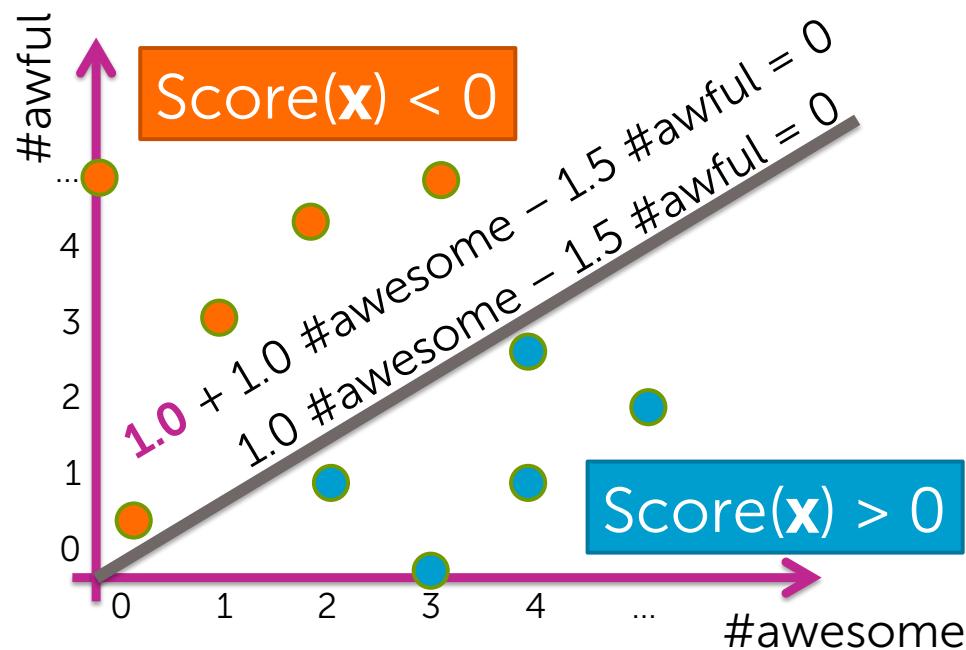
$$\rightarrow \text{Score}(x) = 1.0 \text{ #awesome} - 1.5 \text{ #awful}$$



Decision boundary: effect of changing coefficients

Input	Coefficient	Value
	w_0	1.0
#awesome	w_1	1.0
#awful	w_2	-1.5

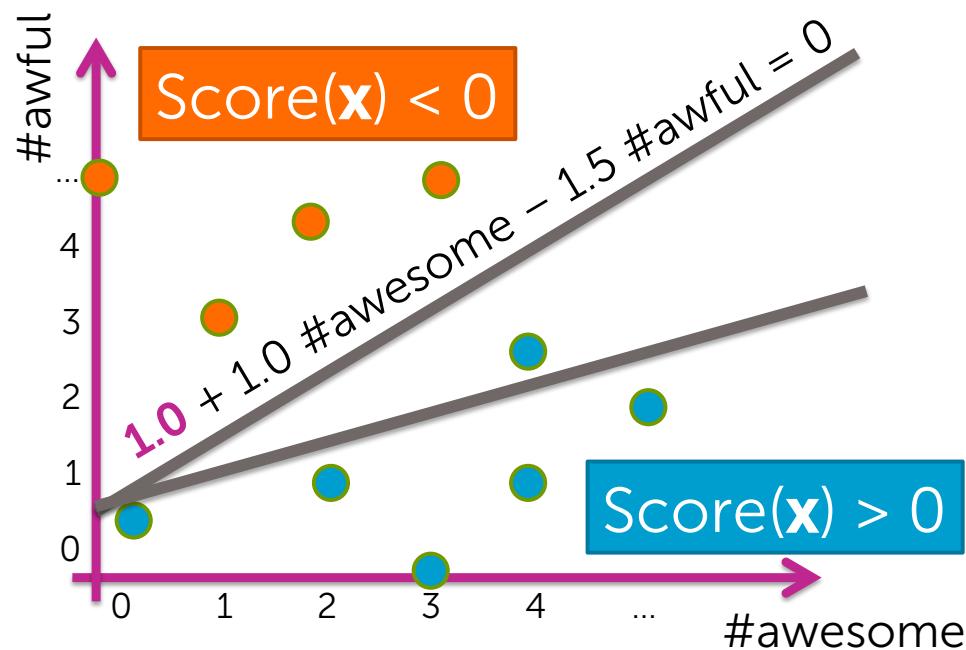
$$\rightarrow \text{Score}(x) = 1.0 \text{ #awesome} - 1.5 \text{ #awful}$$



Decision boundary: effect of changing coefficients

Input	Coefficient	Value
	w_0	1.0
#awesome	w_1	1.0
#awful	w_2	-3.0

$$\rightarrow \text{Score}(x) = 1.0 + 1.0 \text{ #awesome} - 3.0 \text{ #awful}$$



More generic features... D-dimensional hyperplane

Model: $\hat{y}_i = \text{sign}(\text{Score}(\mathbf{x}_i))$

$$\text{Score}(\mathbf{x}_i) = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i)$$

$$= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{h}(\mathbf{x}_i)$$

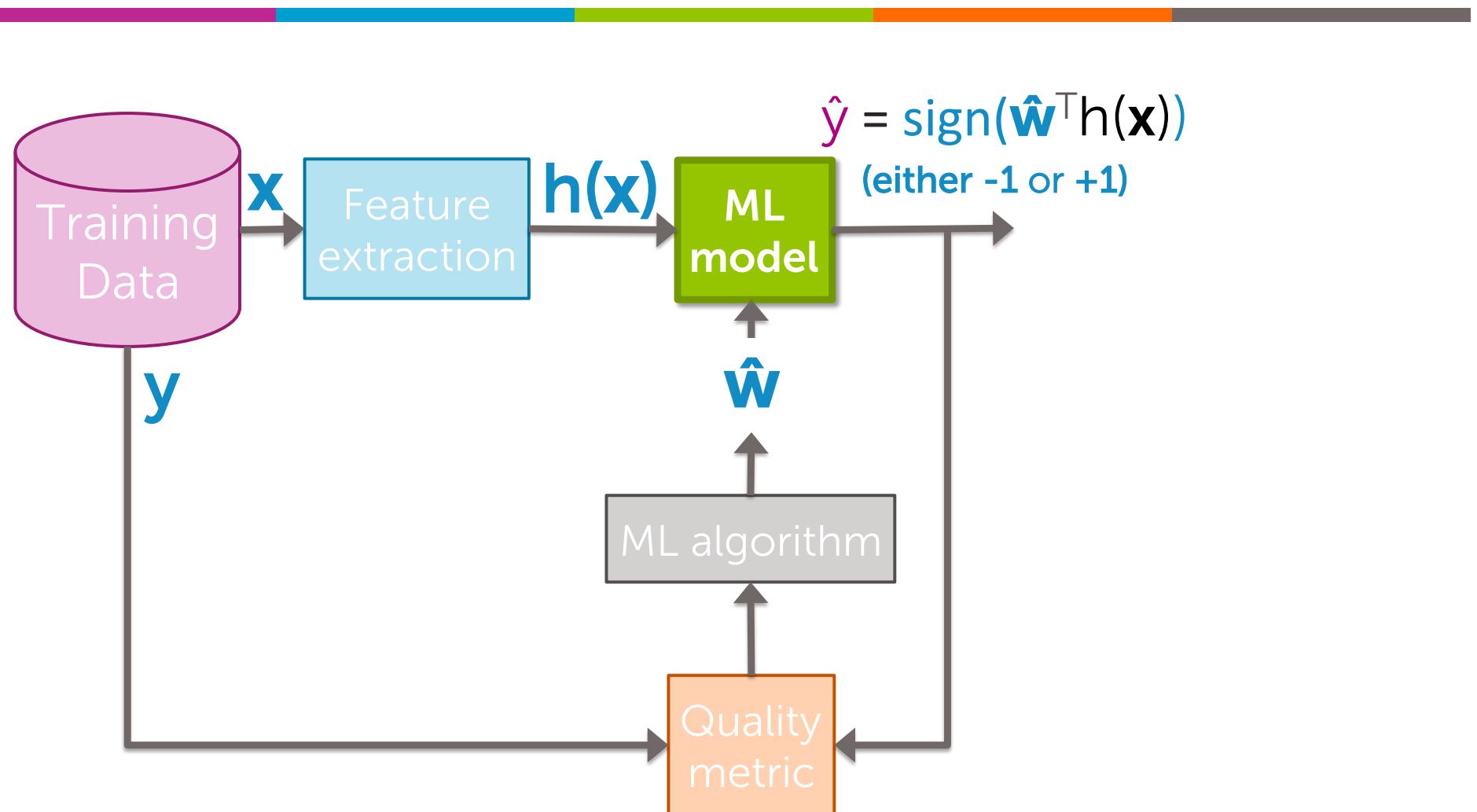
feature 1 = $h_0(\mathbf{x})$... e.g., 1

feature 2 = $h_1(\mathbf{x})$... e.g., $\mathbf{x}[1] = \#\text{awesome}$

feature 3 = $h_2(\mathbf{x})$... e.g., $\mathbf{x}[2] = \#\text{awful}$
or, $\log(\mathbf{x}[7]) \mathbf{x}[2] = \log(\#\text{bad}) \times \#\text{awful}$
or, tf-idf("awful")

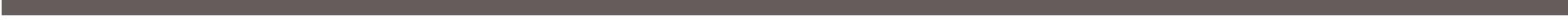
...

feature $D+1 = h_D(\mathbf{x})$... some other function of $\mathbf{x}[1], \dots, \mathbf{x}[d]$



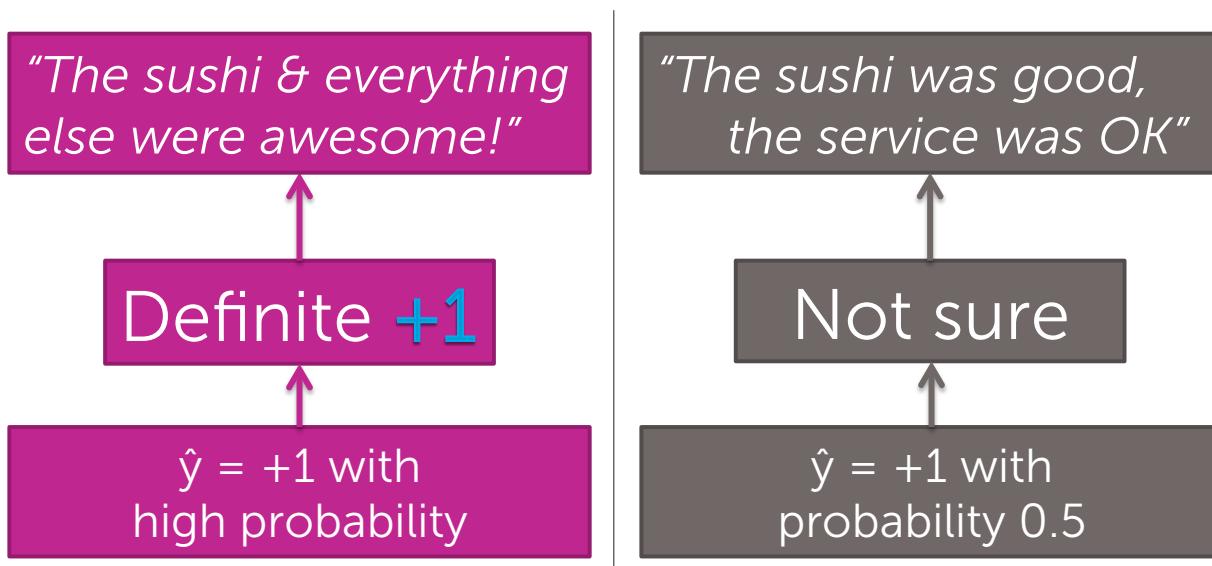


Are you sure about the prediction?
Class probability



How confident is your prediction?

- Thus far, we've outputted a prediction **+1** or **-1**
- But, how sure are you about the prediction?

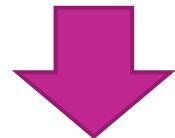




Basics of probabilities – quick review

Basic probability

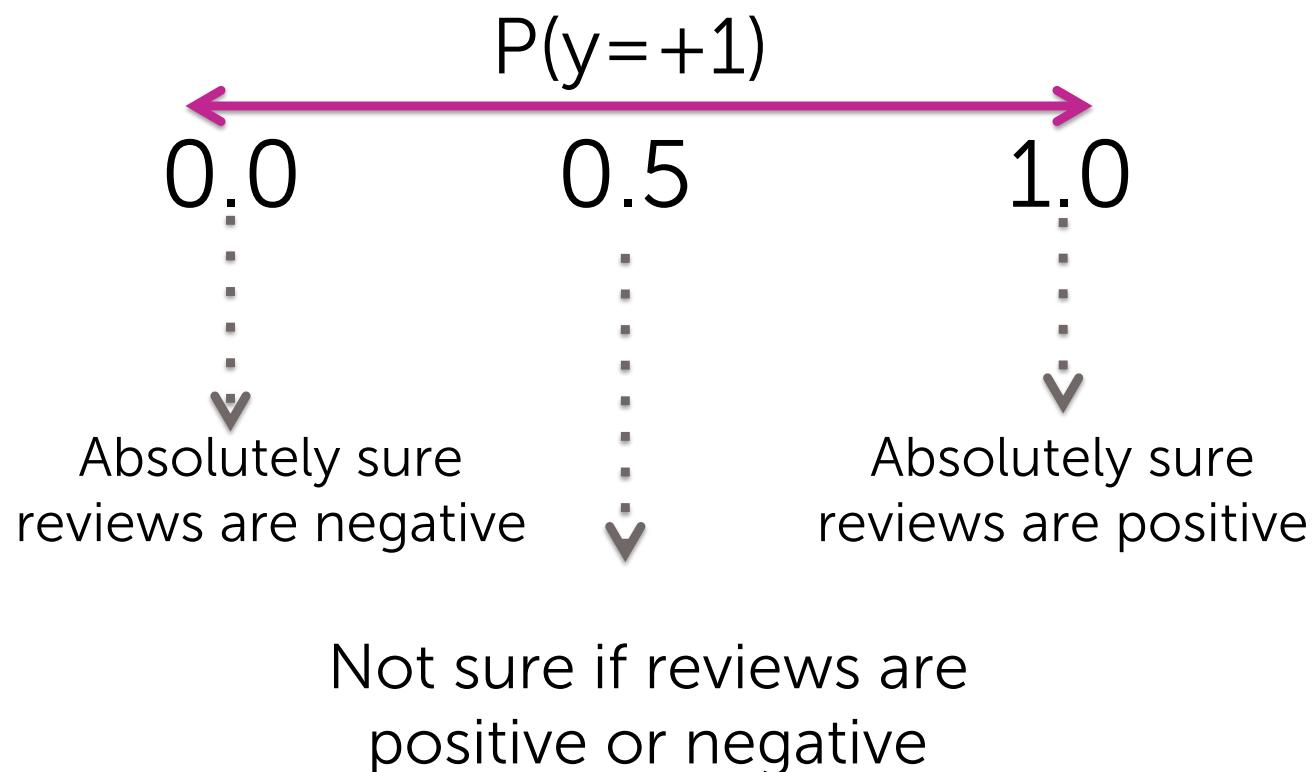
Probability a review is positive is 0.7



$x =$ review text	$y =$ sentiment
All the sushi was delicious! Easily best sushi in Seattle.	+1
The sushi & everything else were awesome!	+1
My wife tried their ramen, it was pretty forgettable.	-1
The sushi was good, the service was OK	+1
...	...

I expect 70% of rows
to have $y = +1$
(Exact number will vary
for each specific dataset)

Interpreting probabilities as degrees of belief



Key properties of probabilities

Property	Two class (e.g., y is +1 or -1)	Multiple classes (e.g., y is dog, cat or bird)
Probabilities always between 0 & 1		
Probabilities sum up to 1		

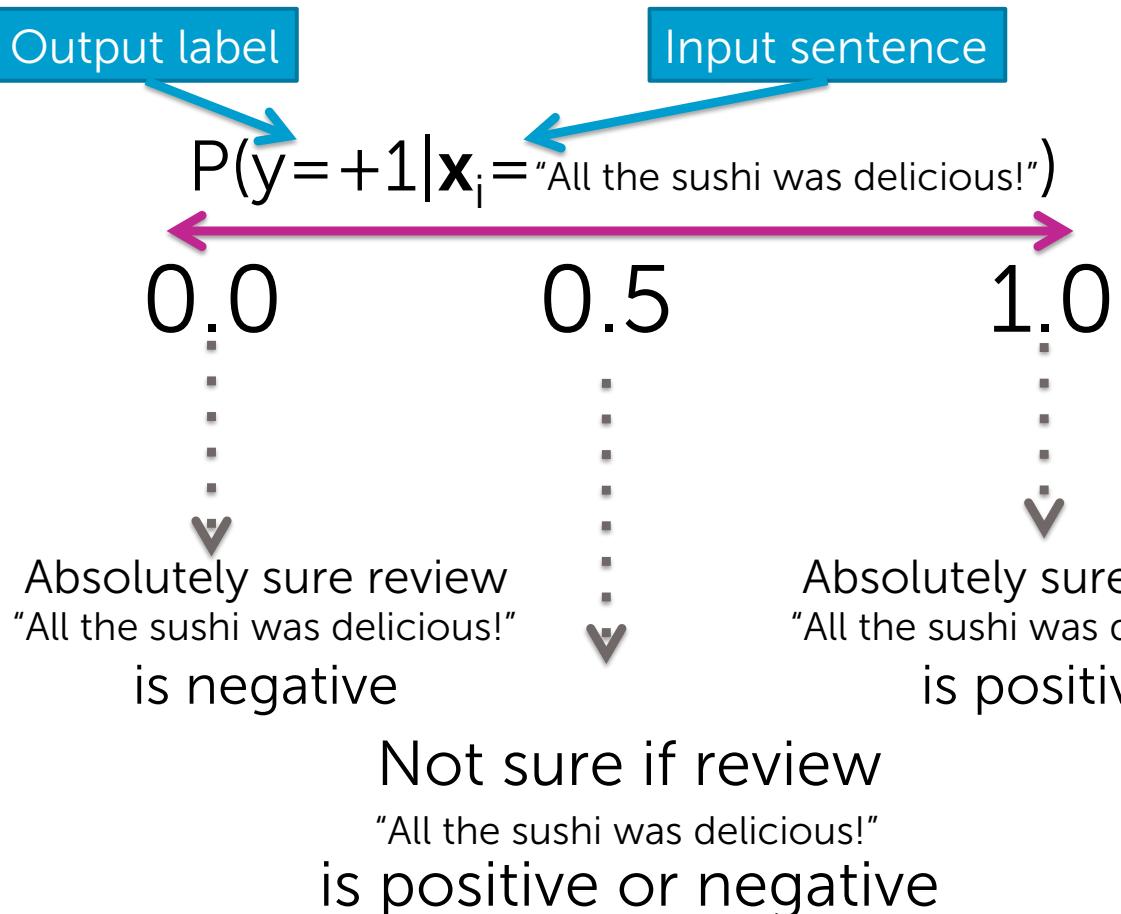
Conditional probability

Probability a review with
3 “awesome” and 1 “awful” is positive is 0.9

x = review text	y = sentiment
All the sushi was delicious! Easily best sushi in Seattle.	+1
Sushi was awesome & everything else was awesome! The service was awful , but overall awesome place!	+1
My wife tried their ramen, it was pretty forgettable.	-1
The sushi was good, the service was OK	+1
...	...
awesome ... awesome ... awful ... awesome	+1
...	...
awesome ... awesome ... awful ... awesome	-1
...	...
...	...
awesome ... awesome ... awful ... awesome	+1

I expect 90% of rows with reviews containing 3 “awesome” & 1 “awful” to have $y = +1$
(Exact number will vary for each specific dataset)

Interpreting conditional probabilities



Key properties of conditional probabilities

Property	Two class (e.g., y is +1 or -1, x_i is review text)	Multiple classes (e.g., y is dog, cat or bird, x_i is image)
Conditional probabilities always between 0 & 1		
Conditional probabilities sum up to 1 over y , but not over x		

Using probabilities in classification

How confident is your prediction?

"The sushi & everything else were awesome!"

Definite +1

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

"The sushi was good, the service was OK"

Not sure

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

Many classifiers provide a degree of certainty:

Output label

$$P(y|x)$$

Input sentence

Extremely useful in practice

Goal: Learn conditional probabilities from data

Training data: N observations (x_i, y_i)

$x[1] = \#awesome$	$x[2] = \#awful$	$y = \text{sentiment}$
2	1	+1
0	2	-1
3	3	-1
4	1	+1
...

Optimize **quality metric**
on training data

Find best model \hat{P}
by finding best \hat{w}

Useful for
predicting \hat{y}

Sentence
from
review

Input: \mathbf{x}

Predict most likely class

$\hat{P}(y|x)$ = estimate of class probabilities

If $\hat{P}(y=+1|x) > 0.5$:

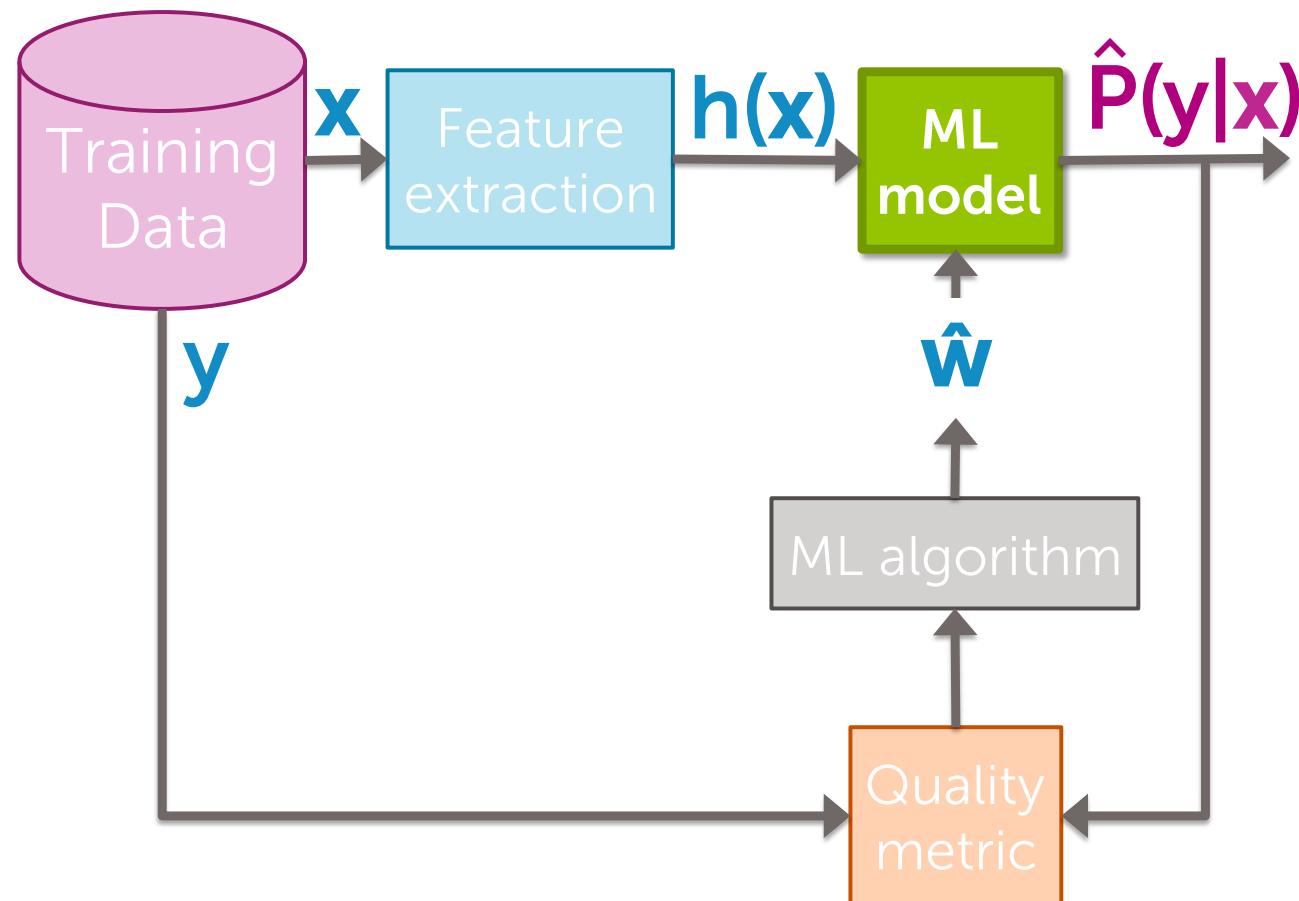
$\hat{y} = +1$

Else:

$\hat{y} = -1$

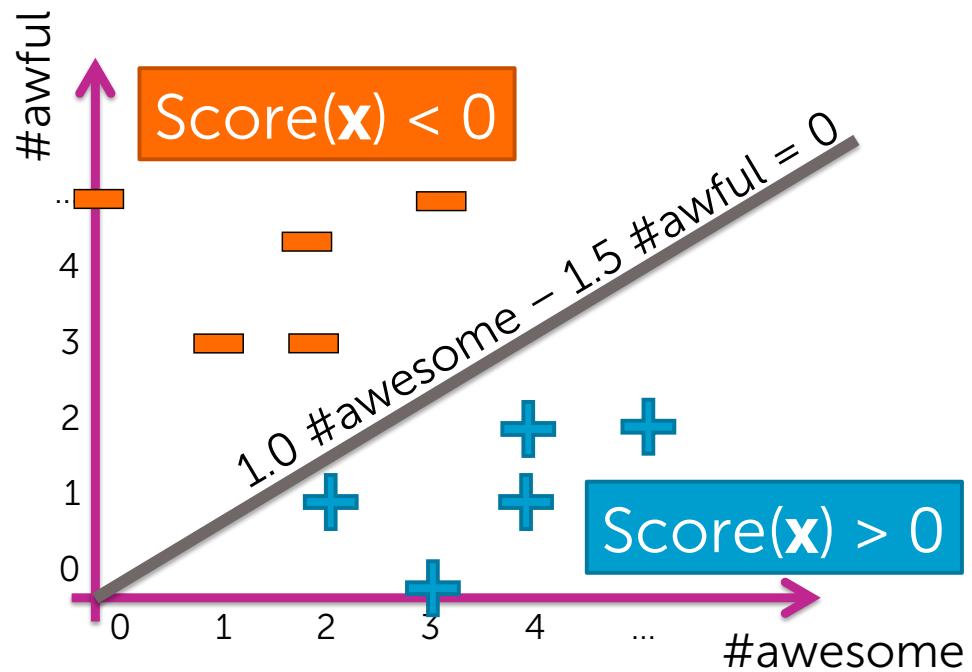
- Estimating $\hat{P}(y|x)$ improves **interpretability**:
 - Predict $\hat{y} = +1$ **and** tell me how sure you are

Predicting class probabilities with generalized linear models



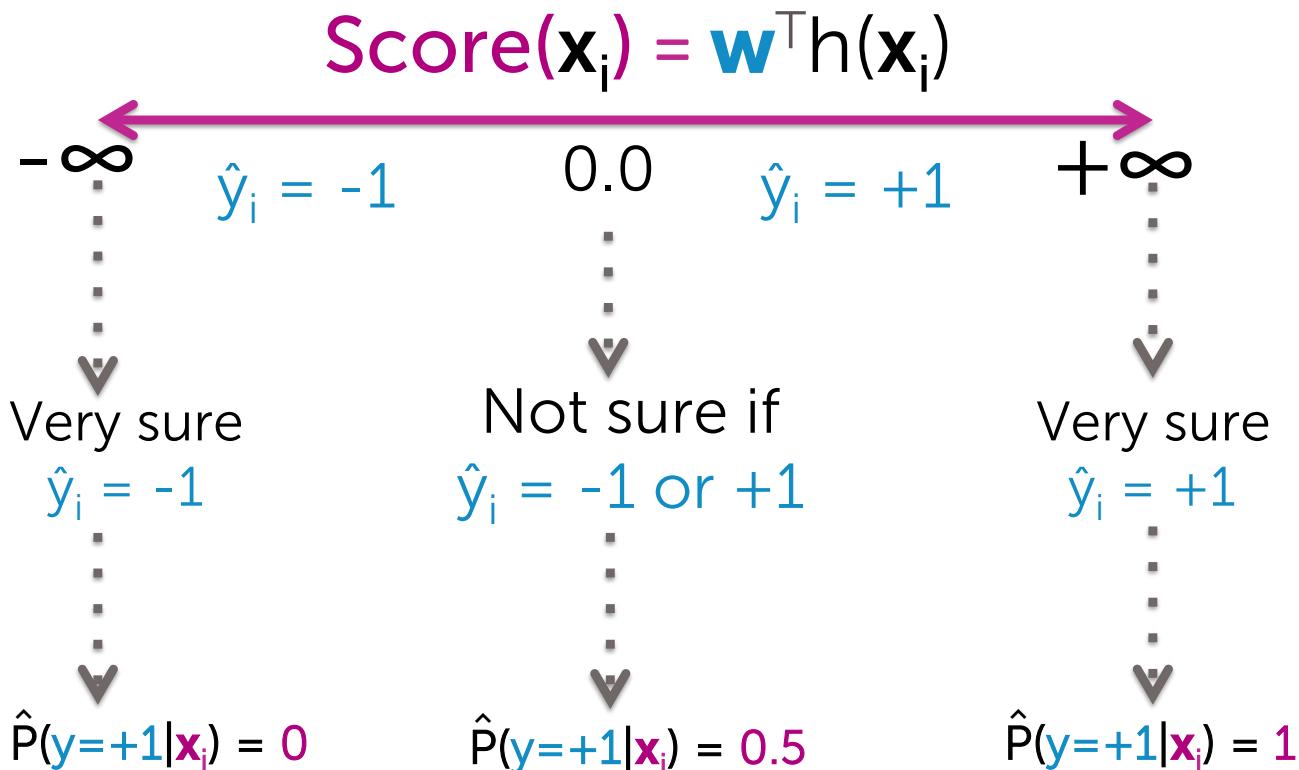
Thus far, we focused on decision boundaries

$$\begin{aligned}\text{Score}(\mathbf{x}_i) &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) \\ &= \mathbf{w}^\top \mathbf{h}(\mathbf{x}_i)\end{aligned}$$

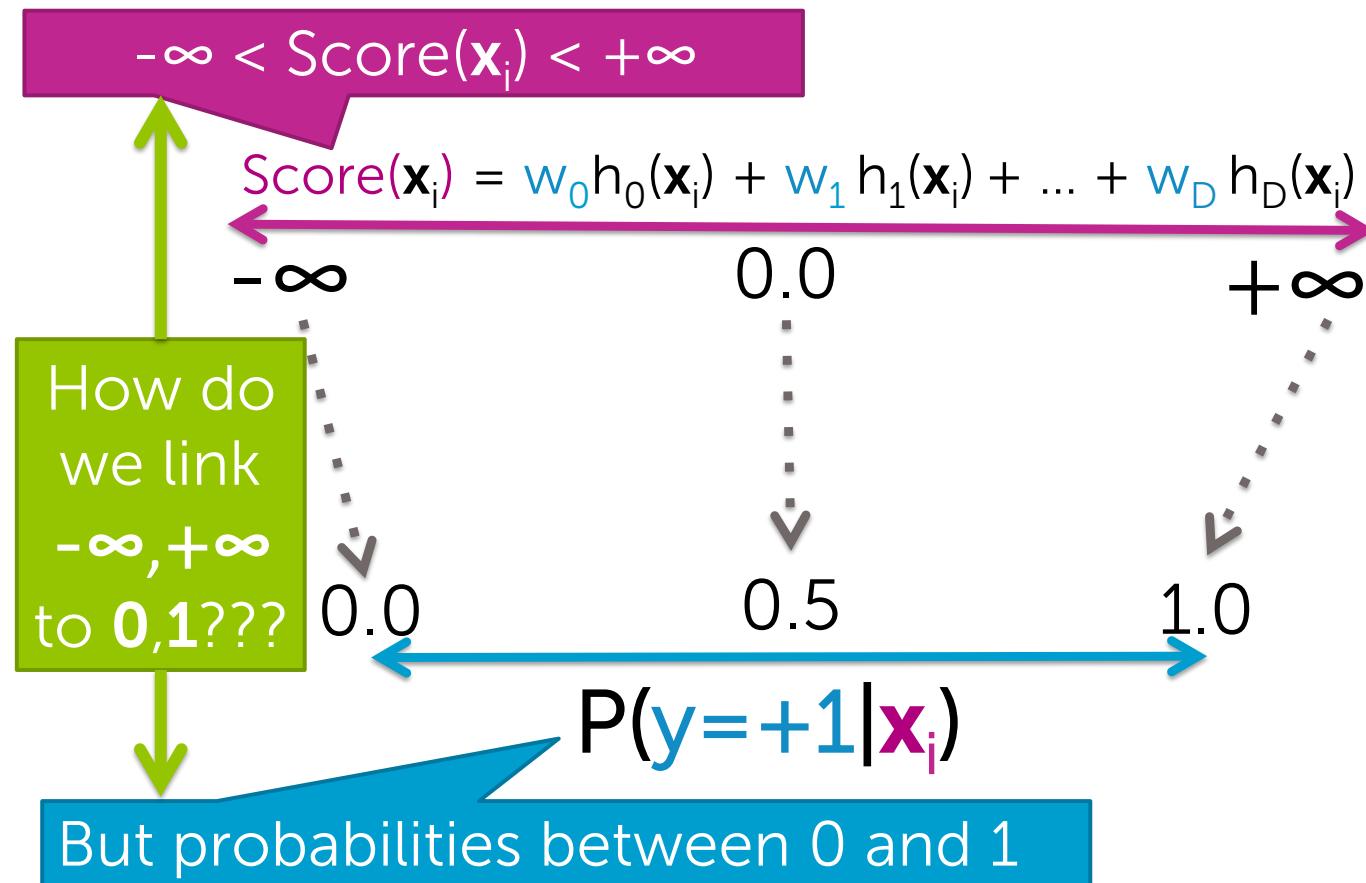


Relate
 $\text{Score}(\mathbf{x}_i)$ to
 $\hat{P}(y=+1|\mathbf{x}, \hat{\mathbf{w}})$?

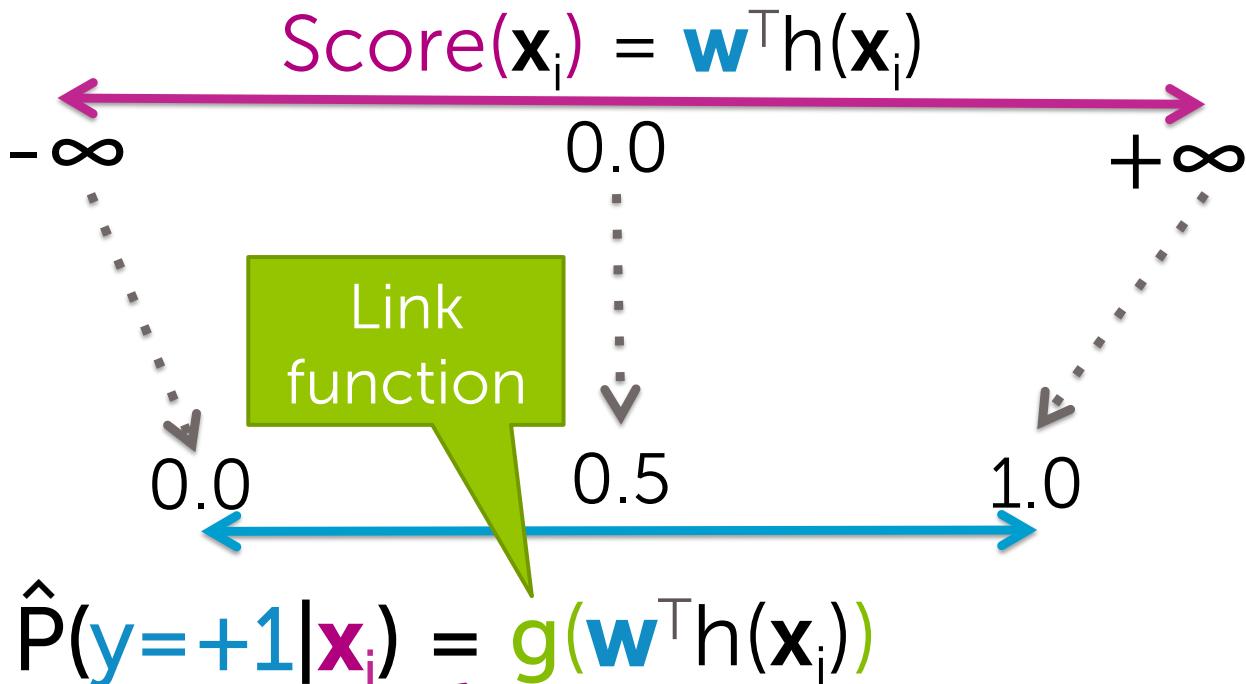
Interpreting $\text{Score}(\mathbf{x}_i)$



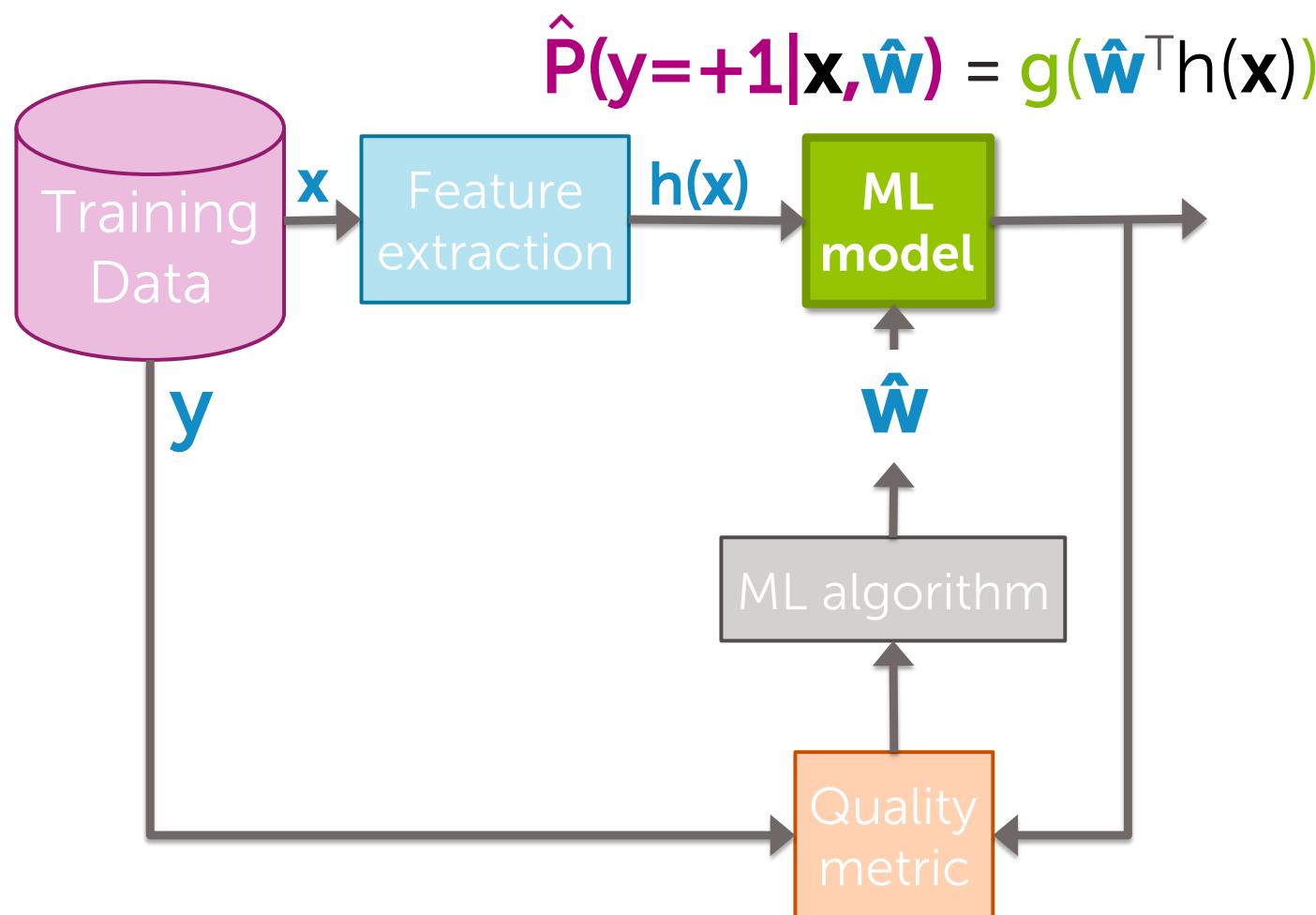
Why not just use regression to build classifier?



Link function: squeeze real line into [0,1]



Generalized linear model

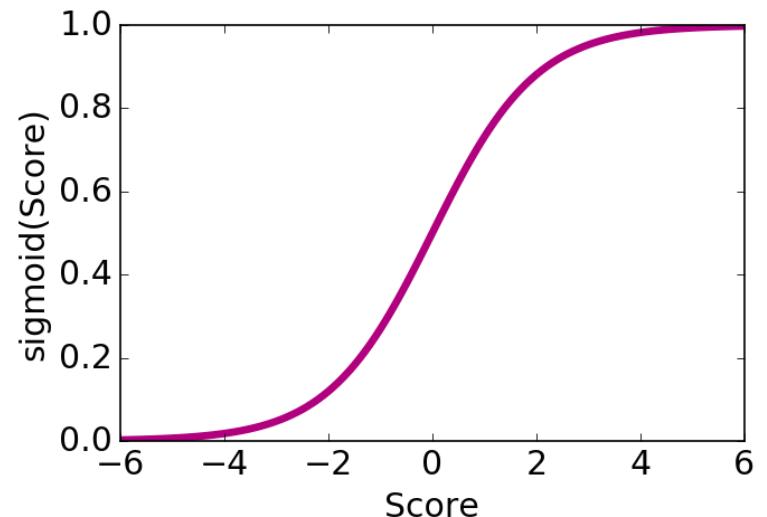


Logistic regression classifier:
linear score with
logistic link function

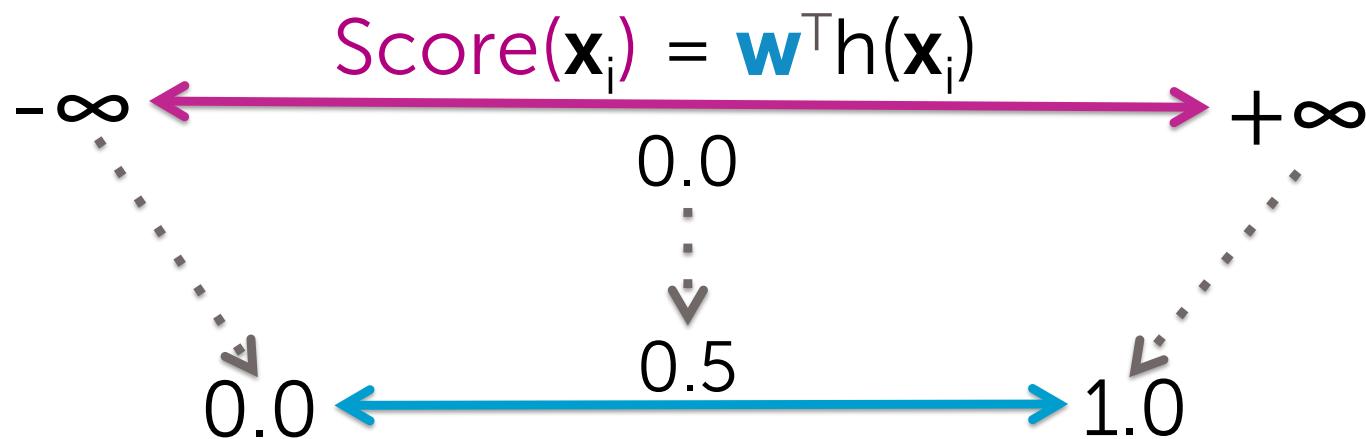
Logistic function (sigmoid, logit)

$$\text{sigmoid}(\text{Score}) = \frac{1}{1 + e^{-\text{Score}}}$$

Score	$-\infty$	-2	0.0	+2	$+\infty$
sigmoid(Score)	0.0	~0.13	0.5	~0.87	1.0



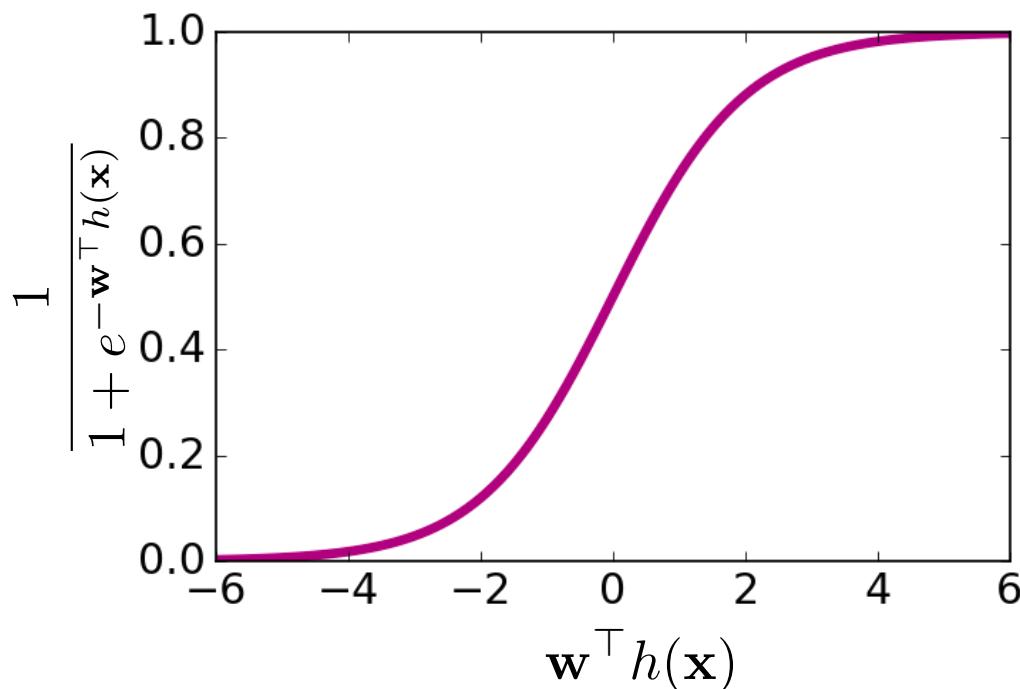
Logistic regression model



$$P(y=+1|x_i, \mathbf{w}) = \text{sigmoid}(\text{Score}(\mathbf{x}_i))$$

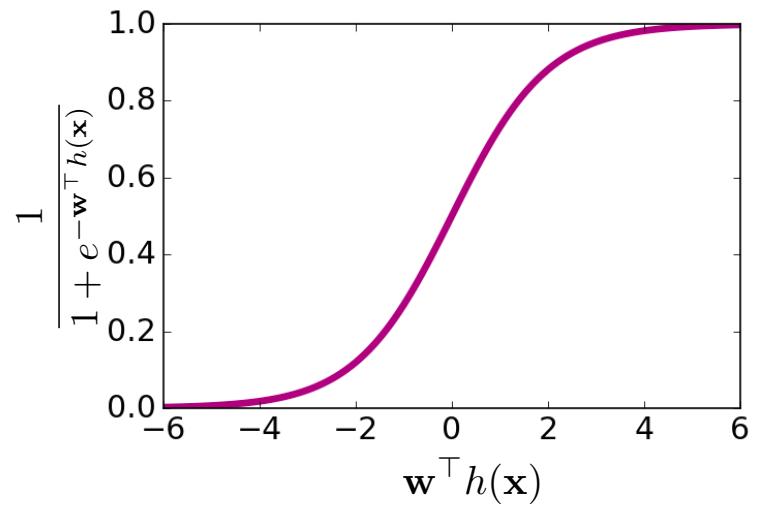
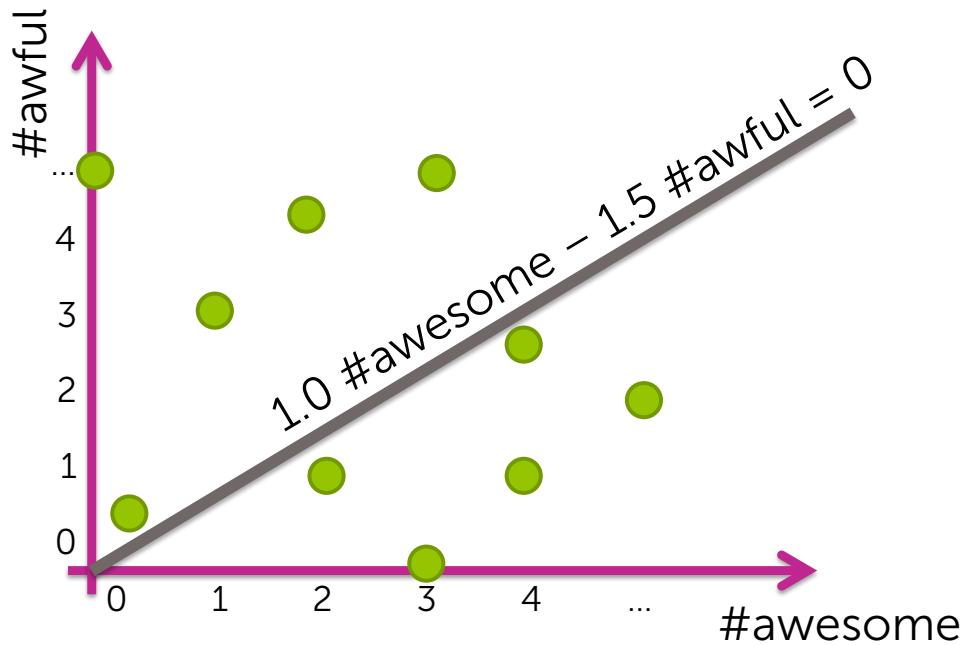
Understanding the logistic regression model

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



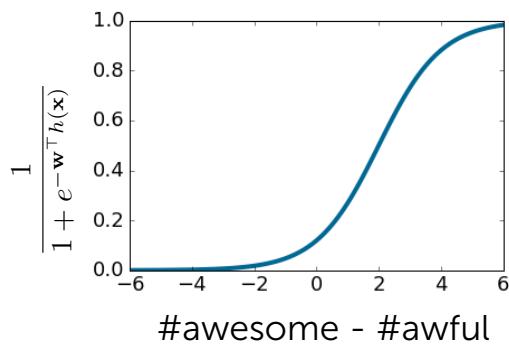
Score(\mathbf{x}_i)	$P(y=+1 \mathbf{x}_i, \mathbf{w})$
0.5	0.5
1.0	0.73
1.5	0.84
2.0	0.91
2.5	0.95
3.0	0.97
3.5	0.98
4.0	0.99
4.5	0.995
5.0	0.998
5.5	0.999
6.0	0.9995
6.5	0.9998
7.0	0.9999
7.5	0.99995
8.0	0.99998
8.5	0.99999
9.0	0.999995
9.5	0.999998
10.0	0.999999

Logistic regression → Linear decision boundary

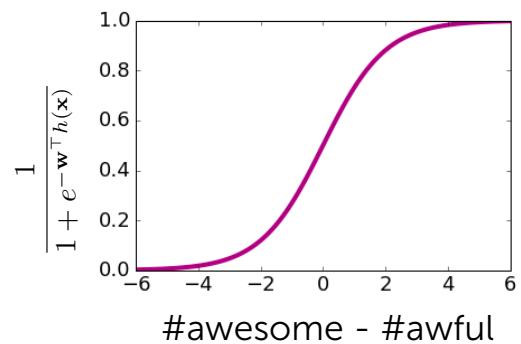


Effect of coefficients on logistic regression model

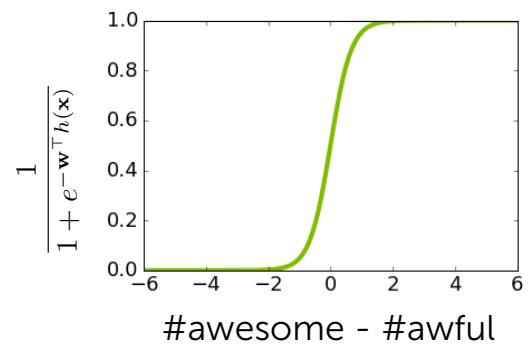
w_0	-2
$w_{\text{#awesome}}$	+1
$w_{\text{#awful}}$	-1



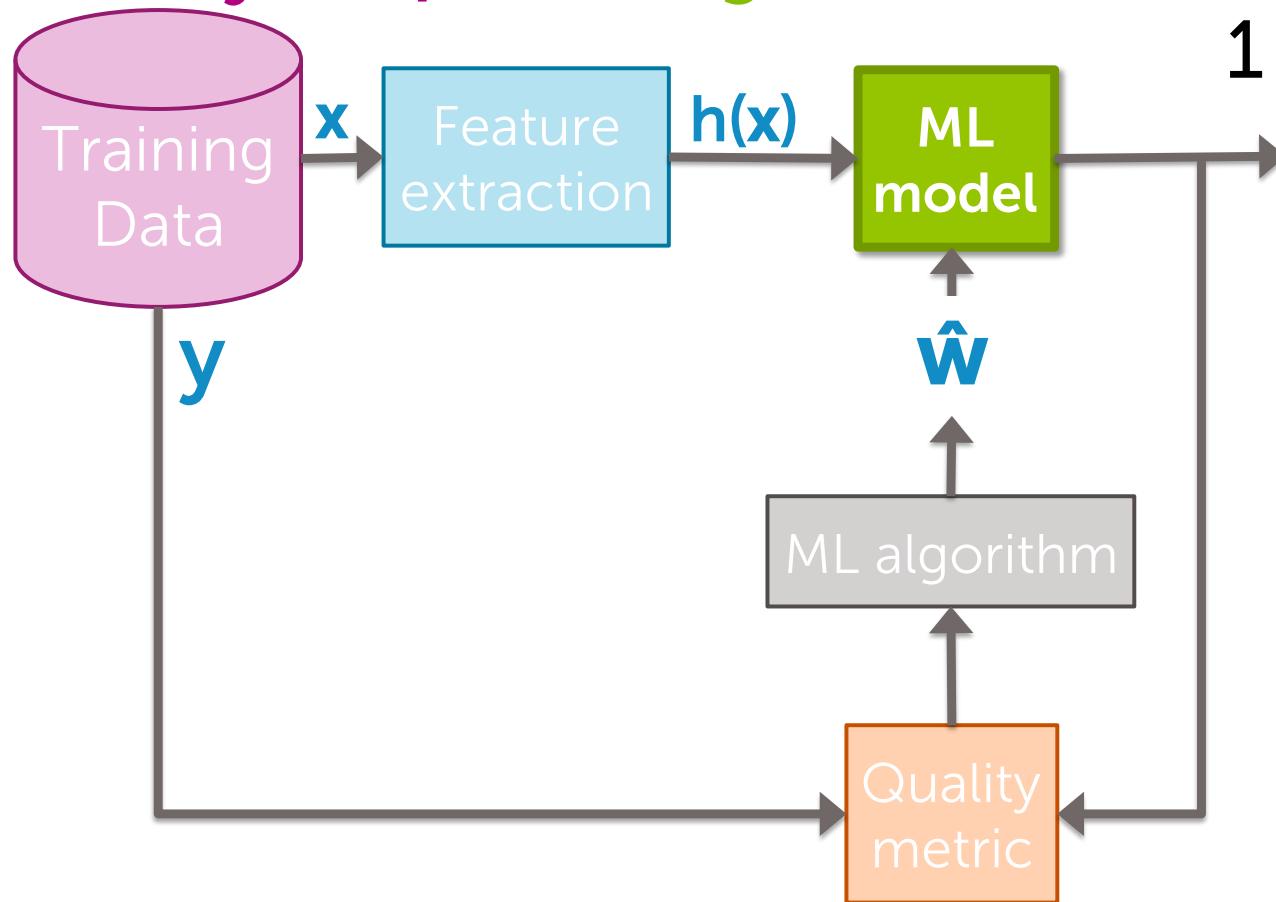
w_0	0
$w_{\text{#awesome}}$	+1
$w_{\text{#awful}}$	-1



w_0	0
$w_{\text{#awesome}}$	+3
$w_{\text{#awful}}$	-3



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

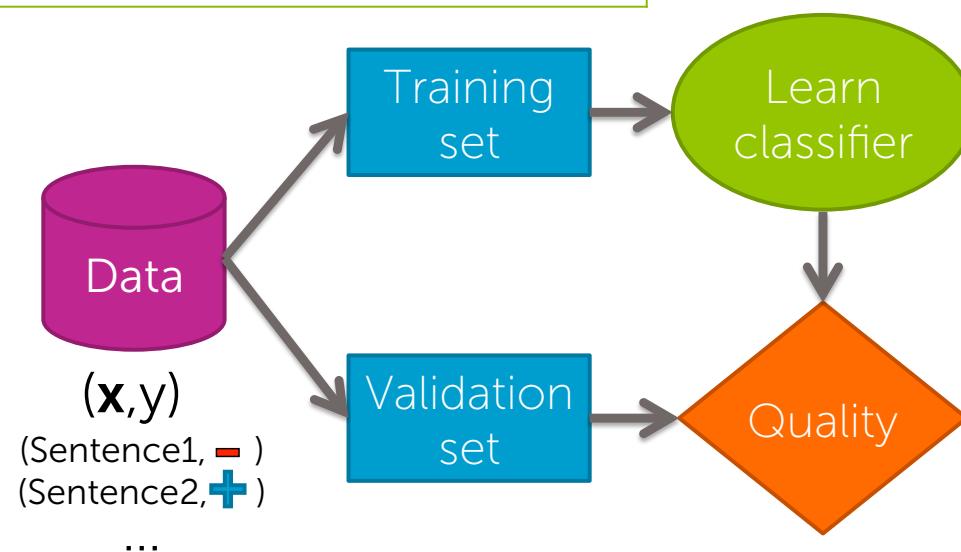


Overview of learning logistic regression model

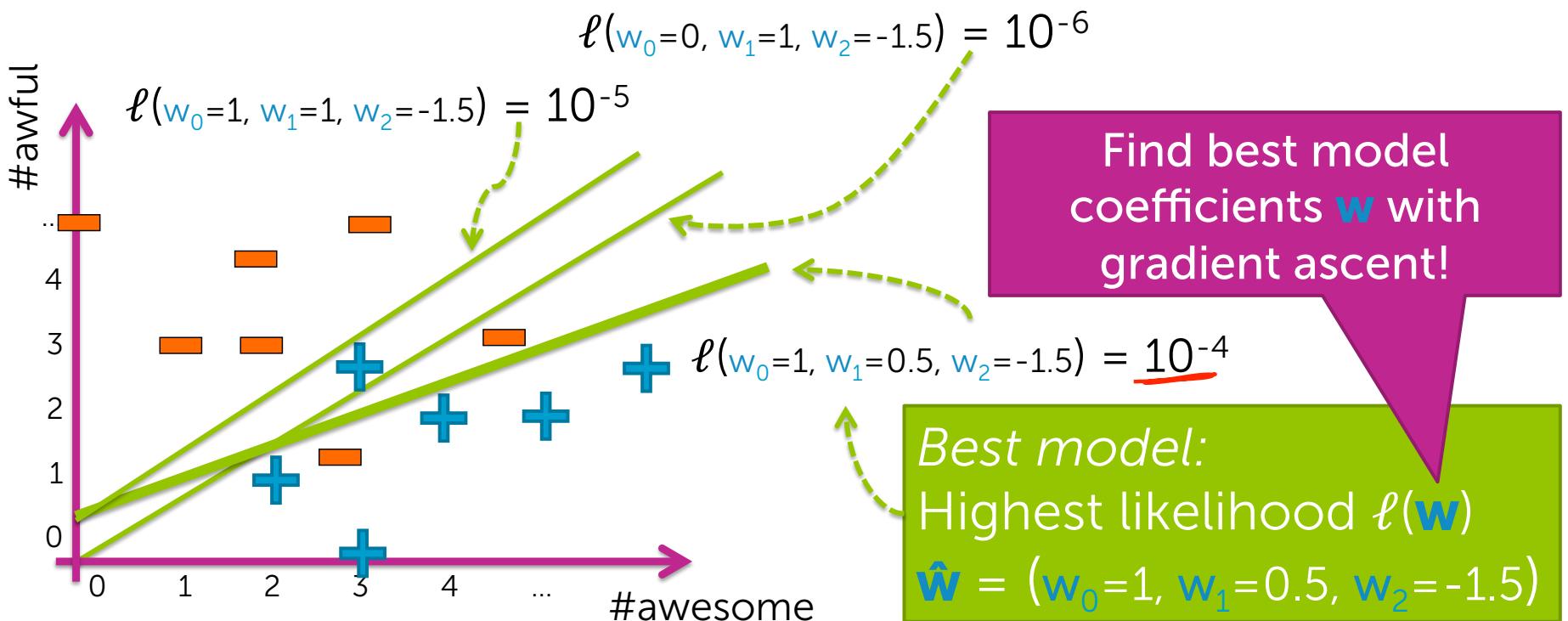
Training a classifier = Learning the coefficients

Word	Coefficient	Value
	\hat{w}_0	-2.0
good	\hat{w}_1	1.0
awesome	\hat{w}_2	1.7
bad	\hat{w}_3	-1.0
awful	\hat{w}_4	-3.3
...

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



Find “best” classifier =
Maximize quality metric over all possible w_0, w_1, w_2
Likelihood $\ell(\mathbf{w})$





Encoding categorical inputs

Categorical inputs

- Numeric inputs:
 - #awesome, age, salary,...
 - Intuitive when multiplied by coefficient
 - e.g., **1.5 #awesome**

Numeric value, but should be interpreted as category
(98195 not about 9x larger than 10005)



Gender
(Male, Female,...)



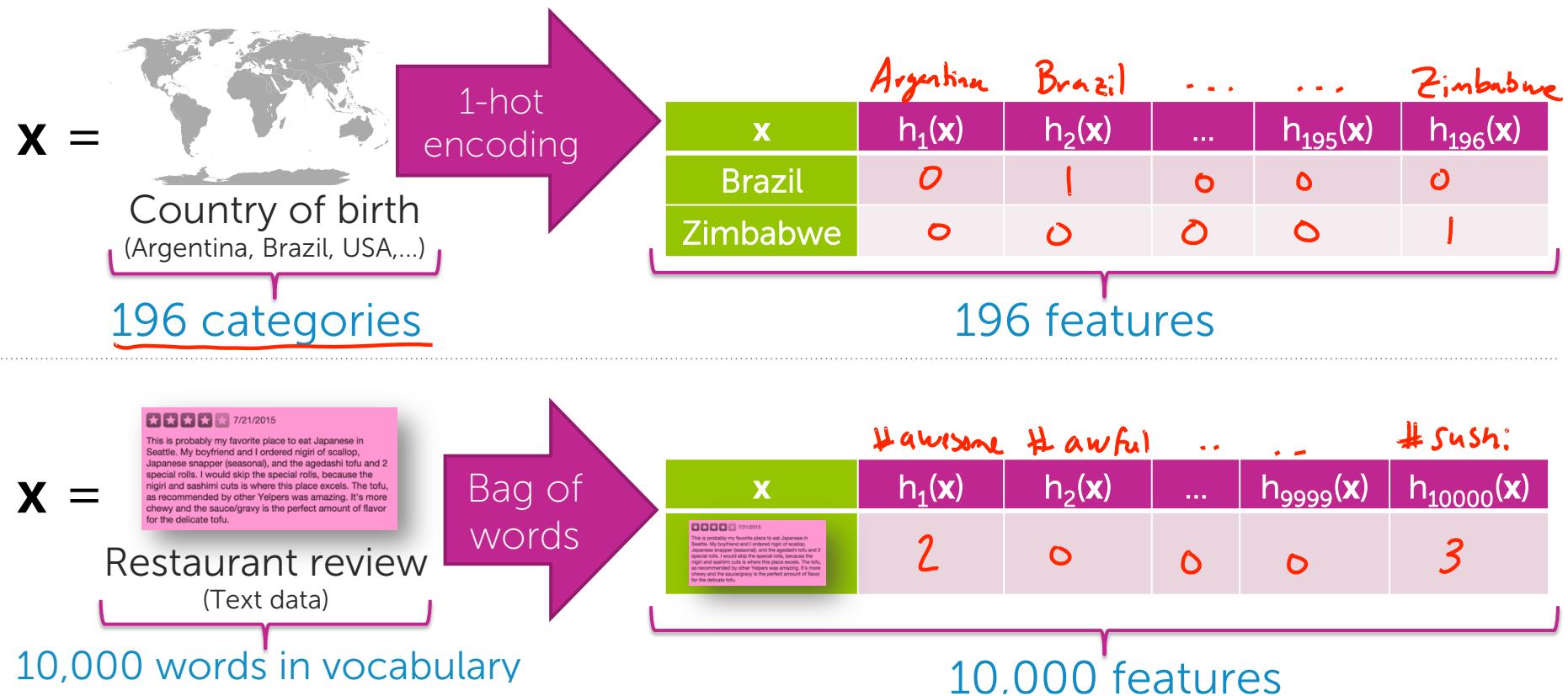
Country of birth
(Argentina, Brazil, USA,...)



Zipcode
(10005, 98195,...)

How do we multiply category by coefficient???
Must convert categorical inputs into numeric features

Encoding categories as numeric features

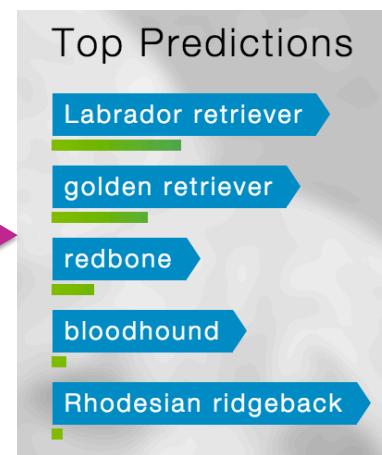




Multiclass classification
using 1 versus all



Multiclass classification



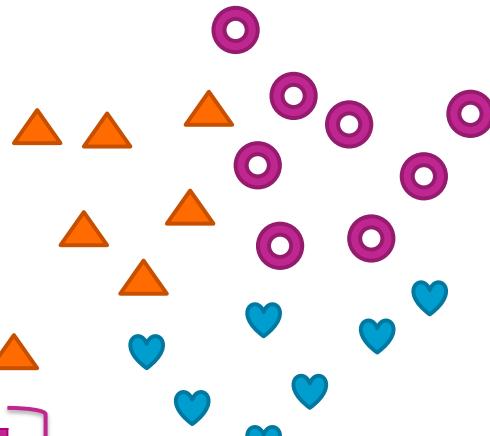
Input: x
Image pixels

Output: y
Object in image

Multiclass classification formulation

- C possible classes:
 - y can be $1, 2, \dots, C$
- N datapoints:

Data point	$x[1]$	$x[2]$	y
\mathbf{x}_1, y_1	2	1	▲
\mathbf{x}_2, y_2	0	2	♥
\mathbf{x}_3, y_3	3	3	○
\mathbf{x}_4, y_4	4	1	○



Learn:

$$\hat{P}(y=\blacktriangle | \mathbf{x})$$

$$\hat{P}(y=\heartsuit | \mathbf{x})$$

$$\hat{P}(y=\circlearrowleft | \mathbf{x})$$

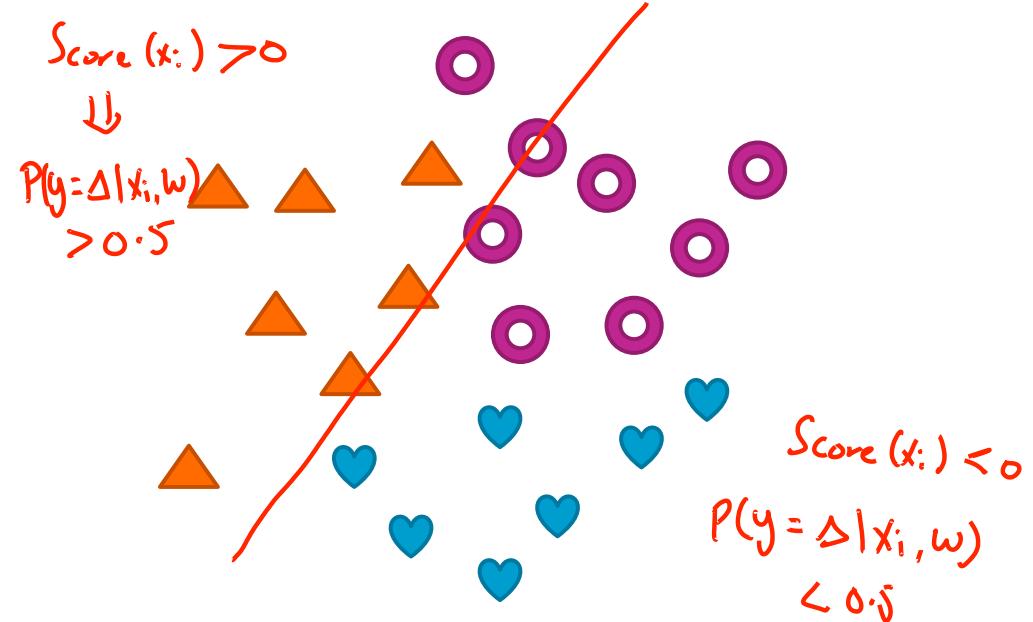
1 versus all:

Estimate $\hat{P}(y=\Delta|x)$ using 2-class model

+1 class: points with $y_i = \Delta$
-1 class: points with $y_i = \heartsuit$ OR \circ

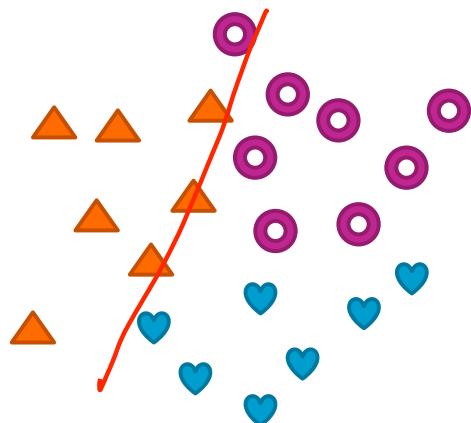
Train classifier: $\hat{P}(\Delta|y=+1|x)$

Predict: $\hat{P}(y=\Delta|x_i) = \hat{P}(\Delta|y=+1|x_i)$

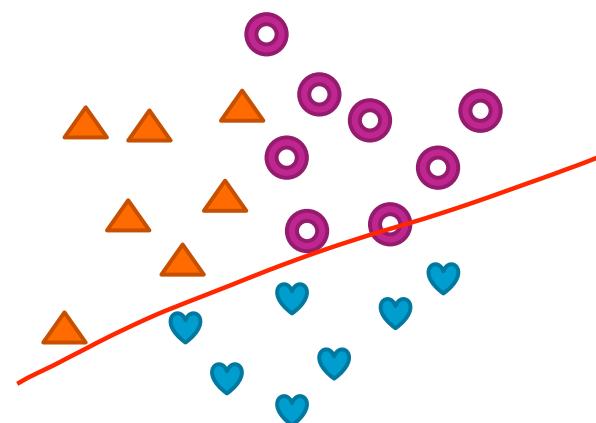


1 versus all: simple multiclass classification using C 2-class models

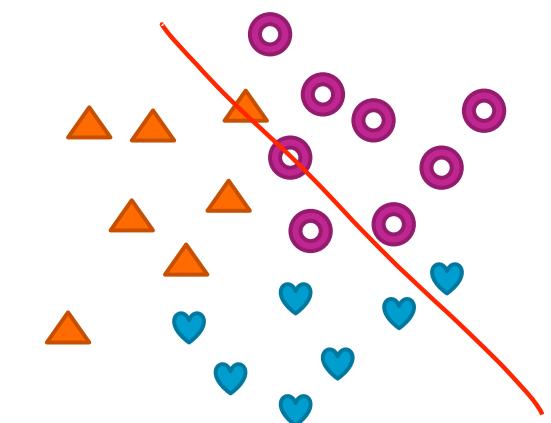
$$\hat{P}(y=\triangle | \mathbf{x}_i) = \hat{P}_{\alpha}(y=+1 | \mathbf{x}_i, w_{\alpha})$$



$$\hat{P}(y=\heartsuit | \mathbf{x}_i) = \hat{P}_{\beta}(y=+1 | \mathbf{x}_i, w_{\beta})$$



$$\hat{P}(y=\circlearrowleft | \mathbf{x}_i) = \hat{P}_{\theta}(y=+1 | \mathbf{x}_i, w_{\theta})$$



Multiclass training

$\hat{P}_c(y=+1|x)$ = estimate of
1 vs all model for each class



Input: x_i

→ Predict most likely class

max_prob = 0; $\hat{y} = 0$

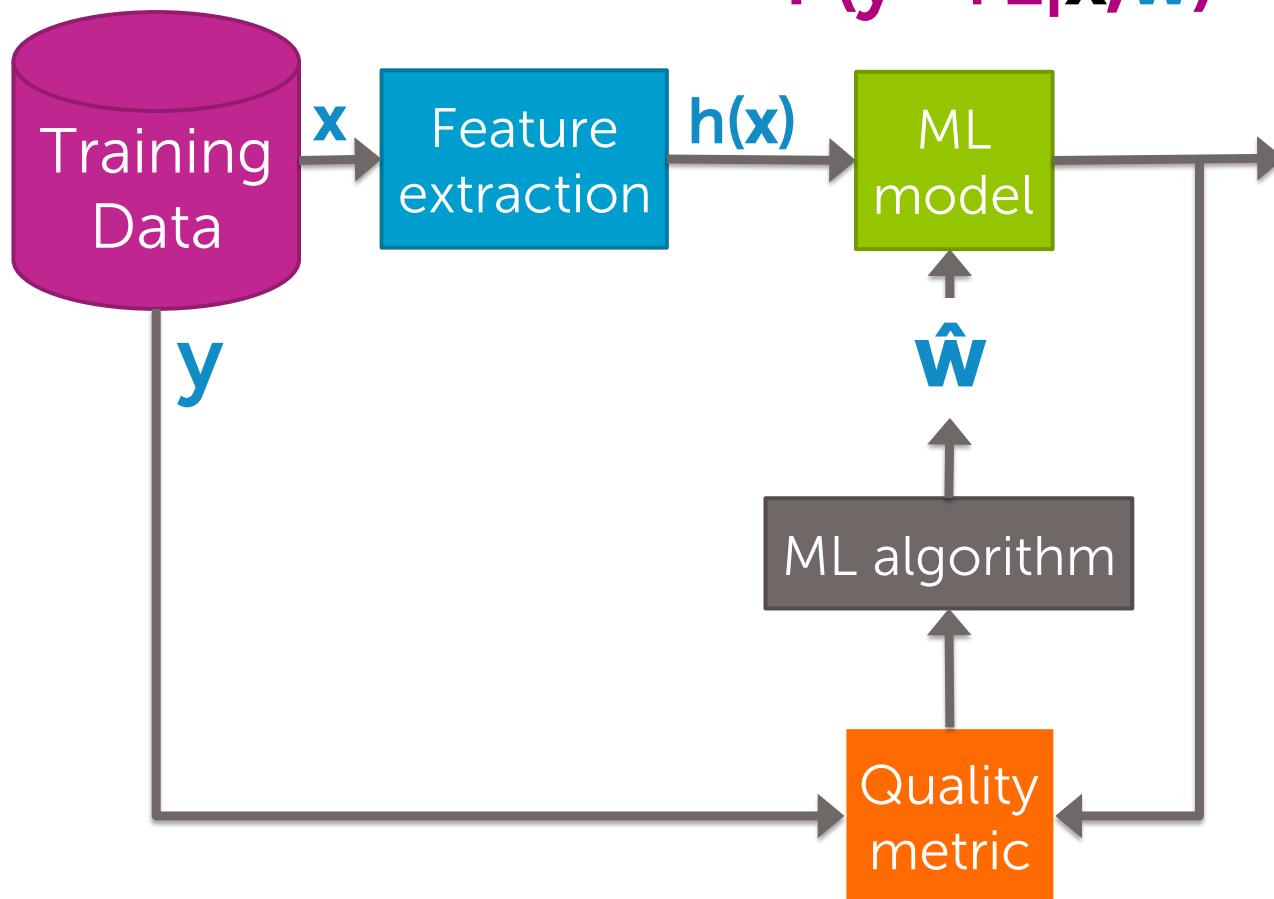
For $c = 1, \dots, C$:

If $\hat{P}_c(y=+1|x_i) >$ max_prob:

$\hat{y} = c$

max_prob = $\hat{P}_c(y=+1|x_i)$

Summary of logistic regression classifier



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



Linear classifiers: Parameter learning



Learn a probabilistic classification model

"The sushi & everything else were awesome!"

Definite +1

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

"The sushi was good, the service was OK"

Not sure

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

Many classifiers provide a degree of certainty:

$$\text{Output label} \quad P(y|x)$$

Input sentence

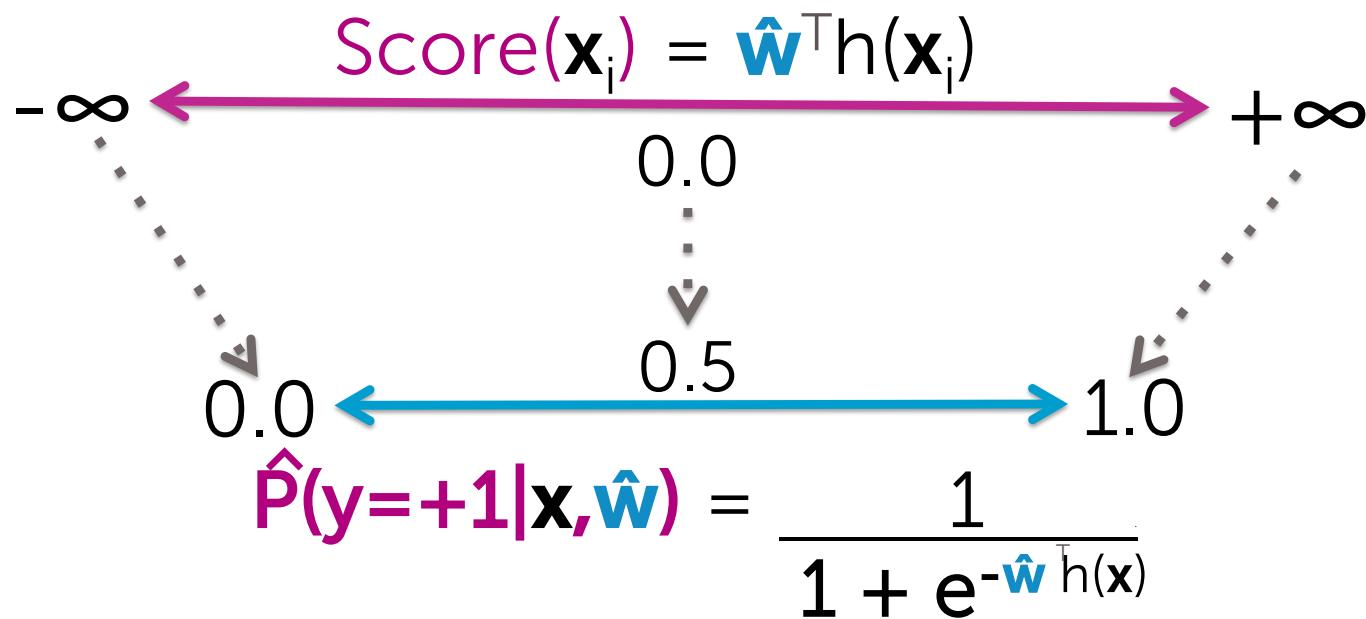
Extremely useful in practice

A (linear) classifier

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

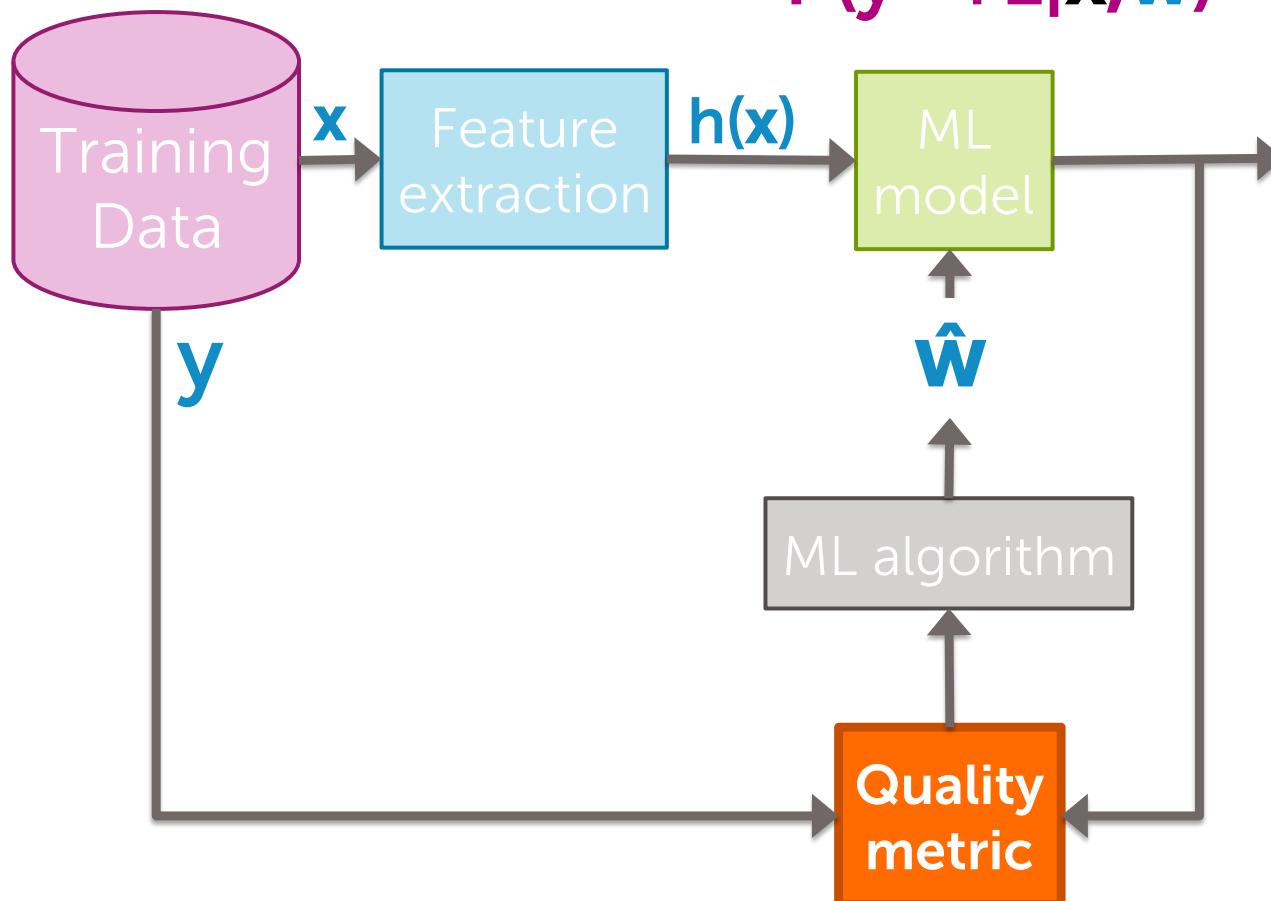
Word	Coefficient	Value
	\hat{w}_0	-2.0
good	\hat{w}_1	1.0
great	\hat{w}_2	1.5
awesome	\hat{w}_3	2.7
bad	\hat{w}_4	-1.0
terrible	\hat{w}_5	-2.1
awful	\hat{w}_6	-3.3
restaurant, the, we, ...	$\hat{w}_7, \hat{w}_8, \hat{w}_9, \dots$	0.0
...		...

Logistic regression model





Quality metric for logistic regression:
Maximum likelihood estimation



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

Learning problem

Training data:

N observations (\mathbf{x}_i, y_i)

$\mathbf{x}[1] = \#awesome$	$\mathbf{x}[2] = \#awful$	$y = \text{sentiment}$
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



$\hat{\mathbf{W}}$

Finding best coefficients

$x[1] = \#awesome$	$x[2] = \#awful$	$y = \text{sentiment}$
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

$x[1] = \#awesome$	$x[2] = \#awful$	$y = \text{sentiment}$
0	2	-1
3	3	-1
2	4	-1
0	3	-1
0	1	-1

$x[1] = \#awesome$	$x[2] = \#awful$	$y = \text{sentiment}$
2	1	+1
4	1	+1
1	1	+1
2	1	+1

Finding best coefficients

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
0	2	-1
3	3	-1
2	4	-1
0	3	-1
0	1	-1

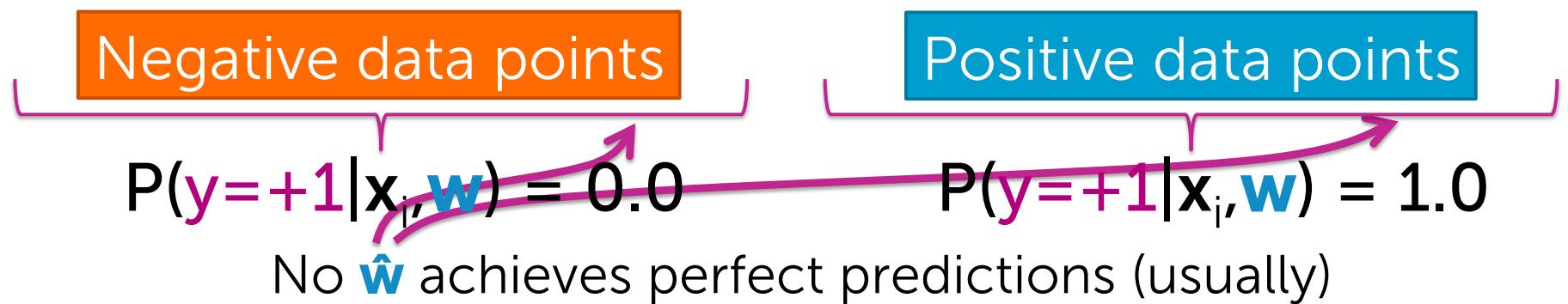
$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
2	1	+1
4	1	+1
1	1	+1
2	1	+1

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

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

Pick \hat{w} that makes

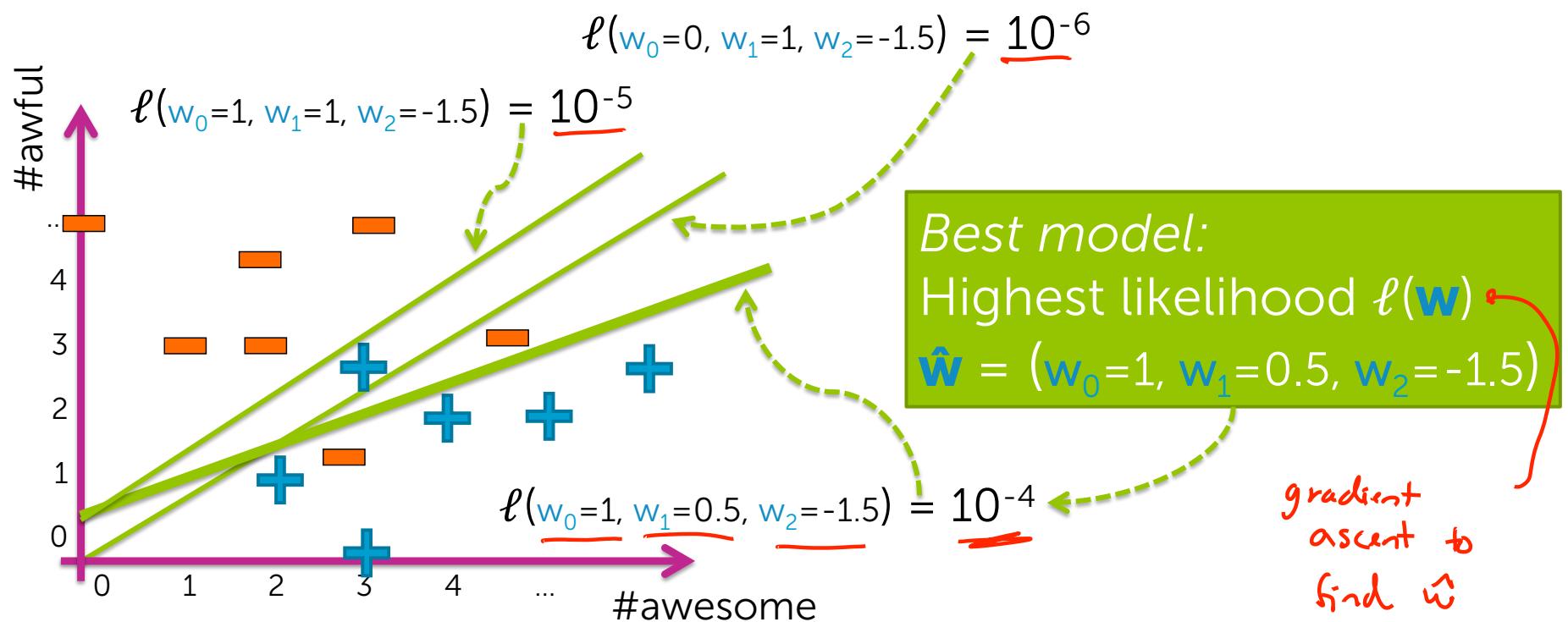
Quality metric = Likelihood function



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

Find “best” classifier

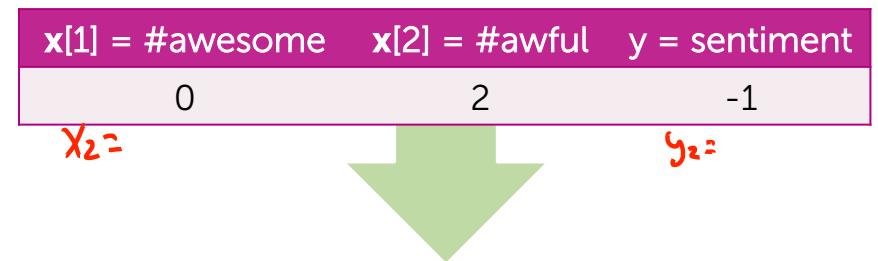
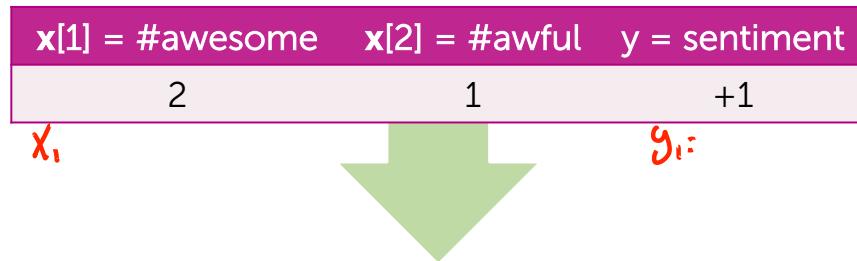
Maximize likelihood over all possible w_0, w_1, w_2





Data likelihood

Quality metric: probability of data



If model good, should predict:
 $\hat{y}_1 = +1$

If model good, should predict:
 $\hat{y}_2 = -1$

Pick w to maximize:
 $P(y=+1|x_1, w) = P(y=+1|x_1=2, x_2=1, w)$

Pick w to maximize:
 $P(y=-1|x_2, w)$

Maximizing likelihood (probability of data)

Data point	x[1]	x[2]	y	Choose w to maximize
\mathbf{x}_1, y_1	2	1	+1	$P(y=+1 \mathbf{x}_1, w) = P(y=+1 x_1=2, x_2=1, w)$
\mathbf{x}_2, y_2	0	2	-1	$P(y=-1 \mathbf{x}_2, w)$
\mathbf{x}_3, y_3	3	3	<u>-1</u>	$P(y=-1 \mathbf{x}_3, w)$
\mathbf{x}_4, y_4	4	1	<u>+1</u>	$P(y=+1 \mathbf{x}_4, w)$
\mathbf{x}_5, y_5	1	1	+1	
\mathbf{x}_6, y_6	2	4	-1	
\mathbf{x}_7, y_7	0	3	-1	
\mathbf{x}_8, y_8	0	1	-1	
\mathbf{x}_9, y_9	2	1	+1	

Must combine into single measure of quality ?

Multiply probabilities

$P(y=+1|\mathbf{x}_1, w) P(y=-1|\mathbf{x}_2, w) P(y=-1|\mathbf{x}_3, w) \dots$

Learn logistic regression model with maximum likelihood estimation (MLE)

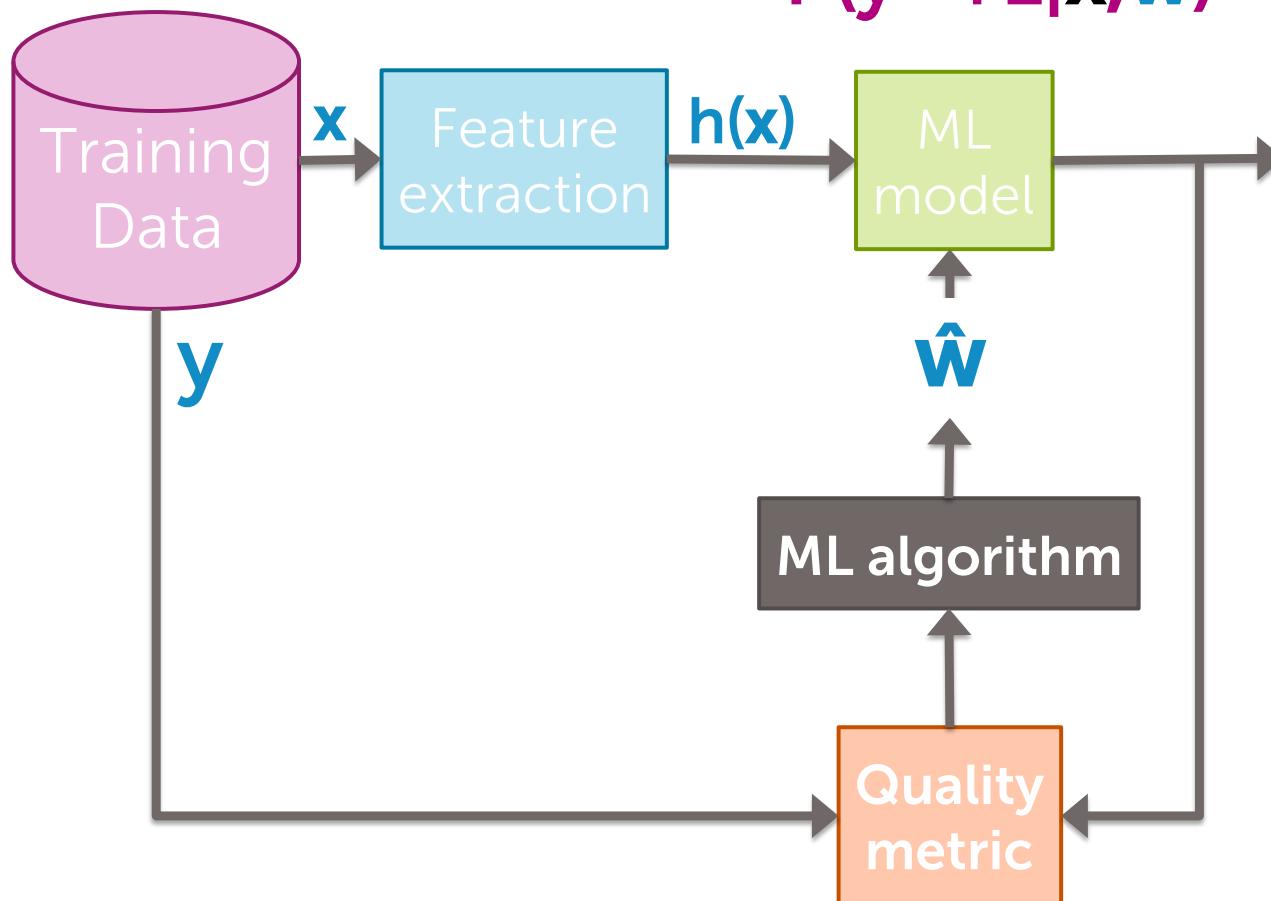
Data point	$\mathbf{x}[1]$	$\mathbf{x}[2]$	y	Choose \mathbf{w} to maximize
\mathbf{x}_1, y_1	2	1	$y_1 = +1$	$P(y_1 = +1 \mathbf{x}_1 = 2, \mathbf{x}_2 = 1, \mathbf{w})$
\mathbf{x}_2, y_2	0	2	-1	$P(y_2 = -1 \mathbf{x}_1 = 0, \mathbf{x}_2 = 2, \mathbf{w})$
\mathbf{x}_3, y_3	3	3	-1	$P(y_3 = -1 \mathbf{x}_1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

$$\ell(\mathbf{w}) = P(y_1 | \mathbf{x}_1, \mathbf{w}) P(y_2 | \mathbf{x}_2, \mathbf{w}) P(y_3 | \mathbf{x}_3, \mathbf{w}) P(y_4 | \mathbf{x}_4, \mathbf{w})$$

Num. of data points → $\ell(\mathbf{w}) = \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$
← pick \mathbf{w} to make this fn. as large as possible.



Finding best linear classifier
with gradient ascent

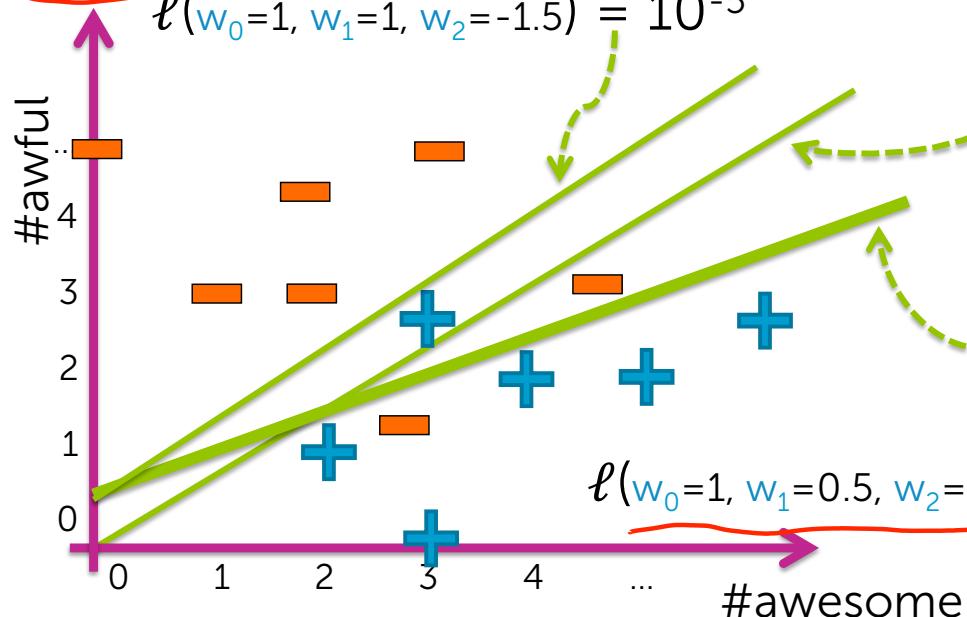


$$\hat{P}(y=+1|x, \hat{w}) = \frac{1}{1 + e^{-\hat{w}^T h(x)}}$$

Find “best” classifier

Maximize likelihood over all possible w_0, w_1, w_2

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



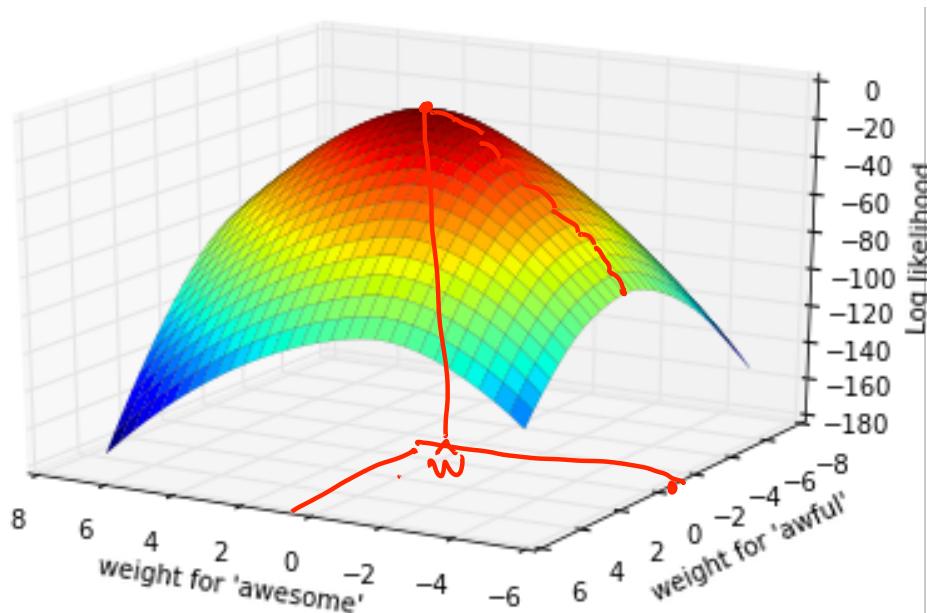
$$\ell(w_0=0, w_1=1, w_2=-1.5) = 10^{-6}$$

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

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

optimize with
gradient ascent

Maximizing likelihood



No closed-form solution → use gradient ascent

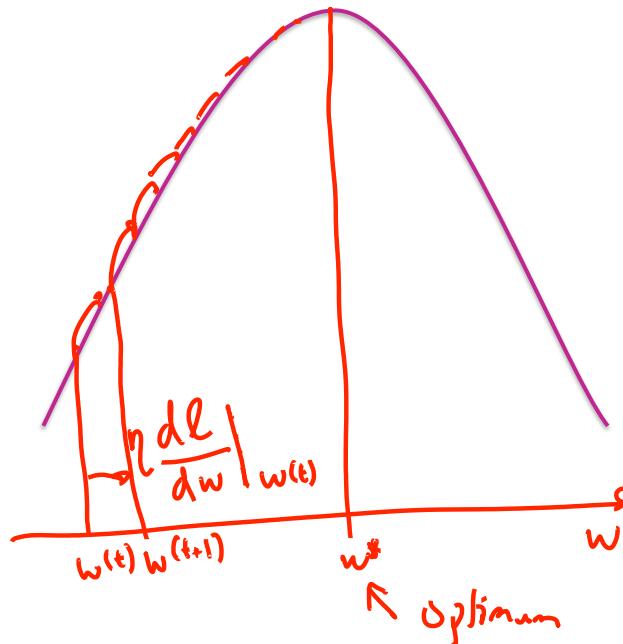
Maximize function over all possible w_0, w_1, w_2

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

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

Review of gradient ascent

Finding the max via hill climbing



Algorithm:

while not converged

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

Convergence criteria

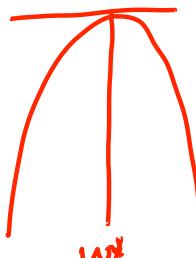
For convex functions,
optimum occurs when

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

In practice, stop when

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

\uparrow
tolerance

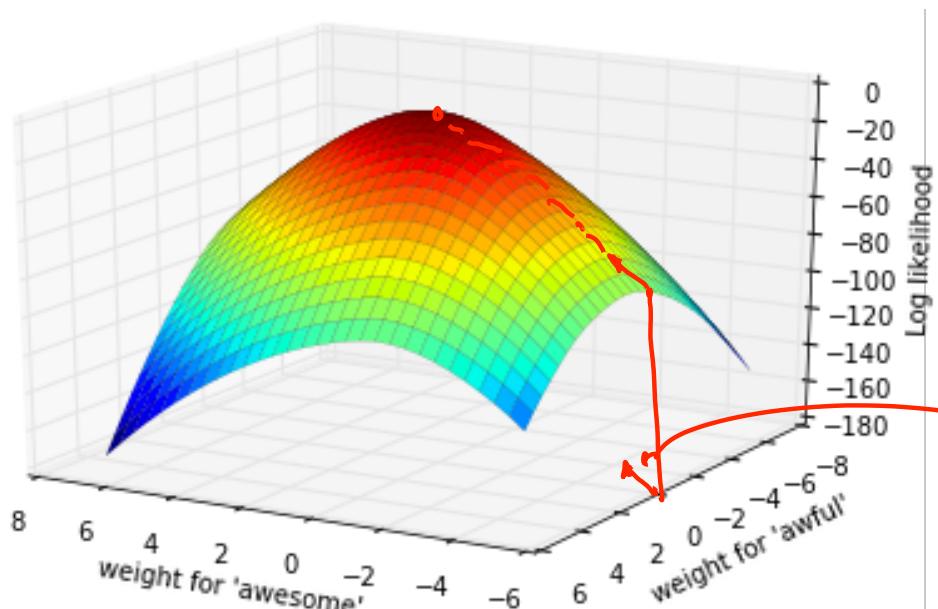


Algorithm:

while not converged

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

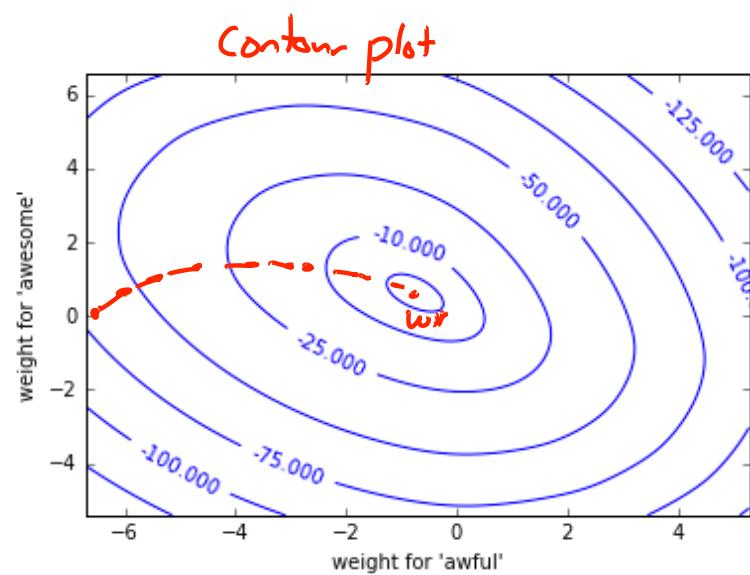
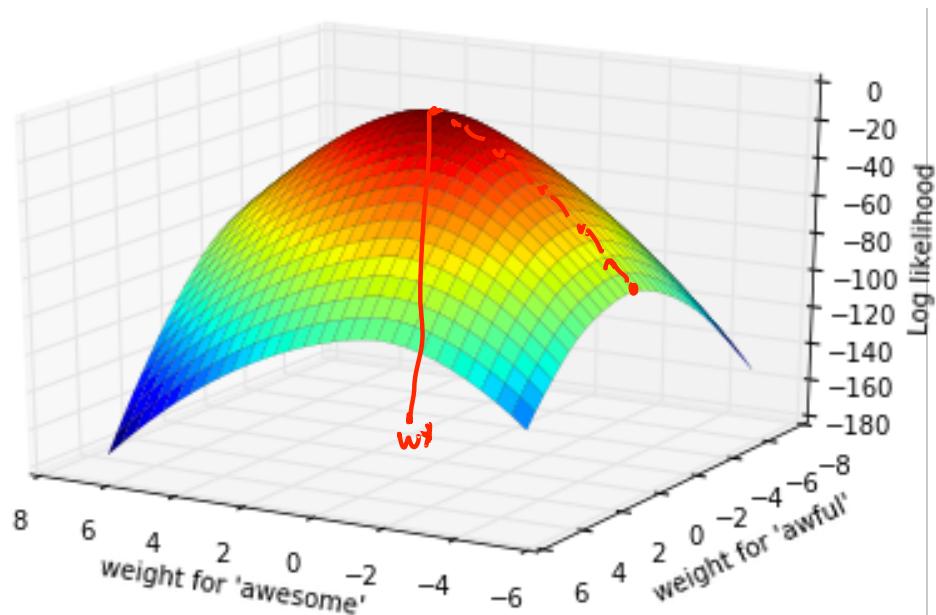
Moving to multiple dimensions: Gradients



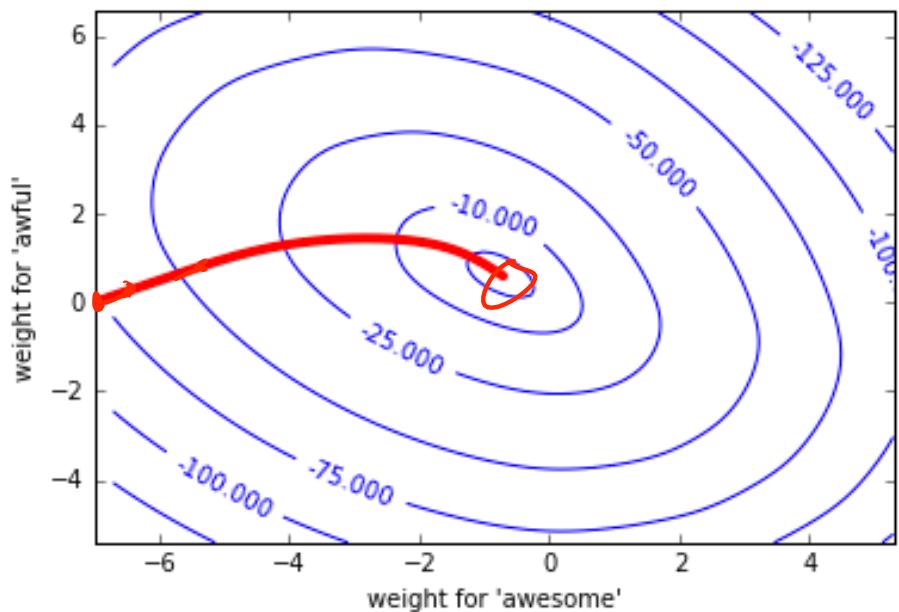
$$\nabla \ell(\mathbf{w}) = \left[\begin{array}{c} \frac{\partial \ell}{\partial w_0} \\ \frac{\partial \ell}{\partial w_1} \\ \vdots \\ \frac{\partial \ell}{\partial w_D} \end{array} \right]$$

D+1 dim vector

Contour plots



Gradient ascent



Algorithm:

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

while not converged

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



Learning algorithm for logistic regression

Derivative of (log-)likelihood

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

Sum over data points

Feature value

Difference between truth and prediction

predict \mathbf{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

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

$w^{(t)}$:

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

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

$h_1(k) = 1$ awesome

$x[1]$	$x[2]$	y	$P(y=+1 \mathbf{x}_i, \mathbf{w})$	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:

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

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

$$= 1 + 0.1 \times 1.33 = 1.133$$

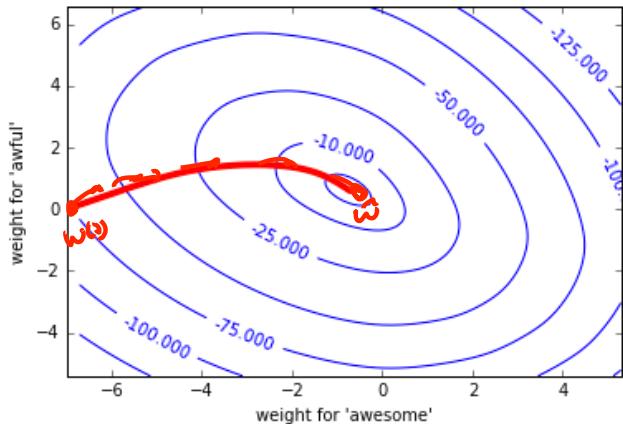
Derivative of (log-)likelihood: Interpretation

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

Difference between truth and prediction

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

Summary of gradient ascent for logistic regression



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

while $\|\nabla \ell(\mathbf{w}^{(t)})\| > \epsilon$
for $j=0, \dots, D$

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

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

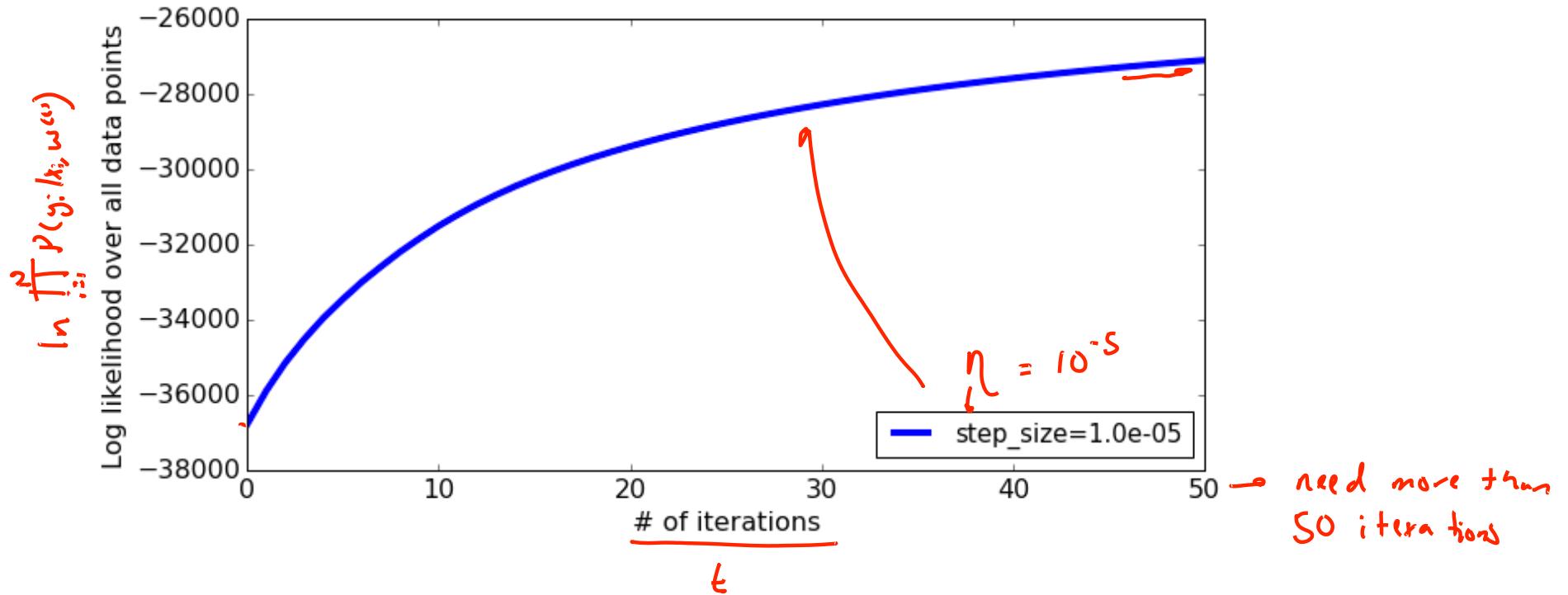
$$t \leftarrow t + 1$$

$$\frac{\partial \ell(\mathbf{w}^{(t)})}{\partial \mathbf{w}_j}$$

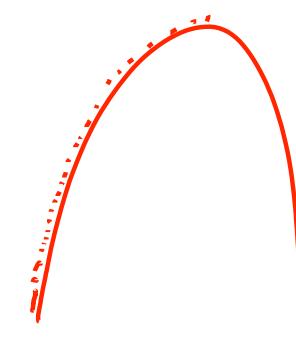
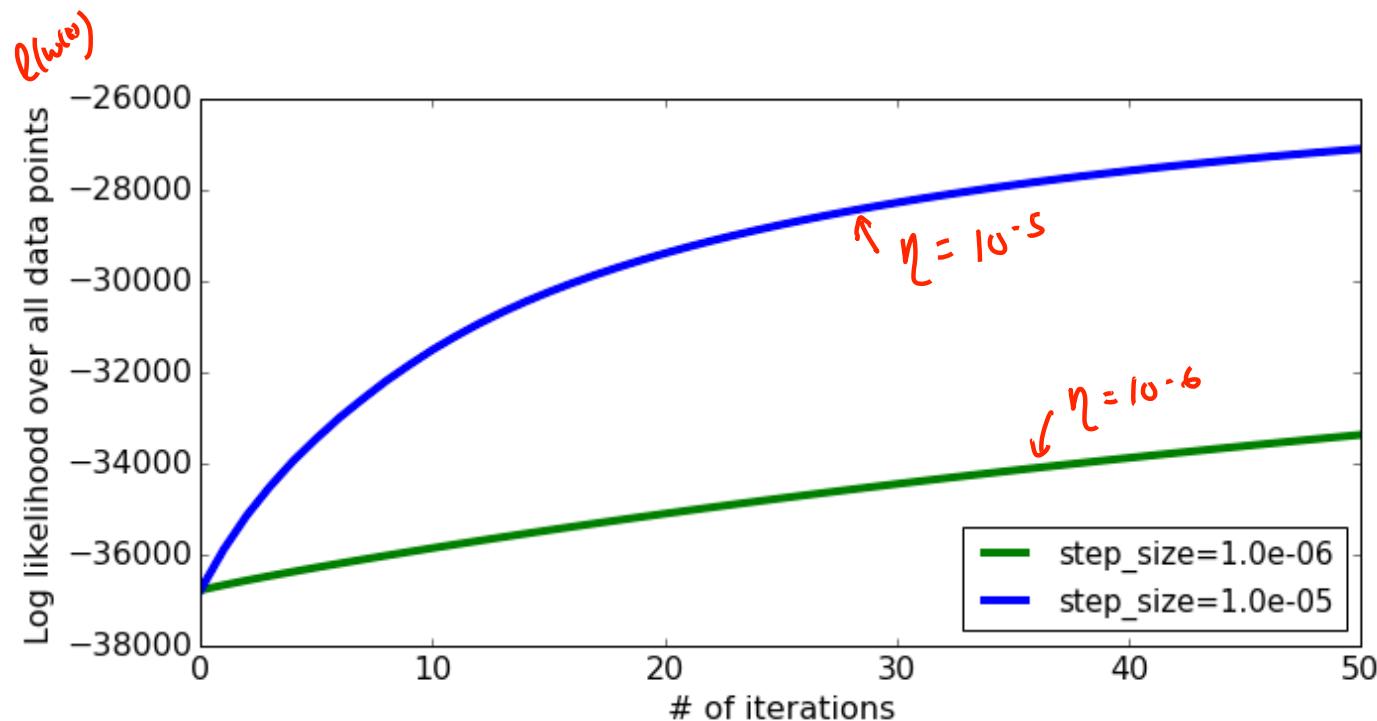
Step size

Choosing the step size η

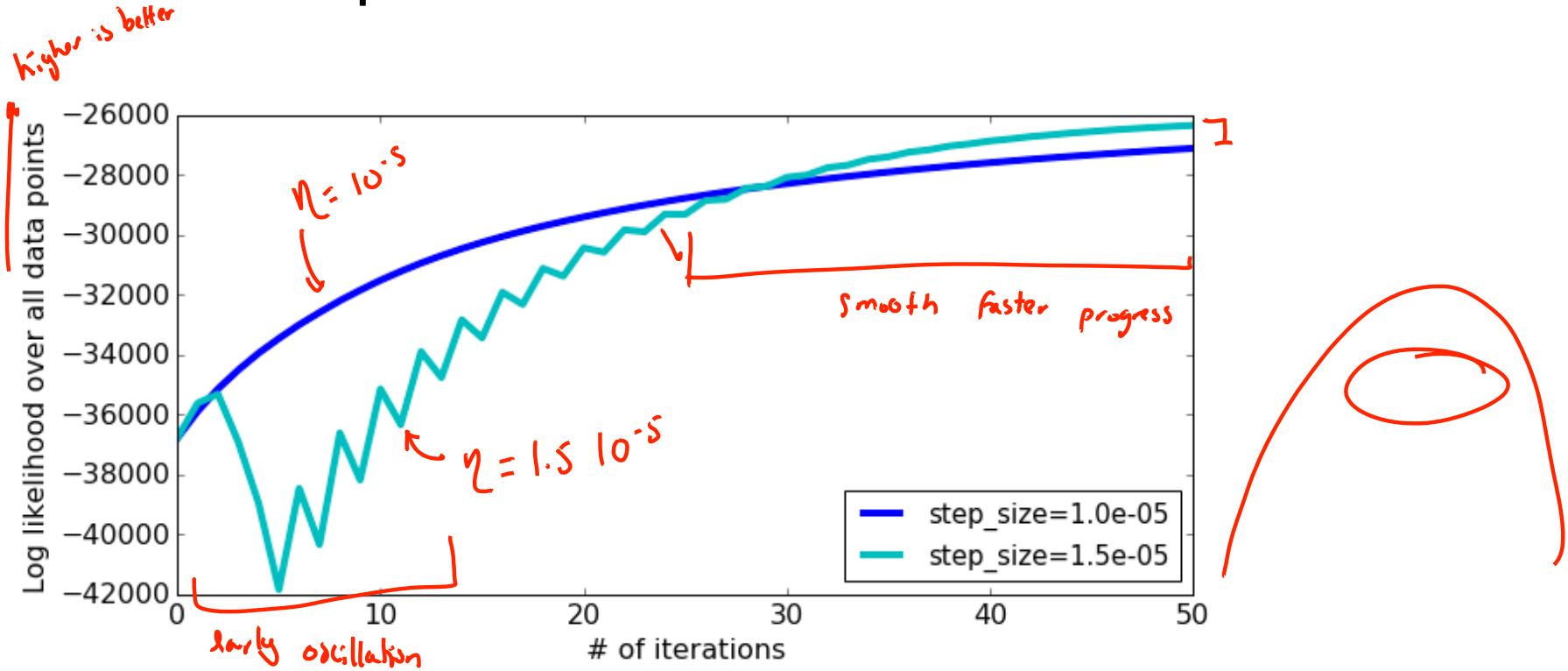
Learning curve: Plot quality (likelihood) over iterations



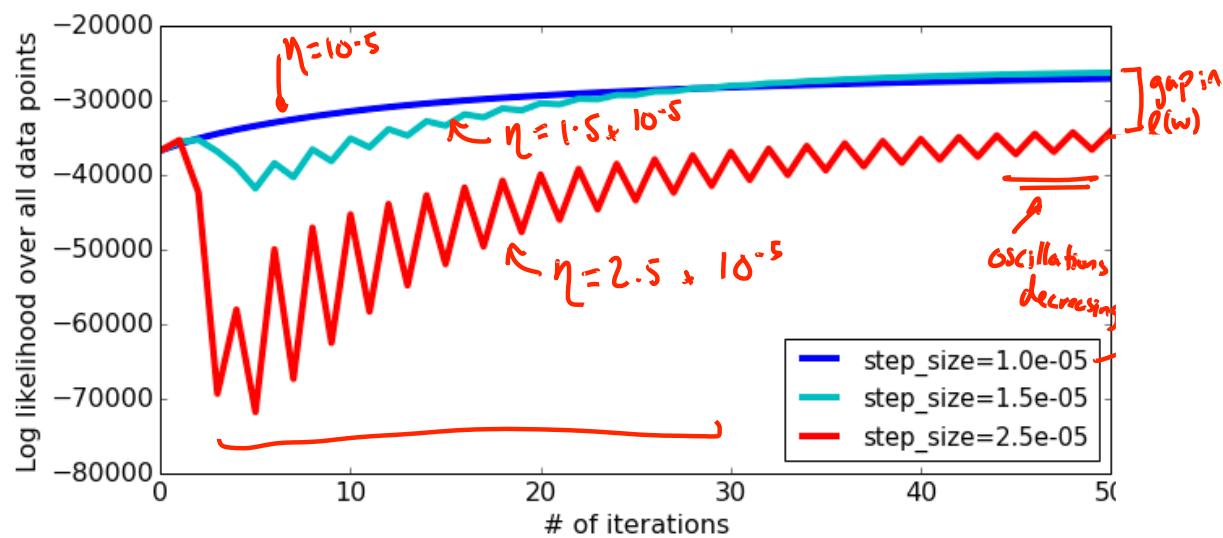
If step size is too small, can take a long time to converge



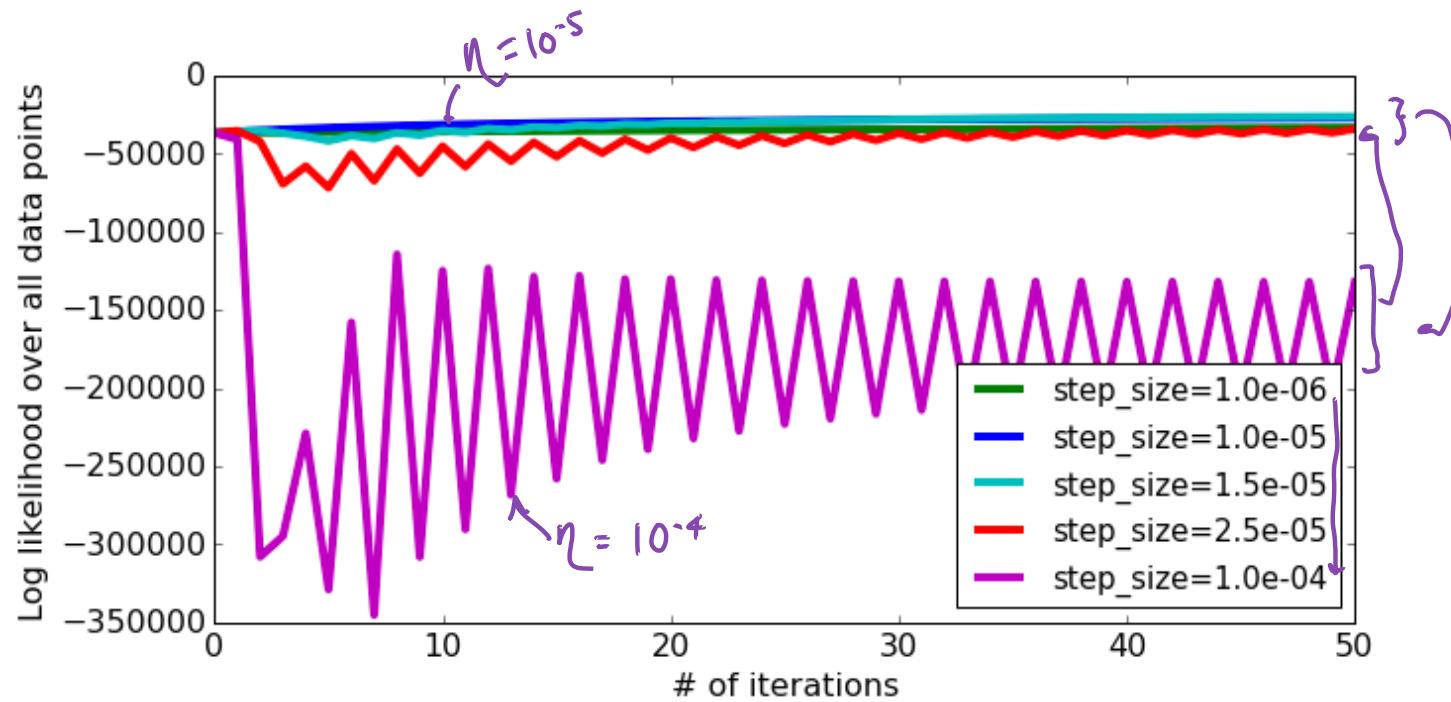
Compare converge with different step sizes



Careful with step sizes that are too large



Very large step sizes can even cause divergence or wild oscillations



Simple rule of thumb for picking step size η

- Unfortunately, picking step size requires a lot of trial and error 😞
- Try several values, exponentially spaced
 - 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 spa~~ce~~, pick one that leads best training data likelihood
- Advanced tip: can also try step size that decreases with iterations, e.g.,

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



Deriving the gradient for
logistic regression

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:

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

*natural
log*

The log trick, often used in ML...

- Products become sums:
 $\ln(a \cdot b) = \ln a + \ln b$ | $\ln \frac{a}{b} = \ln a - \ln b$

- Doesn't change maximum!

- If $\hat{\mathbf{w}}$ maximizes $f(\mathbf{w})$:

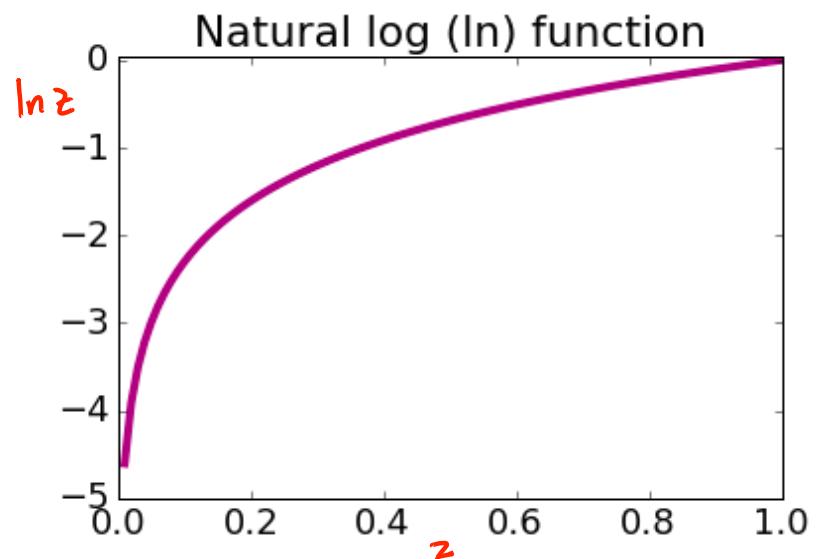
$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} f(\mathbf{w})$$

the \mathbf{w} that makes $f(\mathbf{w})$ largest

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

$$\hat{\mathbf{w}}_{\ln} = \arg \max_{\mathbf{w}} \ln(f(\mathbf{w}))$$

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



Expressing the log-likelihood

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:

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

Rewriting log-likelihood

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

$$\begin{aligned}\ell\ell(\mathbf{w}) &= \sum_{i=1}^N \ln P(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \sum_{i=1}^N [\mathbb{1}[y_i = +1] \ln P(y = +1 \mid \mathbf{x}_i, \mathbf{w}) + \mathbb{1}[y_i = -1] \ln P(y = -1 \mid \mathbf{x}_i, \mathbf{w})]\end{aligned}$$

\checkmark if $y_i = +1$

✓

o

if $y_i = -1$

o

✓



Deriving probability that $y=-1$ given x

Logistic regression model: $P(y=-1|x, \mathbf{w})$

- Probability model predicts $y=+1$:

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

- Probability model predicts $y=-1$:

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

$$\begin{aligned} &= \frac{1 + e^{-\mathbf{w}^\top h(\mathbf{x})} - 1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}} = \frac{e^{-\mathbf{w}^\top h(\mathbf{x})}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}} \end{aligned}$$

Rewriting the log-likelihood

Plugging in logistic function for 1 data point

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

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

$$\begin{aligned} &= \mathbb{1}[y_i = +1] \ln \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} + (1 - \mathbb{1}[y_i = +1]) \ln \frac{e^{-\mathbf{w}^\top h(\mathbf{x}_i)}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} \\ &= -\mathbb{1}[y_i = +1] \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) + (1 - \mathbb{1}[y_i = +1])[-\mathbf{w}^\top h(\mathbf{x}_i) - \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)})] \end{aligned}$$

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

Simpler form

$$\begin{aligned} \ln e^a &= a \\ \mathbb{1}[y_i = +1] &= 1 - \mathbb{1}[y_i = -1] \\ \ln \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} &= -\ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) \\ \ln \frac{e^{-\mathbf{w}^\top h(\mathbf{x}_i)}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} &= -\ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) \\ \ln e^{-\mathbf{w}^\top h(\mathbf{x}_i)} &- \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) \\ -\mathbf{w}^\top h(\mathbf{x}_i) & \end{aligned}$$

Deriving gradient of log-likelihood

Gradient for 1 data point

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

$$\begin{aligned}\frac{\partial \ell\ell}{\partial w_j} &= - (1 - \mathbb{1}[y_i = +1]) \frac{\partial}{\partial w_j} w^\top h(\mathbf{x}_i) - \frac{\partial}{\partial w_j} \ln(1 + e^{-w^\top h(\mathbf{x}_i)}) \\ &= - (1 - \mathbb{1}[y_i = +1]) h_j(\mathbf{x}_i) + h_j(\mathbf{x}_i) P(y_i = -1 | \mathbf{x}_i, \mathbf{w}) \\ &= h_j(\mathbf{x}_i) \left[\mathbb{1}[y_i = +1] - P(y_i = +1 | \mathbf{x}_i, \mathbf{w}) \right]\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial w_j} w^\top h(\mathbf{x}_i) &= h_j(\mathbf{x}_i) \\ \frac{\partial}{\partial w_j} \ln(1 + e^{-w^\top h(\mathbf{x}_i)}) &= -h_j(\mathbf{x}_i) \frac{e^{-w^\top h(\mathbf{x}_i)}}{1 + e^{-w^\top h(\mathbf{x}_i)}} \\ P(y_i = -1 | \mathbf{x}_i, \mathbf{w}) &\end{aligned}$$

Finally, gradient for all data points

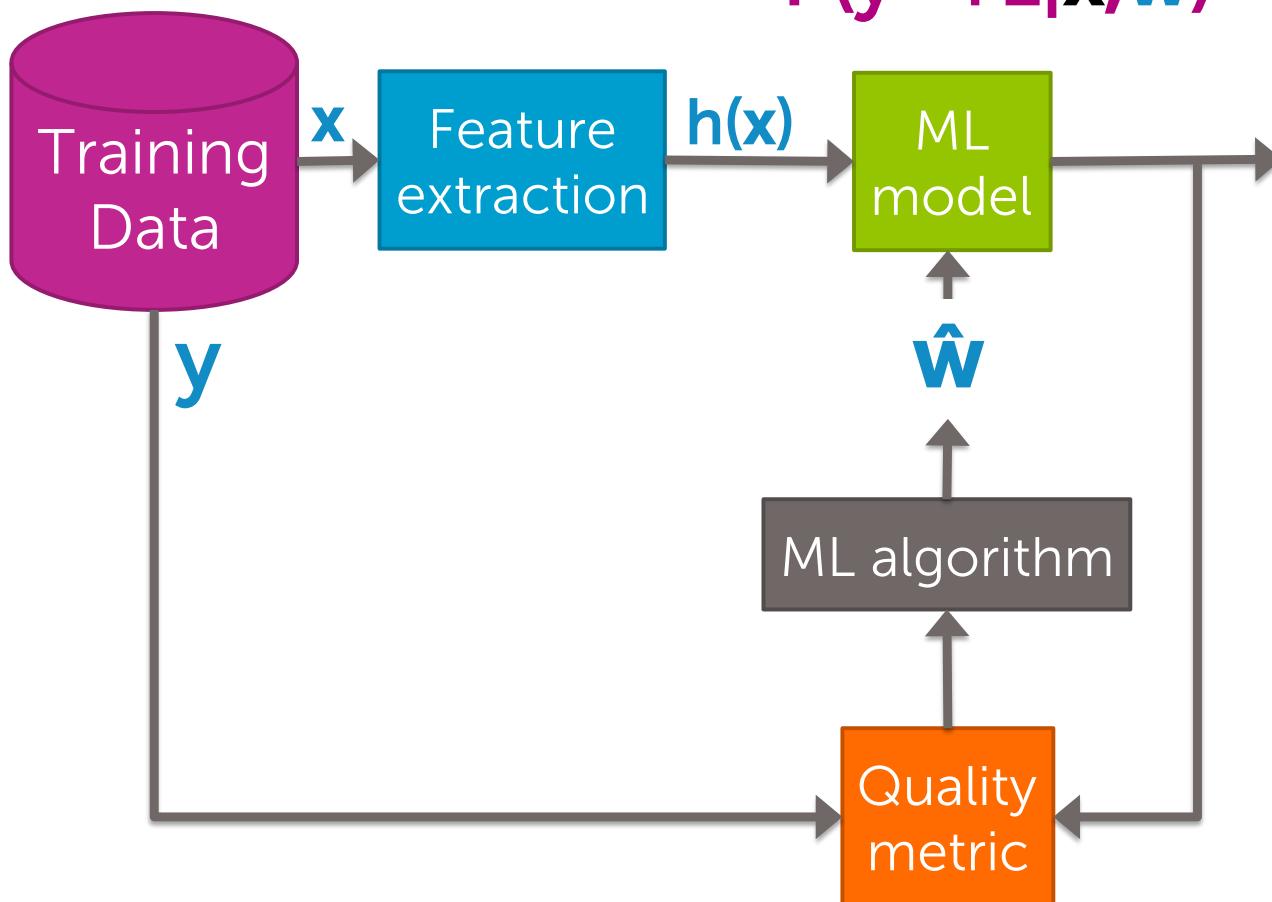
- Gradient for one data point:

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

- Adding over data points:

$$\frac{\partial \ell}{\partial w_j} = \sum_{i=1}^N h_j(\mathbf{x}_i) \left(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w}) \right) \quad \bigg\} \quad \dot{\cup}$$

Summary of logistic regression classifier



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



Learn logistic regression model with maximum likelihood estimation (MLE)

- Choose coefficients \mathbf{w} that maximize likelihood:

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

- No closed-form solution → use gradient ascent



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