

Linear Classifiers & Logistic Regression

Linear classifiers: A motivating example

Classifying sentiment of review

"Easily best sushi in Seattle"

Sentence Sentiment
Classifier

Good

Bad

Linear classifier Intuition

Sentence
from
review

Input X :

Classifier
NODE

Output: y
predicted
class

$$\hat{y} = +1$$

$$\hat{y} = -1$$

A (linear) classifier

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

Word

good
great
awesome
bad
terrible
awful
restaurant, the, we, where

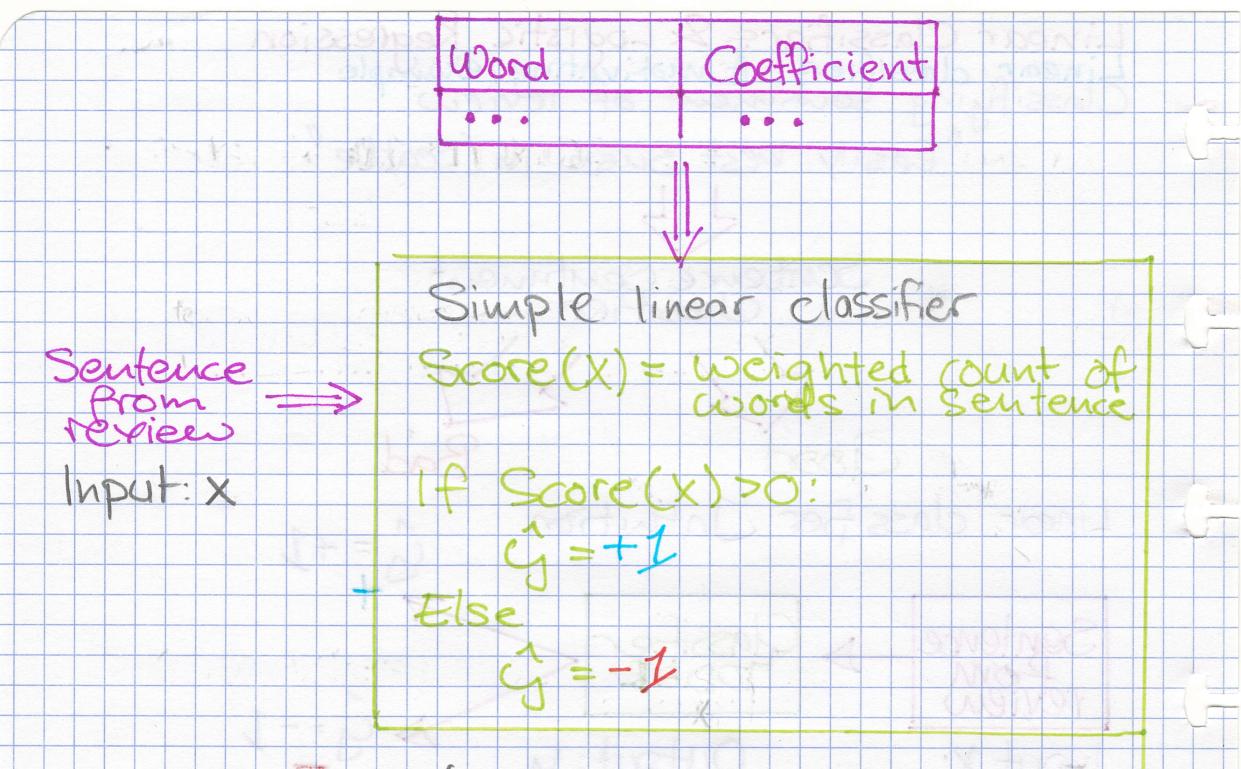
Coefficient

1.0
1.5
2.7
-1.0
-2.1
-3.3
0.0

Input x_i :

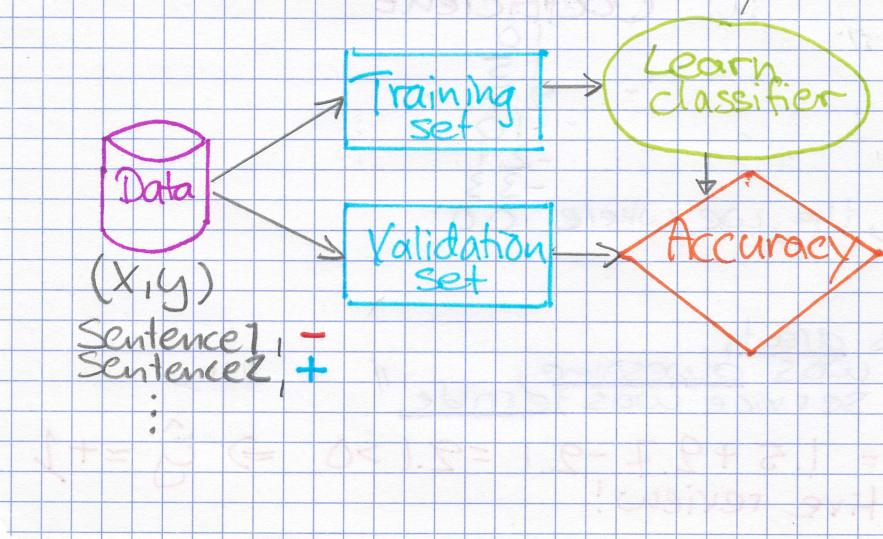
Sushi was great,
the food was awesome,
but the service was terrible.

Score (x_i) = $1.5 + 2.7 - 2.1 = 2.1 > 0 \Rightarrow \hat{y} = +1$
positive review!



Training a classifier = Learning the coefficients

Word	Coefficient
good	1.0
awesome	1.7
bad	-1.0
awful	-3.3
...	...



Decision boundaries

Suppose only two words had non-zero coefficient.

Word Coefficient

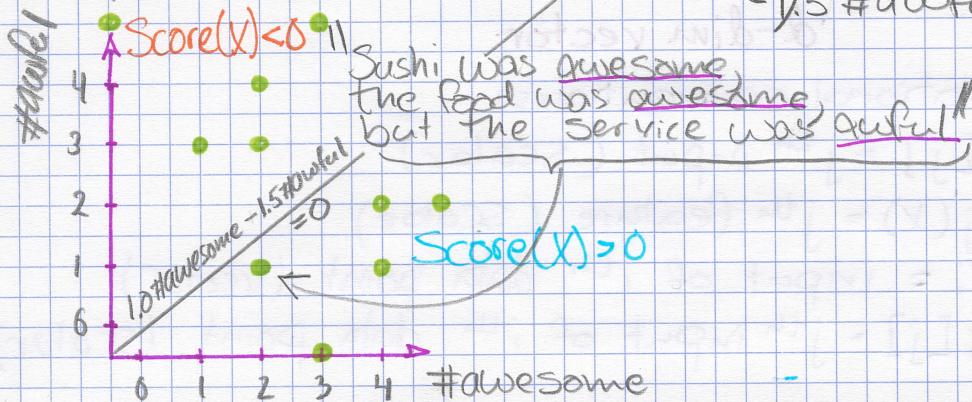
#awesome

1.0

#awful

-1.5

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

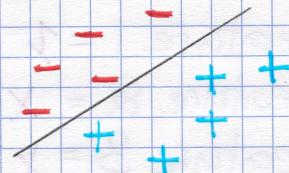


Decision boundary separates positive and negative predictions

- For linear classifiers
 - When 2 coefficients are non-zero
→ line

- When 3 coefficients are non-zero
→ plane

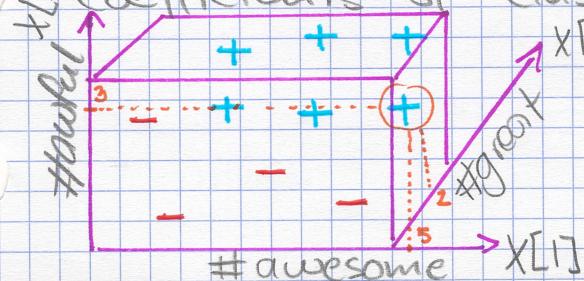
- When many coefficients are non-zero
→ hyperplane



- For more general classifiers
→ more complicated shapes

Linear classifier: Model

(Coefficients of classifier)



$$\text{score}(x) = w_0$$

$$+ w_1 \# \text{awesome}^5$$

$$+ w_2 \# \text{awful}^3$$

$$+ w_3 \# \text{great}^2$$

General notation

Output $y \in \{-1, +1\}$

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

\uparrow
 d -dim vector

Notational conventions:

$x[j]$ = j^{th} input (scalar)

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

x_i = input of i^{th} data point (vector)

$x_i[j]$ = j^{th} input of i^{th} data point (scalar)

Simple hyperplane

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

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

feature 1 = $x[1]$

feature 2 = $x[2]$... e.g. #awesome

feature 3 = $x[3]$... e.g. #awful

⋮
feature $d+1 = x[d]$

More generic features... D-dimensional hyperplane

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

$$\begin{aligned}\text{Score}(x_i) &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) \\ &= \sum_{j=0}^D w_j h_j(x_i) = \vec{w}^T \vec{h}(x_i)\end{aligned}$$

feature 1 = $h_0(x)$... eg., 1

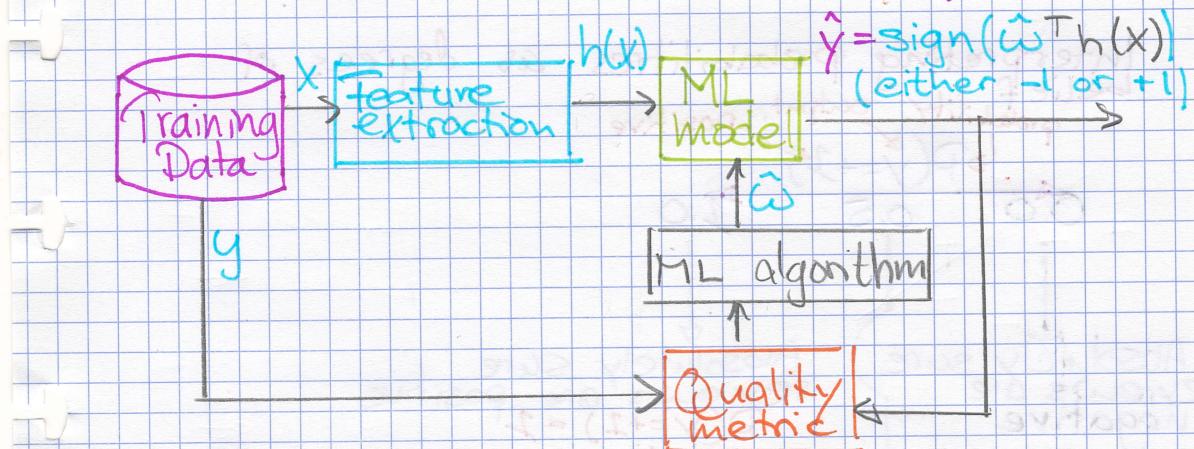
feature 2 = $h_1(x)$... eg $x[1] - \# \text{awesome}$

feature 3 = $h_2(x)$... eg $x[2] = \# \text{awful}$

or, $\log(x[7]) x[2] = \log(\#\text{bad})$
 $x \#\text{awful}$

or, tf-idf("awful")

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



Predicting class probabilities

How confident is your prediction?

- Thus far we've outputted a prediction: $+1$ or -1

- But, how sure are you about the prediction?

"The sushi & everything else were awesome"

Definite $+1$

"The sushi was good
the service was ok"

not sure

Basics of probabilities

Probability a review is positive is 0.7



$x =$
review text
all the sushi was
delicious! Easily best
sushi in Seattle.

The sushi & everything
else were awesome

My wife tried their ramen,
it was pretty forgettable

The sushi was good,
the service was OK

$y =$
sentiment

+1

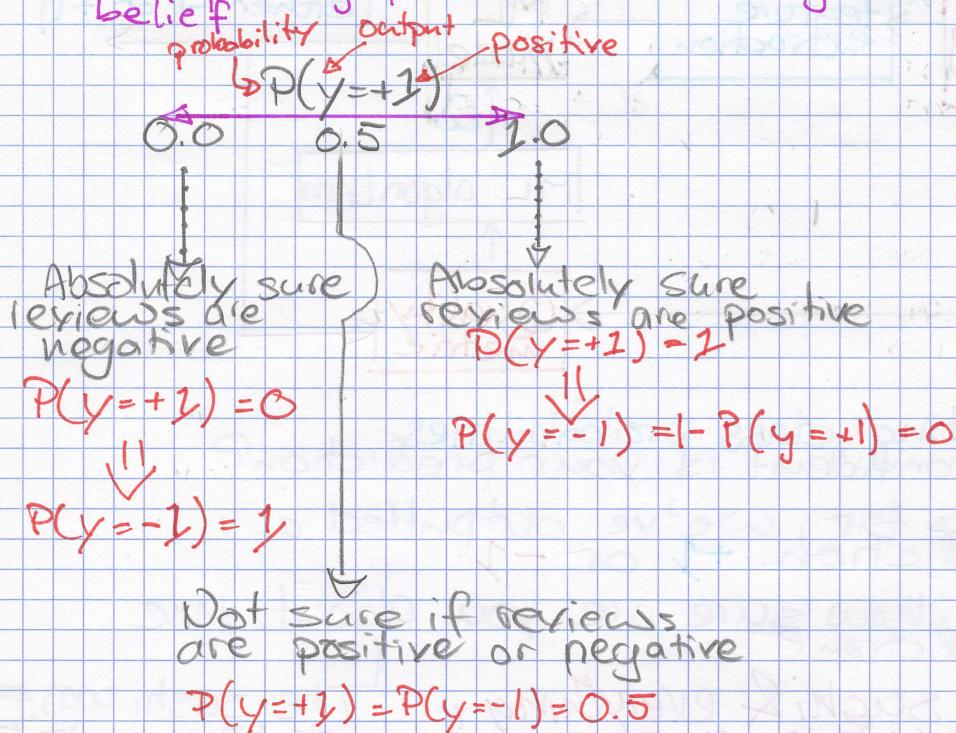
+1

-1

+1

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

Interpreting probabilities as degrees of
belief



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

$$0 \leq P(y=+1) \leq 1$$

$$0 \leq P(y=\text{dog}) \leq 1$$

$$0 \leq P(y=-1) \leq 1$$

$$0 \leq P(y=\text{cat}) \leq 1$$

Probabilities sum
up to 1

$$P(y=+1) + P(y=-1) = 1$$

$$\begin{aligned} P(y=\text{dog}) + P(y=\text{cat}) \\ + P(y=\text{bird}) = 1 \end{aligned}$$

Conditional probabilities

Probability of review with
3 "awesome" and 1 "awful" is positive is 0.9

x = review text

All the sushi was delicious!
Easily best Sushi in Seattle

y = sentiment

+1

{ Sushi was awesome & everything
else was awesome! The service
was awful, but overall awesome }

+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

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

probability output positive given input
 $P(y=+1|x_i = \text{"all the sushi was delicious"})$

Output label

0,0 0,5 1,0

Absolutely sure review "All the sushi was delicious" is negative

$$P(y=+1|x_i) = 0$$

$$\Downarrow$$

$$P(y=-1|x_i) = 1$$

Absolutely sure review "All the sushi was delicious" is positive

$$P(y=+1|x_i) = 1$$

$$\Rightarrow P(y=-1|x_i) = 1 - P(y=+1|x_i) = 0$$

Not sure if review "All the sushi was delicious" positive or negative

$$P(y=+1|x_i) = P(y=-1|x_i) = 0.5$$

Property

Conditional probability always between 0 & 1

Conditional probabilities sum up to 1 over y , but not over x

Two classes
 (e.g. y is +1 or -1,
 x_i is review text)

$$0 \leq P(y=+1|x_i) \leq 1$$

$$0 \leq P(y=-1|x_i) \leq 1$$

$$P(y=+1|x_i) + P(y=-1|x_i) = 1$$

But

$$\sum_x P(y=+1|x) \neq 1$$

Multiple classes
 (e.g. y is dog, cat or bird,
 x_i is image)

$$0 \leq P(y=\text{dog}|x_i) \leq 1$$

$$0 \leq P(y=\text{cat}|x_i) \leq 1$$

$$0 \leq P(y=\text{bird}|x_i) \leq 1$$

$$P(y=\text{dog}|x_i) +$$

$$P(y=\text{cat}|x_i) +$$

$$P(y=\text{bird}|x_i) = 1$$

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)

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

2	1	+1
0	-1	-1
3	-1	-1
4	+1	+1

Optimize quality metric on training data

Find best model by finding best \hat{w}

useful for predicting \hat{y}

Sentence from review \Rightarrow

Input: 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:

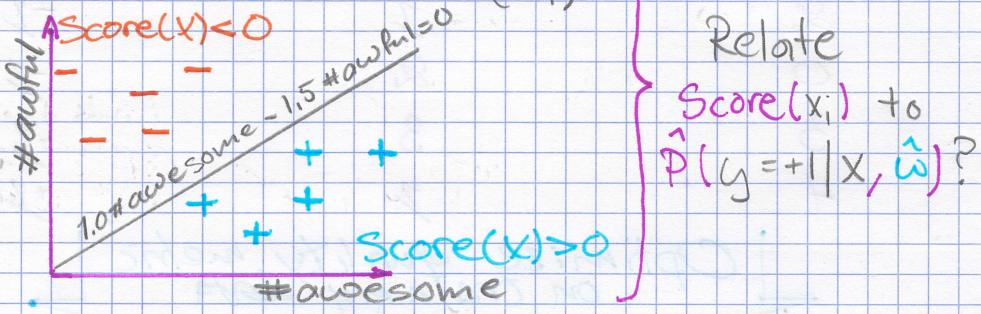
$$\hat{y} = -1$$

- 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

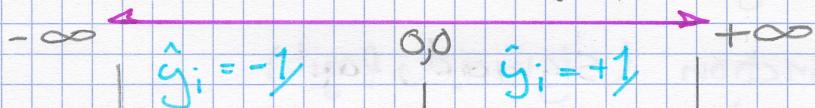
Thus far, we focused on decision boundaries

$$\text{Score}(x_i) = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_d h_d(x_i) \\ = \mathbf{w}^\top h(x_i)$$



Interpreting Score(x_i)

$$\text{Score}(x_i) = \omega^T h(x_i)$$



Very Sure

$$\hat{y}_i = -1$$

Not sure if

$$\hat{y}_i = -1 \text{ or } +1$$

Very Sure

$$\hat{y} = +1$$

$$\hat{P}(y = +1 | x_i) = 0$$

$$P(y = +1 | x_i) = 0.5$$

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

Why not just use regression to build classifier?

$$\text{Score}(x_i) = \omega_0 h_0(x_i) + \omega_1 h_1(x_i) + \dots + \omega_D h_D(x_i)$$

How do we link $-\infty$ to $0, 1$?

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

$$g(\omega^T h(x_i))$$

Link function:
Squeeze real
time into $[0, 1]$

$$\hat{P}(y = +1 | x_i) = g(\omega^T h(x_i))$$

Generalized Linear model

Logistic regression classifier:
linear score with logistic link function

Logistic function (Sigmoid, Logit)

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

input

Score	$-\infty$	-2	0	$+\infty$
Sigmoid	$\frac{1}{1 + e^{-\infty}} = 0$	0.12	$\frac{1}{1 + e^0} = 0.5$	1
	$e^{-\infty} = 0$		$e^0 = 1$	



Score	$+2$	$+\infty$
Sigmoid	0.88	$\frac{1}{1 + e^{-\infty}} = 1$
		$e^{-\infty} = 0$

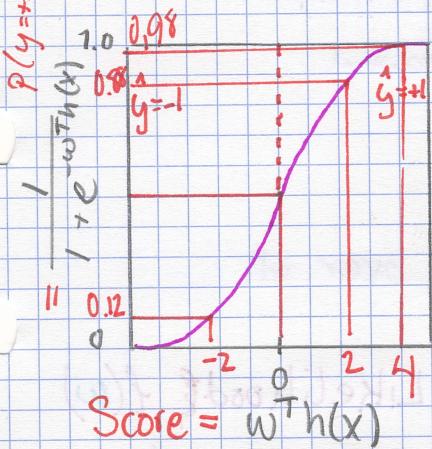
Logistic regression model

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

$$= \frac{1}{1 + e^{-\text{Score}(x_i)}} = \frac{1}{1 + e^{-w^T h(x_i)}}$$

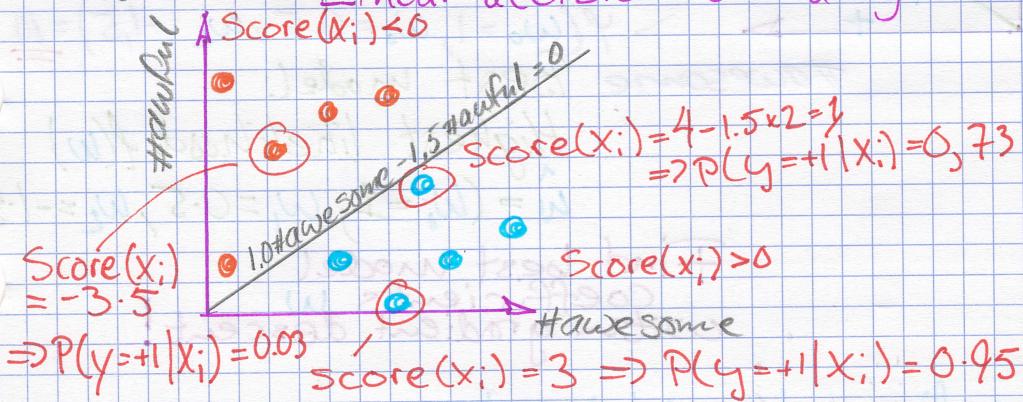
Understanding the logistic regression model

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

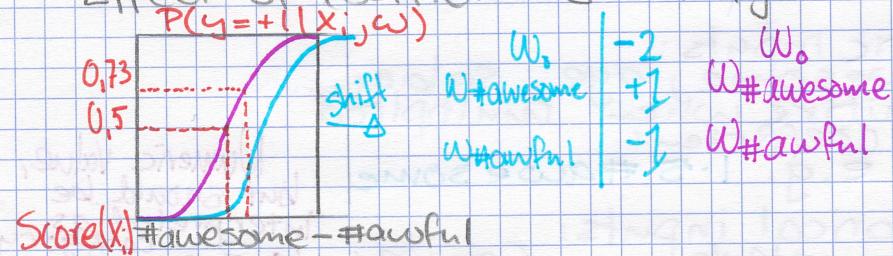


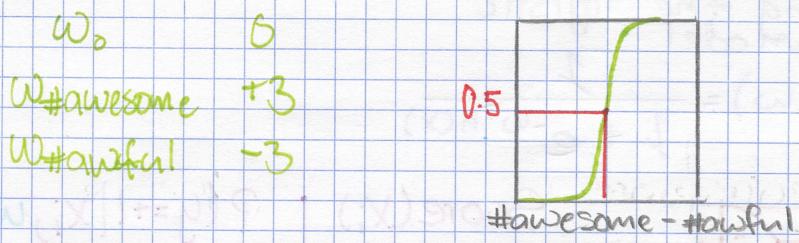
Score(x_i)	$P(y=+1 x_i, \omega)$
0	0.5
-2	$0.12 < 0.5 \Rightarrow y=-1$
2	0.88
4	$0.98 \Rightarrow \hat{y}=+1$

Logistic regression \rightarrow Linear decision boundary



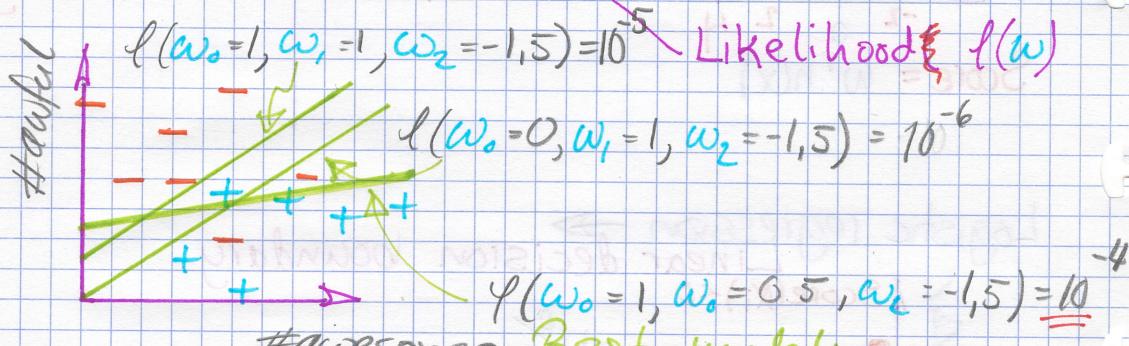
Effect of coefficients on logistic regression model





Overview of learning logistic regression model

Find "best" classifier =
Maximize quality metric over all
possible (w_0, w_1, w_2)



Find best model
coefficients w
with gradient descent!

Encoding categorical inputs

Categorical inputs

- Numeric inputs:
 - = #awesome, age, salary...
 - = intuitive when multiplied by coefficient
 - e.g. 1.5 #awesome
- Categorical inputs:
 - Gender: Male, Female
 - Country of birth

Numeric value,
but should be
interpreted as
category
↳ Zipcode

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

Encoding categories as numeric features

$X = \text{Country of birth}$
(Argentina, Brazil, USA, ...)

1-hot
encoding

196 categories

x	Argentina $h_1(x)$	Brazil $h_2(x)$	Zimbabwe $h_{196}(x)$
Brazil	0	1	0	0	0
Zimbabwe	0	0	0	0	1

196 features

$X = \text{Restaurant reviews}$
(text data)

bag of words

10,000 words
in vocabulary

x	#awesome $h_1(x)$	#awful $h_2(x)$...	#Sushi: $h_{9999}(x)$	$h_{10000}(x)$
Text	2	0	0	0	3

10,000 features

Multiclass classification using "Y" versus all

Multiclass classification

"dog"

Input: x
Image pixels

Top Predictions

Labrador retriever
golden retriever
red bone
bloodhound

Output: y
Object in image

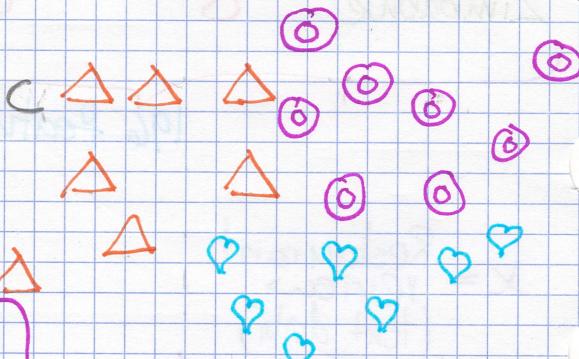
Multiclass classification formulation

- C possible classes

- y can be $1, 2, \dots, C$

- N data points:

Data point	$x[1]$	$x[2]$	y
x_1, y_1	2	1	△
x_2, y_2	0	2	♡
x_3, y_3	3	3	○
x_4, y_4	4	1	○



Learn:

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

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

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

y 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 \circlearrowleft



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



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

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

$$\hat{P}(y=\Delta|x_i)$$

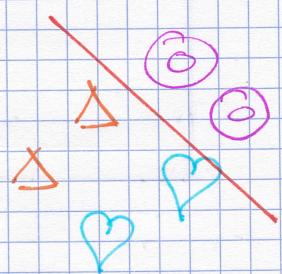
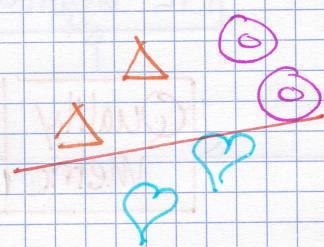
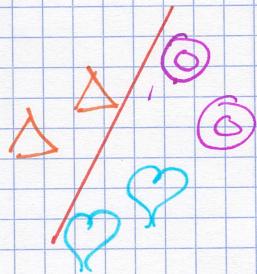
$$= \hat{P}_\Delta(y=+1|x_i, w)$$

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

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

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

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



Multiclass training

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



Input x_i

Predict most likely class

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

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

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

$$\hat{y} = c$$

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

$$\hat{P}(y=+1|x, \hat{\omega}) =$$

