## What is classification?

From features to predictions
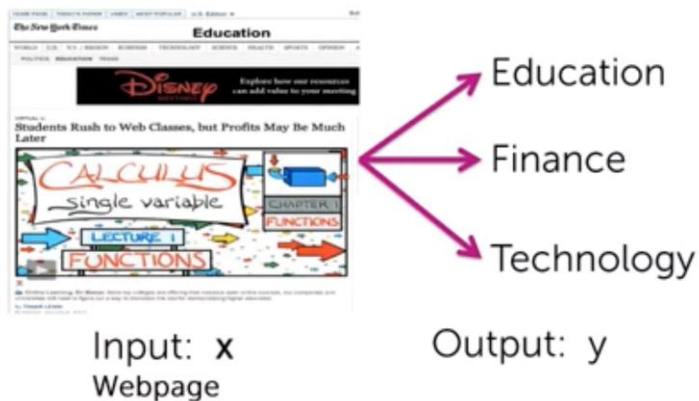
Data → Classifier → Intelligence

Input x:
features derived from data

Learn x→y relationship

Predict y:
categorical "output", class or label

5    ©2015-2016 Emily Fox & Carlos Guestrin

**Multi-class classification**

Given a web page, we have to find out whether that page belongs to 'Education', 'finance', 'technology' based on the content of that webpage.



## Example multiclass classifier
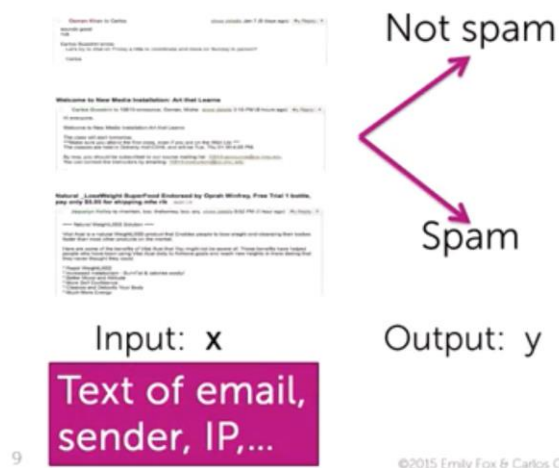*Output y has more than 2 categories*

Education

Finance

Technology

Input: x
Webpage

Output: y

8    ©2015 Emily Fox & Carlos Guestrin

**Famous Example** : Spam filtering

This classifier is for all the patients ir-respective of their personal habits



But personalized medical diagnosis tests our DNA, food-habits, genetical problems, total body metabolism etc to decide what treatment is going to be most effective for that individual person instead a giving a routine flat general treatment.

This is the real world example of classification problem.

**Mind-Reading :**

A person sees an hammer. A FMRI is taken for that particular person, From that Scanned images we can conclude that what a person is seeing whether an 'hammer' or a 'house'

# Linear Classifier

## 1) Linear Classifier Model

## Simple hyperplane

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

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

*awesome*          *awful*

feature 1 = 1
feature 2 = **x**[1] ... e.g., #awesome
feature 3 = **x**[2] ... e.g., #awful
...
feature d+1 = **x**[d] ... e.g., #ramen

26                                    ©2015-2016 Emily Fox & Carlos Guestrin          Machine Learning Specialization

$\mathbf{w}^\top \quad \mathbf{x}_i$

$w_0$
$w_1$
$w_2$
$\vdots$
$w_d$

$\mathbf{x}_i[1]$
$\mathbf{x}_i[2]$
$\vdots$
$\mathbf{x}_i[d]$

For a single row representing → xi
Score for that single row → score(xi)
Sign of the score of that single row → y^(i) = Sign(score(xi))

i from 1 to N (row wise)   [xi]


## The Role of *sign*

If the score>0 then predict → +1

If the score<0 then predict → -1

At 0, we have the choice to predict either -1/+1.  You make an arbitrary choice

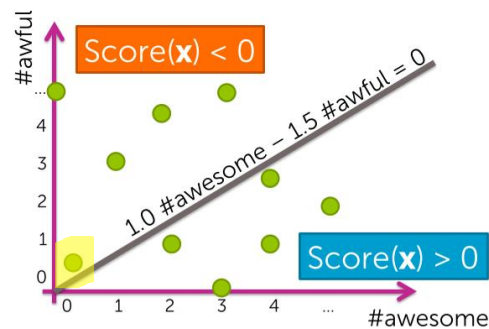## 2) Effects of coefficient values on decision boundary

Initially the intercept → 0



Now the intercept → 1 and the line slightly shifts up, so the orange point which was -ve before, now becomes +ve

After changing the w2=-3.0, the line gets modified, and the blue point which was +ve Now becomes -ve.

## Decision boundary: effect of changing coefficients

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

$\Rightarrow$ Score(x) = **1.0** + 1.0 #awesome − **3.0** #awful

Score(x) < 0

Score(x) > 0

#awful

#awesome

29    ©2015-2016 Emily Fox & Carlos Guestrin    Learning Specialization

From this we can conclude that the coefficients are playing a very important role in the classification.

**3) Using features of inputs**

## More generic features... D-dimensional hyperplane

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

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

$= \sum_{j=0}^{D} w_j h_j(x_i) = w^T h(x_i)$

feature 1 = $h_0(x)$ ... e.g., 1
feature 2 = $h_1(x)$ ... e.g., x[1] = #awesome
feature 3 = $h_2(x)$ ... e.g., x[2] = #awful
　　　　　　　　or, log(x[7]) x[2] = log(#bad) x #awful
　　　　　　　　or, tf−idf("awful")

...

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

30    ©2015-2016 Emily Fox & Carlos Guestrin    Machine Learning Specialization

# Probabilities and its basics

Not all our output will be exactly +1 or -1, especially the output of logistic regression will be like 0.432 (or) 0.211. So in-order to conclude them as +1 or -1 (Here probability comes into picture)

## 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!" | "The sushi was good, the service was OK" |
|---|---|
| Definite **+1** | Not sure |
| $\hat{y} = +1$ with high probability | $\hat{y} = +1$ with probability 0.5 |

©2015-2016 Emily Fox & Carlos Guestrin

## 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)

©2015-2016 Emily Fox & Carlos Guestrin    Machine Learning Specialization

From the above it is assumed on an average,
70% → +ve reviews and the
remaining 30% → -ve reviews

# Interpreting probabilities as degrees of belief

Probability → output → positive

$$P(y=+1)$$

0.0 —————— 0.5 —————— 1.0

Absolutely sure reviews are negative

$P(y=+1)=0$
⇓
$P(y=-1)=1$

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

Not sure if reviews are positive or negative

Absolutely sure reviews are positive

$P(y=+1)=1$
⇓
$P(y=-1)=1-P(y=+1)$
$=0$

37

P(y=+1)->1[Complete +ve reviews],P(y=-1)→0 [No negative reviews]

P(y=-1)->1[Complete -ve reviews],P(y=1)→0  [No positive reviews]

# 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 | $0 \leq P(y=+1) \leq 1$  $0 \leq P(y=-1) \leq 1$ | $0 \leq P(y=dog) \leq 1$  $0 \leq P(y=cat) \leq 1$  $0 \leq P(y=bird) \leq 1$ |
| Probabilities sum up to 1 | $P(y=+1)+P(y=-1)=1$ | $P(y=dog)+P(y=cat)+P(y=bird)=1$ |

3

# Conditional Probability



**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)

39   ©2015-2016 Emily Fox & Carlos Guestrin   Machine Learning Specialization

The given condition is that ("3 awesome and 1 awful"). In this given condition, it is observed that 90% of the reviews are +ve and the remaining 10% are -ve.



week-1 logistic-regression-model-annotated.pdf - Drawboard PDF

## Interpreting conditional probabilities

Output label   positive   Given   Input sentence

Probability   $P(y=+1|\mathbf{x}_i = $ "All the sushi was delicious!"$)$

0.0      0.5      1.0

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

$P(y=-1|xi) = 1$
$P(y=+1 |xi) = 1- P(y=-1|xi)$
$= 1-1$
$= 0$
$P(y=-1|xi) = 0$

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

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

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

$P(y=+1|xi) = 1$
$P(y=-1 |xi) = 1- P(y=+1|xi)$
$= 1-1$
$= 0$
$P(y=-1|xi) = 0$

40   ©2015-2016 Emily Fox & Carlos Guestrin   Machine Learning Specialization

20:28
02-07-2020
ENG

## 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 | $0 \le P(Y=+1 \mid x_i) \le 1$ <br> $0 \le P(y=-1 \mid x_i) \le 1$ | $0 \le P(Y=dog \mid x_i) \le 1$ <br> $0 \le P(y=cat \mid x_i) \le 1$ <br> $0 \le P(Y=bird \mid x_i) \le 1$ |
| Conditional probabilities sum up to 1 over y, but not over **x** | $P(Y=+1 \mid x_i) + P(Y=-1 \mid x_i) = 1$ <br> But <br> $\sum_{x} P(Y=+1 \mid x) \ne 1$ <br> $\sum_{i=1}^{N} P(Y=+1 \mid x_i) \ne 1$ | $P(Y=dog \mid x_i) + P(Y=cat \mid x_i) + P(Y=bird \mid x_i)$ <br> $= 1$ |

41

Machine Learning Specialization

## Probabilities used in Classification



P(A|B) → The probability of How much A is in B

# How confident is your prediction?

| "The sushi & everything else were awesome!" | "The sushi was good, the service was OK" |
|---|---|

| Definite +1 | Not sure |
|---|---|

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

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

Many classifiers provide a degree of certainty:

Output label

Input sentence

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

Extremely useful in practice

---

# Goal: Learn conditional probabilities from data

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 |
| ... | ... | ... |

Optimize **quality metric** on training data

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

Useful for predicting $\hat{y}$

## Predict most likely class

Sentence from review

Input: $\mathbf{x}$

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

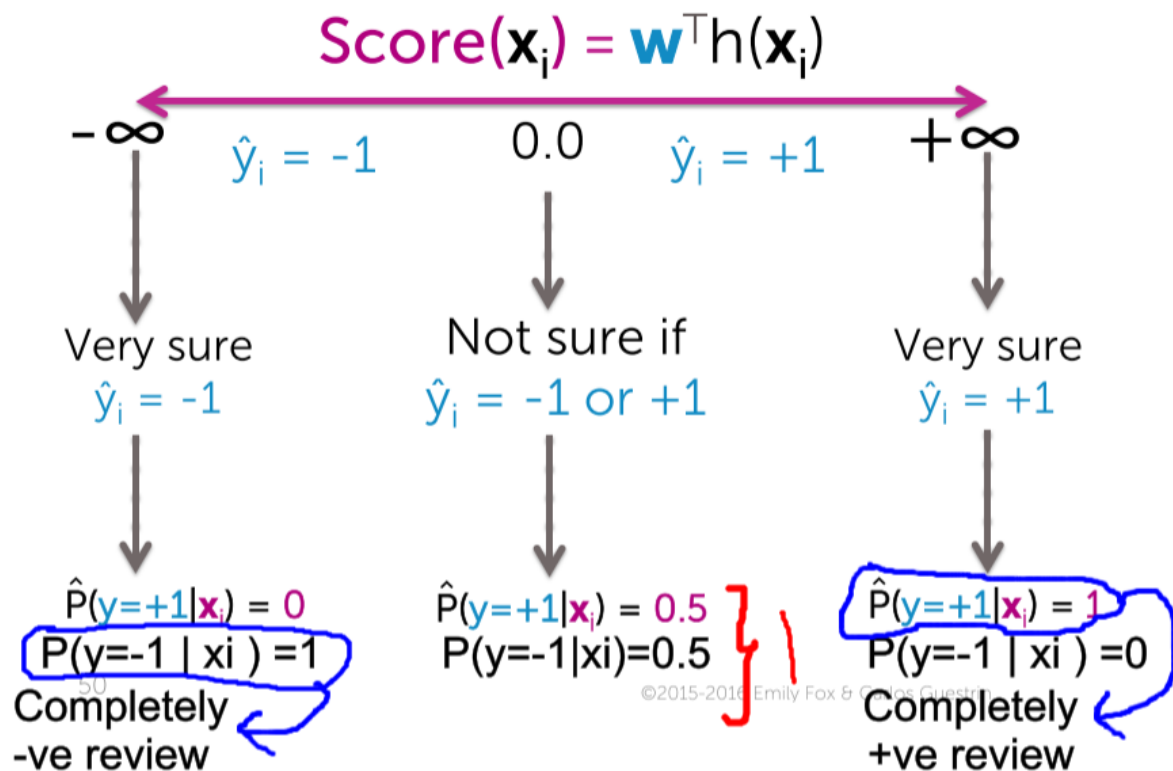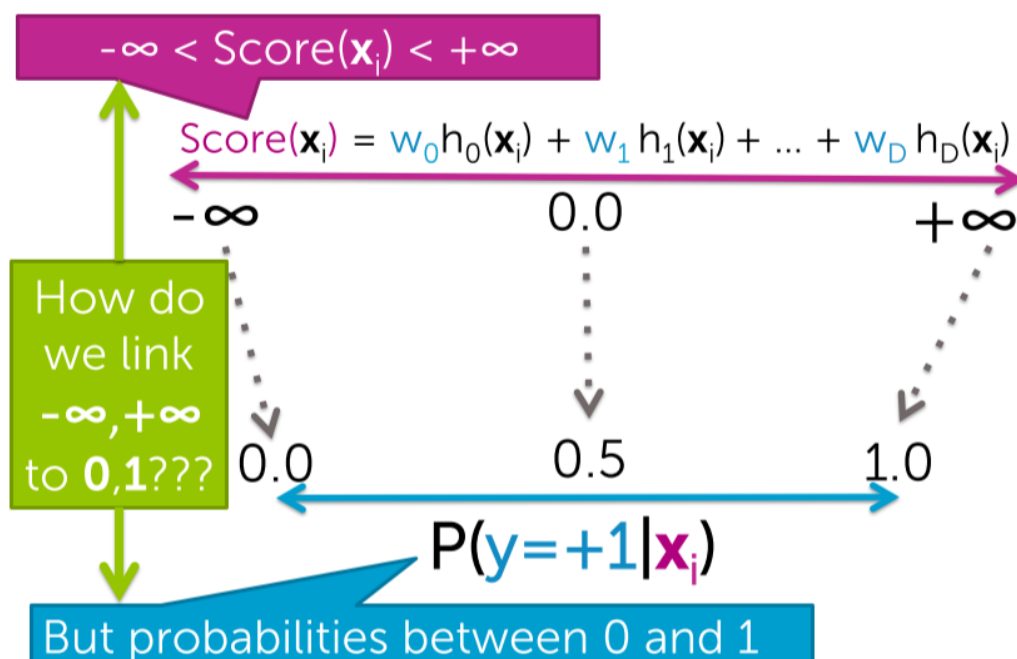### Logistic Regression

$$\text{Score}(\mathbf{x}_i) = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \ldots + w_D h_D(\mathbf{x}_i)$$
$$= \mathbf{w}^\top h(\mathbf{x}_i)$$



#awful

Score($\mathbf{x}$) < 0

1.0 #awesome − 1.5 #awful = 0

Score($\mathbf{x}$) > 0

#awesome

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

We know below the line the score(x)>0 and above the line the score(x)<0. But we don't know how much far is it less/greater than 0

## 1) Predicting the class probabilities with (generalized) linear model

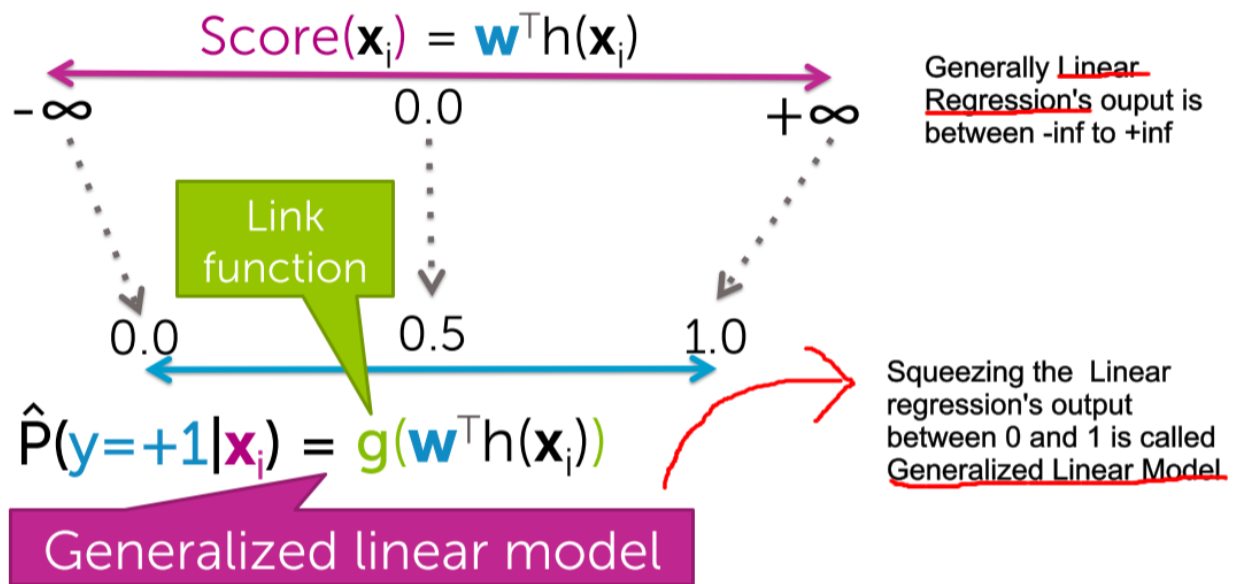## Interpreting Score($x_i$)

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

$-\infty$ ⟵ $\hat{y}_i = -1$     0.0     $\hat{y}_i = +1$     $+\infty$

| Very sure | Not sure if | Very sure |
|---|---|---|
| $\hat{y}_i = -1$ | $\hat{y}_i = -1$ or $+1$ | $\hat{y}_i = +1$ |

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

$P(y=-1 \mid x_i) = 1$

Completely
-ve review

$\hat{P}(y=+1|x_i) = 0.5$
$P(y=-1|x_i) = 0.5$

©2015-2016 Emily Fox & Carlos Guestrin

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

$P(y=-1 \mid x_i) = 0$

Completely
+ve review

## Why not just use regression to build classifier?

$-\infty < \text{Score}(x_i) < +\infty$

$$\text{Score}(x_i) = w_0 h_0(x_i) + w_1 h_1(x_i) + \ldots + w_D h_D(x_i)$$

$-\infty$     0.0     $+\infty$

How do
we link
$-\infty, +\infty$
to **0,1**???    0.0     0.5     1.0

$P(y=+1|x_i)$

But probabilities between 0 and 1

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

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

$$-\infty \qquad\qquad 0.0 \qquad\qquad +\infty$$

Generally Linear Regression's ouput is between -inf to +inf

**Link function**

$$0.0 \qquad\qquad 0.5 \qquad\qquad 1.0$$

$$\hat{P}(y=+1|\mathbf{x}_i) = g(\mathbf{w}^\top h(\mathbf{x}_i))$$

Squeezing the Linear regression's output between 0 and 1 is called Generalized Linear Model

**Generalized linear model**

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization



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

©2015-2016 Emily Fox & Carlos Guestrin

©2015-2016 Emily Fox & Carlos Guestrin

## 2) The sigmoid(logistic) Link function

### Logistic function (sigmoid, logit)

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

input

| Score | $-\infty$ | -2 | 0.0 | +2 | $+\infty$ |
|---|---|---|---|---|---|
| sigmoid(Score) | $\frac{1}{1+e^{\infty}}$ $= \frac{1}{1+\infty}$ $= 0$ | 0.12 | $Sigmoid(0)$ $=\frac{1}{1+e^{0}}$ $=\frac{1}{1+1}$ $= 0.5$ | 0.88 | $\frac{1}{1+e^{-\infty}}$ $= 1$ |

$e^{\infty} = \infty$     $e^{0} = 1$     $e^{-\infty} = 0$

56
©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

We are giving the score whose output is (-inf to +inf)

Wants to change the output   (0 to +inf)

This is done using Link function which is Sigmoid function

## 2)Logistic Regression Model

### Understanding the logistic regression model

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

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

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

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

>0.5       <0.5

58
©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

-ve reviews       +ve reviews

Taking Score(x-axis) and Probability(y-axis)

**Effect of coefficient values on predicted probabilities**



**4) Effect of coefficient on Logistic Regression Model**

1) $w_0 \Rightarrow 0$

$w_{\# awesomes} \Rightarrow +1$

$w_{\# awful} \Rightarrow -1$

$\boxed{Score \,(x_i) \Rightarrow w_1 \# awesome + w_2 \# awful + constant}$

$\Rightarrow +1 -1 +0$

$Score \,(z_i) \Rightarrow 0$

$f(x) \Rightarrow \dfrac{1}{1+e^{-(Score \,(x_i))}}$

$\Rightarrow \dfrac{1}{1+e^{-0}} \Rightarrow \dfrac{1}{2} \Rightarrow 0.5$

$f(x) \Rightarrow 0.5$

2) $w_0 \Rightarrow 0$

$w_{\# awesomes} \Rightarrow 1$

$w_{\# awfuls} \Rightarrow -1$

Increasing # of awesomes

$Score \,(x_i) \Rightarrow 2 -1 +0 \Rightarrow 1$

$f(x) \Rightarrow \dfrac{1}{1+e^{-1}} \Rightarrow \dfrac{1}{1.3678} \Rightarrow 0.73$

$w_0 \Rightarrow -2$

$w_{\# awesome} \Rightarrow +1$

$w_{\# awful} \Rightarrow -1$

$Score \,(x_i) \Rightarrow (w_1 \# awesomes + w_2 \# awfuls + w_0)$

$Score \,(x_i) \Rightarrow 2$

Here if my make the $Score \,(x_i) \Rightarrow 2$, then only I will get probability $\Rightarrow \underline{0.5}$

$w_0 \Rightarrow 0$

$w_{\# awesome} \Rightarrow +3$

$w_{\# awful} \Rightarrow -3$

$Score \,(x_i) \Rightarrow (w_1 \# awesome + w_2 \# awful + w_0)$

$\Rightarrow 3-3 \Rightarrow 0$

$f(x) \Rightarrow 0.5$

Increasing # awesomes

Then the probability $\Rightarrow 0.9$ (Nearly 1)
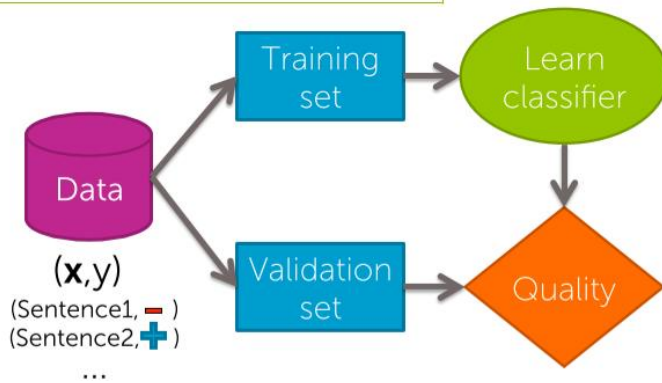
## Conclusion

- See the effect of coefficients affecting the probability. We can conclude if the model has some bigger coefficients then the probabilities can be found more quickly.
- Changing the constant will shift the line to the left and to the right.

**Now we want to find the coefficients that best fits**

# 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 |
| ... | ... | ... |

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



The data is splitted up into training and validation set.

In training set, running a learning algorithm and output the parameter w^.

This w^ is fitted into the model, to estimate the probability that the input sentence either +ve (or) -ve.

Now we can take the validation set and fit into the model . And can predict the quality metric, error etc..

# How to choose w^ ????

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

Likelihood $\ell(w)$

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

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

Find best model coefficients **w** with gradient ascent!

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

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

#awful

#awesome

66

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

Quality metric → likelihood function l(w)
We are saying among the three 10^-4 → best likelihood
( Since for the best classifier  → Maximise the likelihood l(w)  )

We want the required w^, in which l(w) → maximum. So in-order to choose that w^
gradient ascent comes to the picture.

# 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)

- Categorical inputs:

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

69

©2015-2016 Emily Fox & Carlos Guestrin

---

Basically in numerical data,
in score function → we will multiply the numeric with the coefficient.

The zip-code is 10005,98195 etc. It is not meant that 98195 is 9 times of 10005.
These are different postal codes representing different parts of the country.
Hence they are not numerical features, they are categorical features.

But how to multiply the coefficient with the categorical values???
This is achieved by **encoding technique**

# Encoding Categorical inputs



Encoding categories as numeric features

If somebody is born in Brazil → then for Brazil put 1 and for everything put 0
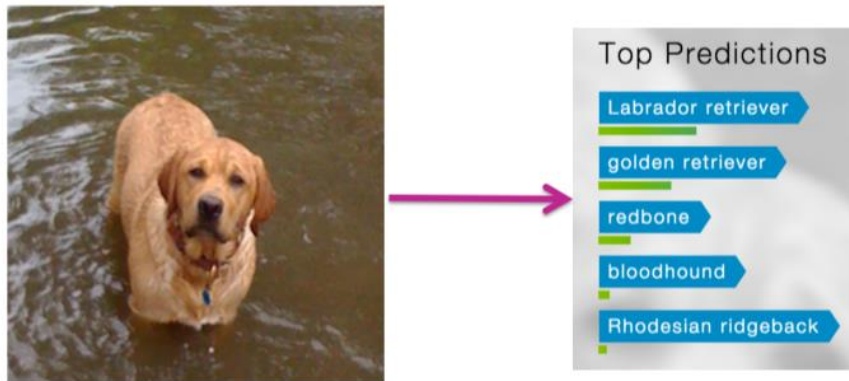If somebody is born in Zimbabwe → then for Zimbabwe put 1 and for everything put 0

How to encode a restaurant's review???
Take 1 review, and put down the word count in the respective places
In the above case , 1st review contains (2 → awesome) , (0 → awful) ,....,(3→sushi)

# Multiclass classification (1 versus all)

## Multiclass classification



Input: x
Image pixels

Output: y
Object in image

In this image → there is a dog
Our aim is to predict 1) Whether is it a dog ?? 2) What kind of dog is it???
Here we are not having only 2 categories . There a nearly 1000's of categories.

**How to solve this???** → **One versus ALL**

Eg: Triangle, donut and hearts are different classes

## Multiclass classification formulation

- C possible classes:
  – y can be 1, 2,..., C
- N datapoints:

| 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:

$\hat{P}(y=▲|x)$

$\hat{P}(y=♥|x)$

$\hat{P}(y=◎|x)$

75

©2015-2016 Emily Fox & Carlos Guestrin

There are 3 classes,

What I need to know is for a particular input, whether is it a triangle, hearts or donut.

Now I want to classify the triangle from the rest



**1 versus all:**
Estimate $\hat{P}(y=\triangle \,|x)$ using 2-class model

+1 class: points with $y_i = \triangle$
-1 class: points with $y_i = \heartsuit$ OR $\bigcirc$

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

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

*More likely to be △*

Score (x:) > 0
⇓
P(y=△|x_i, w) > 0.5

*Not more likely to be a △*

Score (x:) < 0
P(y = △|x_i, w) < 0.5

77

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

Now we are going to make a classifier to learn → that separates the triangle from the donuts and hearts. This train classifier outputs +1, if the input x is more likely to be a triangle.
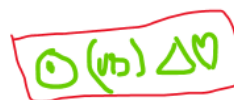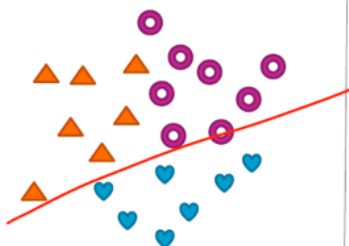


**1 versus all:** simple multiclass classification
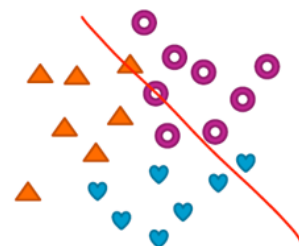using C 2-class models

△ (vs) ⊙ ♡     ♡ (vs) ⊙ △     ⊙ (vs) △♡

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

78

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

In the same data-set

1st classifier → triangle . Note down the probability

2nd classifier -> hearts. Note down the probability

3rd classifier -> donot . Note down the probability

Among these three classifier , capture a classifier which has the highest probability. For eg: If the 2nd classifier (i.e heart classifier) has the highest probability means that classifier ( i.e w (heart) ) make the classification more exact and accurate.

**Take the dog classification example**

Iterate over all the classifier and finally note which classifier has the highest probability and that classifier classifies correctly when compared to the rest of the classifer.