

DataStar Machine Learning



ADAX

Modules

- 01. Introduction
- 02. Regression
- 03. Classification**
- 04. Ensemble Methods & Cross-Validation
- 05. Machine Learning Algorithms
- 06. Regularization Techniques
- 07. Introduction to Unsupervised ML
- 08. Dimensionality Reduction Techniques
- 09. Clustering Techniques
- 10. Introduction to Natural Language Processing



Session 3: Classification

18 Sept 2017

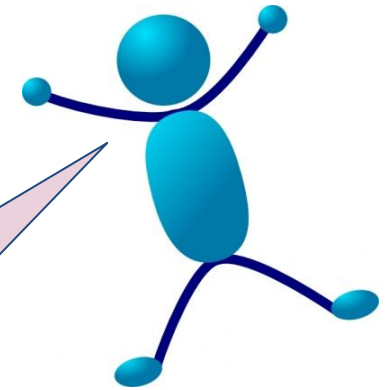


ADAX

A Case Study



Kuala Lumpur has many 5-star sushi restaurants



What are people saying about the food? The ambiance?



A Case Study

Positive reviews are NOT positive everytime, and maybe about something but NOT everything...



Marufuku Ramen
★★★★★ 465 reviews
\$\$ - Ramen



Kui Shin Bo
★★★★☆ 550 reviews
\$\$ - Japanese, Sushi Bars



Ichido
★★★★★ 55 reviews
\$\$\$ - Seafood, Japanese, Asian Fusion

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 Seattle.

Experience



A Case Study

Classifying Sentiment of Review

Easily best sushi in Seattle.



Sentence Sentiment
Classifier



A Case Study

Classifying Sentiment of Review

All reviews
for restaurant

★★★★★ 7/21/2015
This is probably my favorite place to eat Japanese in Seattle. My boyfriend and I ordered nigiri of scallop, Japanese snapper (seasonal), and the agedashi tofu and 2 special rolls. I would skip the special rolls, because the nigiri and sashimi cuts is where this place excels. The tofu, as recommended by other helpers was amazing. It's more chewy and the sauce/soy is the perfect amount of flavor for the delicate tofu.

★★★★★ 6/11/2015
Dining here at the sushi bar made me feel like sitting front row to an amazing performance. We didn't have much, banged down to the 10 after work, got here breathlessly at 8:10pm, and got the last two seats in the place.

★★★★★ 9/19/2016
I came here having high expectations due to the reviews of this place, but I was let disappointed. The restaurant is small so do make reservations when you come here. Dishes cost from \$4-26 each and dishes are small.

Break all reviews
into sentences

The seaweed salad was just OK,
vegetable salad was just ordinary.

I like the interior decoration and
the blackboard menu on the wall.

All the sushi was delicious.

My wife tried their ramen and
it was pretty forgettable.

The sushi was amazing, and
the rice is just outstanding.

The service is somewhat hectic.

Easily best sushi in Seattle.

Sentence
Sentiment
Classifier



Good



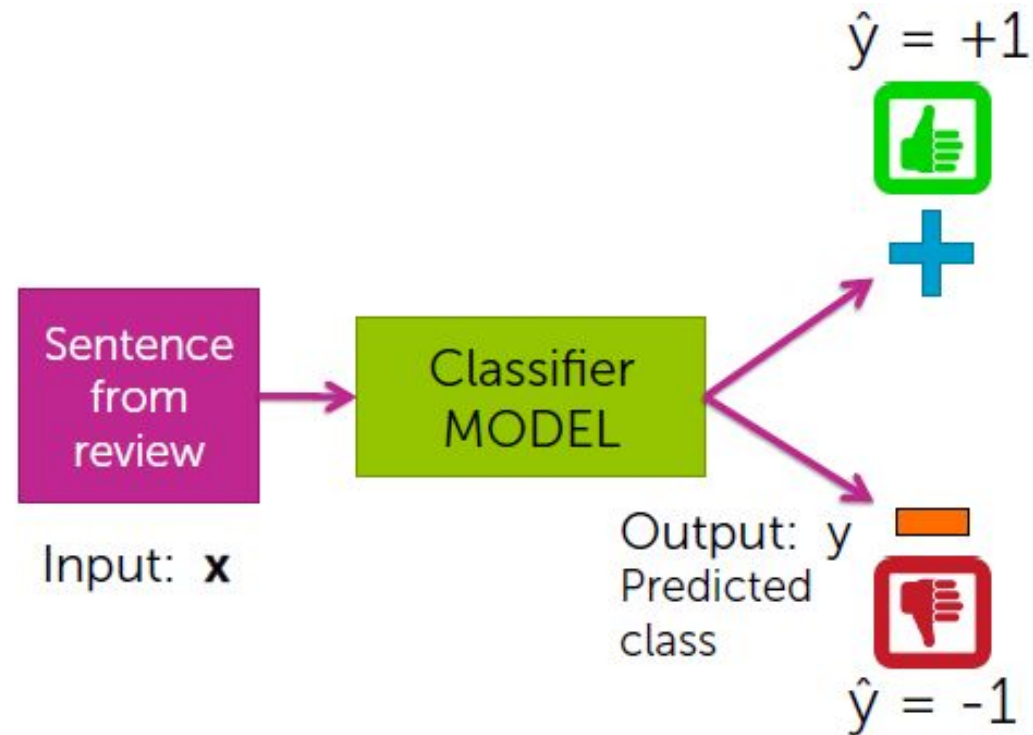
Bad



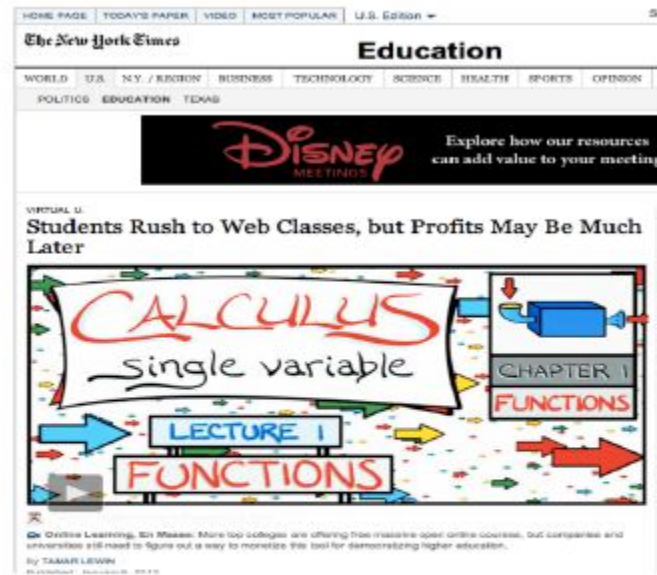
The background is a solid red color. In the four corners, there are triangular regions filled with a halftone pattern of small, dark red dots.

What is Classification?

What is Classification?



Classification can have more than 2 categories



Input: x
Webpage

Education

Finance

Technology

Output: y



Spam Filtering

Dorian Khan to Carlos show details Jan 7 (6 days ago) My Reply

seems good
rob

Carlos Guestrin wrote:
Let's try to chat on Friday a little to coordinate and meet on Sunday in person?

Carlos

Not spam

Welcome to New Media Installation: Art that Learns

Carlos Guestrin to 10015-announce, Carmen, Mohr show details 3:10 PM (8 hours ago) My Reply

Hi everyone,

Welcome to New Media Installation: Art that Learns

The class will start tomorrow.

Make sure you attend the first class, even if you are on the West List!

The classes are held in Doremy Hall C219, and will be Tue, Thu 01:30-4:20 PM.

By now, you should be subscribed to our course mailing list: 10015-announce@cs.stanford.edu

You can contact the instructors by emailing: 10015.instructors@cs.stanford.edu

Natural LoseWeight SuperFood Endorsed by Oprah Winfrey, Free Trial 1 bottle, pay only \$5.95 for shipping mfw rtk show details

Jacquelyn Halley to nherlein, bco, bwherney, bco, ang show details 9:52 PM (1 hour ago) My Reply

== Natural WeightLoss Solution ==

Vital Acai is a natural WeightLoss product that Enables people to lose weight and cleansing their bodies faster than most other products on the market.

Here are some of the benefits of Vital Acai that You might not be aware of. These benefits have helped people who have been using Vital Acai daily to Achieve goals and reach new heights in their doing that they never thought they could.

- Rapid WeightLoss
- Increased metabolism - Burnfat & calories easily!
- Better Mood and Attitude
- More Self Confidence
- Cleanse and Detoxify Your Body
- Much More Energy

Spam

Input: x

Text of email,
sender, IP,...

Output: y



Image Classification



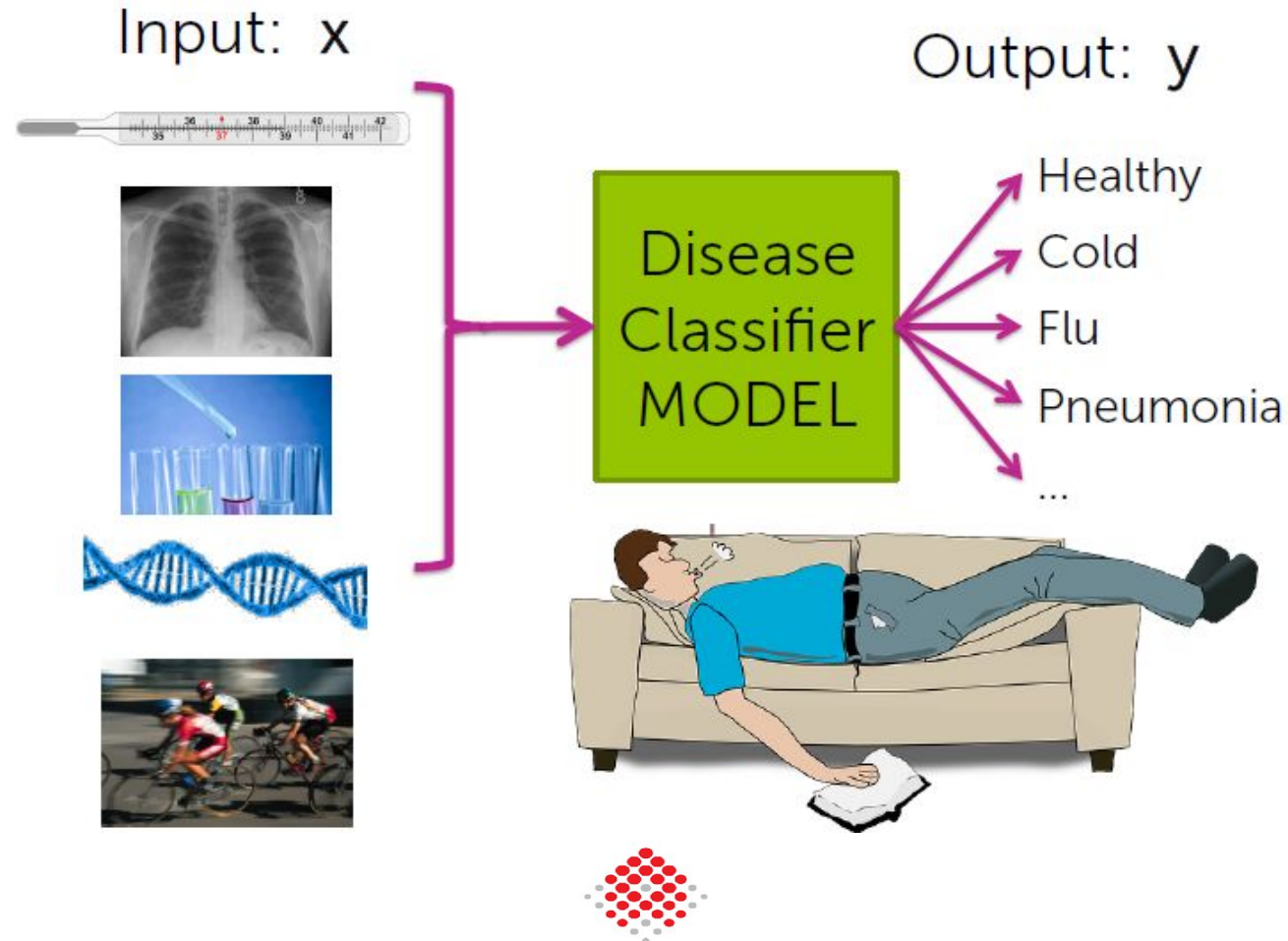
Input: x
Image pixels



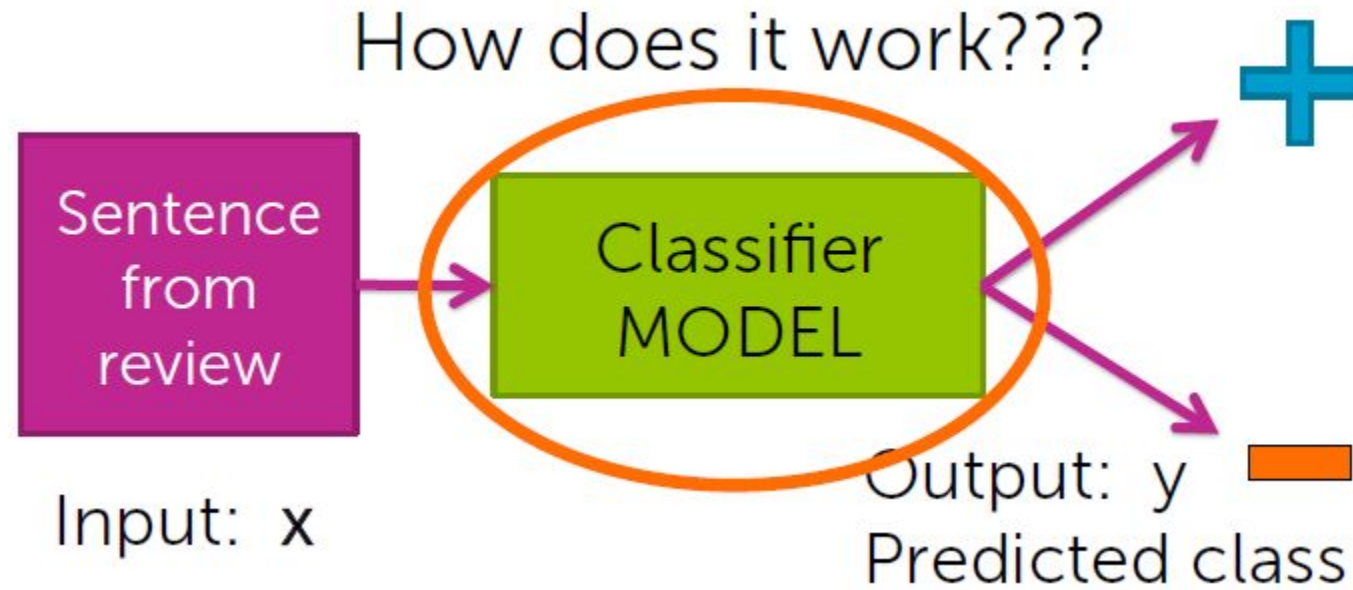
Output: y
Predicted object



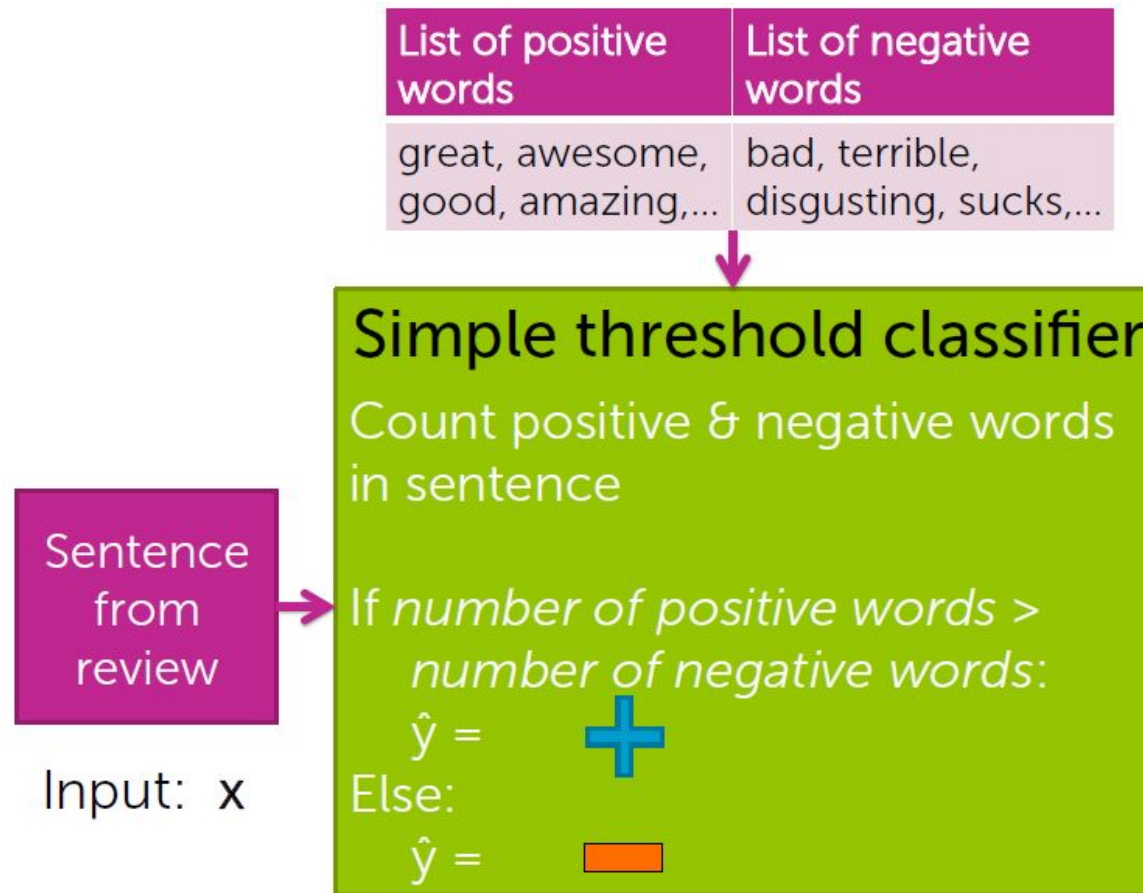
Medical Diagnosis



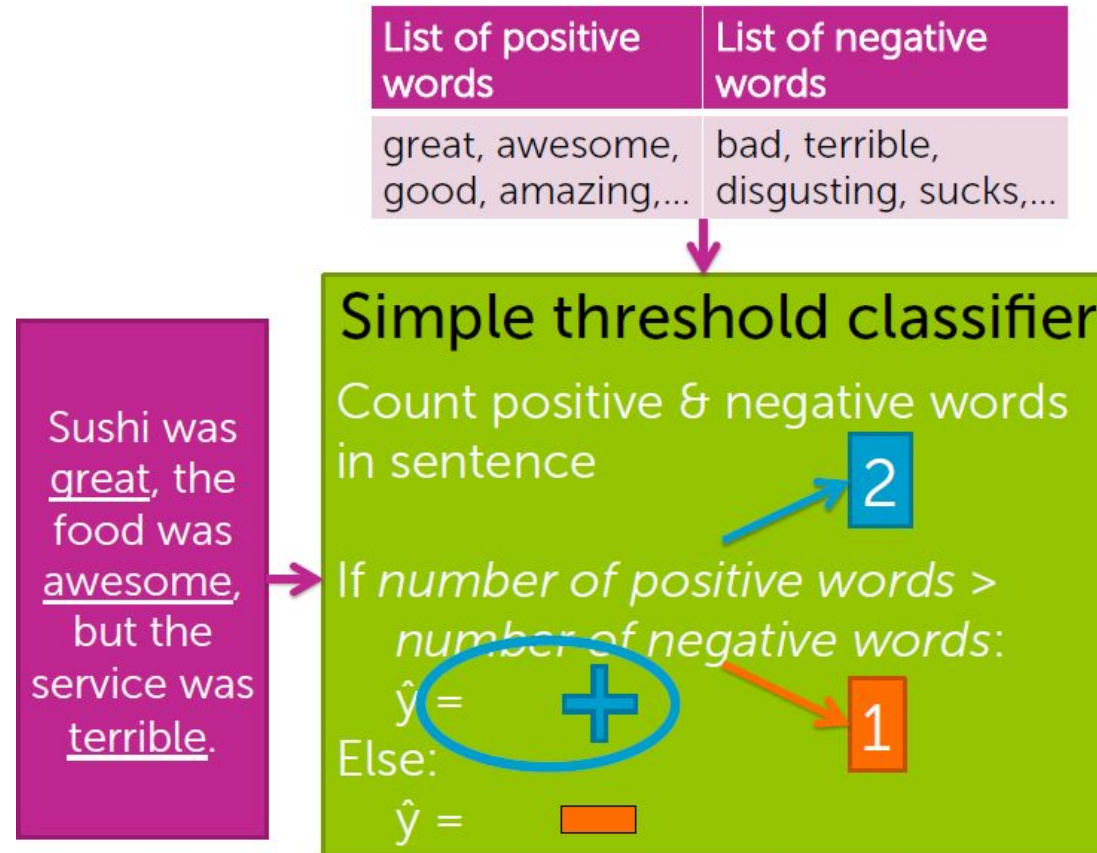
Intuition of Classifiers



Intuition of Classifiers



Intuition of Classifiers



Problem with threshold classifiers

- How do we get list of positive/negative words?
- Words have different degrees of sentiment:
 - Great > good
 - How do we weigh different words?
- Single words are not enough:
 - *Good* → Positive
 - *Not good* → Negative

Addressed
by learning
a classifier

Addressed
by more
elaborate
features



Linear Classifiers



Simple (linear) classifier

Will use training data to learn a weight for each word

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

Input x:

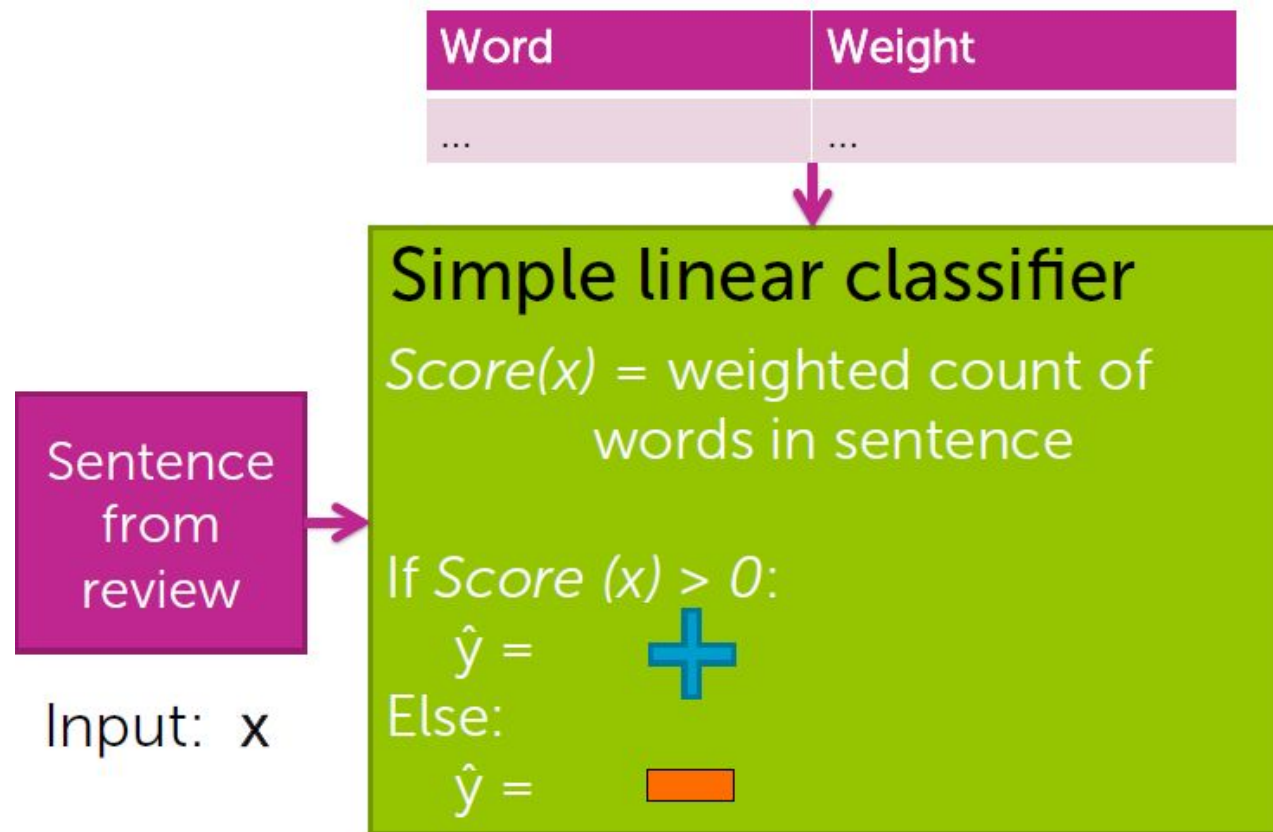
Sushi was great,
the food was awesome,
but the service was terrible.

Score(x) = ?

Linear classifier:
Output is weighted sum of
input



Scoring a sentence



Decision boundary

Suppose only 2 words had non-zero weights...

Word	Weight
awesome	1.0
awful	-1.5

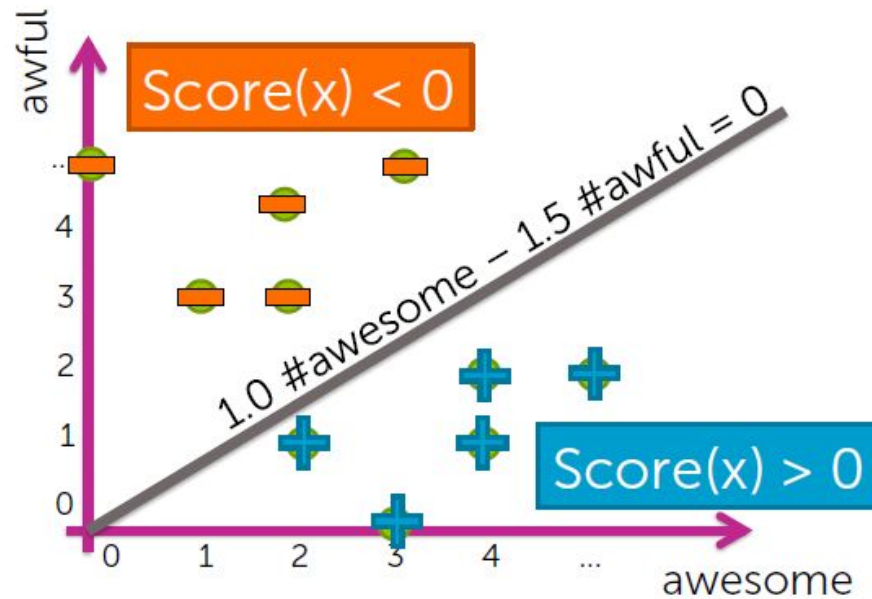
→ $\text{Score}(x) = 1.0 \# \text{awesome} - 1.5 \# \text{awful}$



Decision boundary

Word	Weight
awesome	1.0
awful	-1.5

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

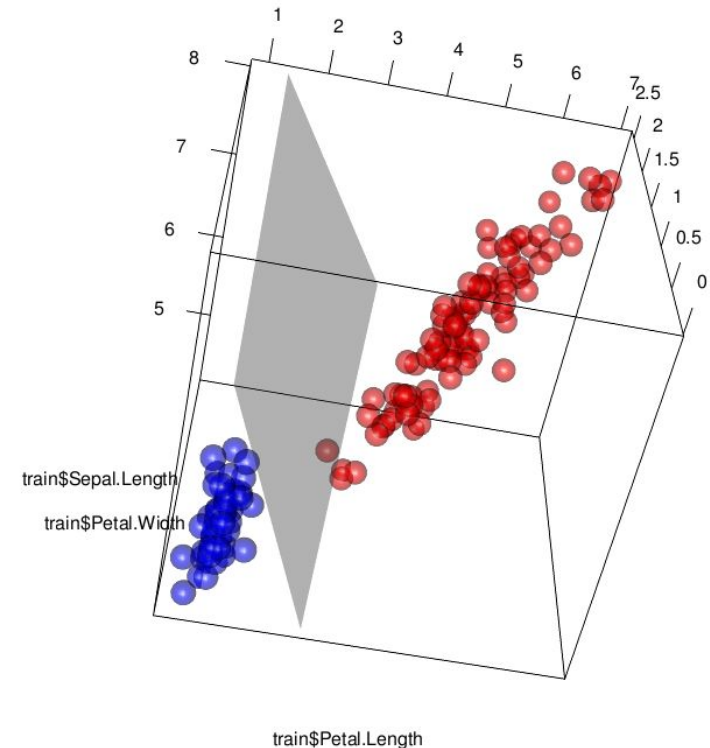


A boundary is to...separate

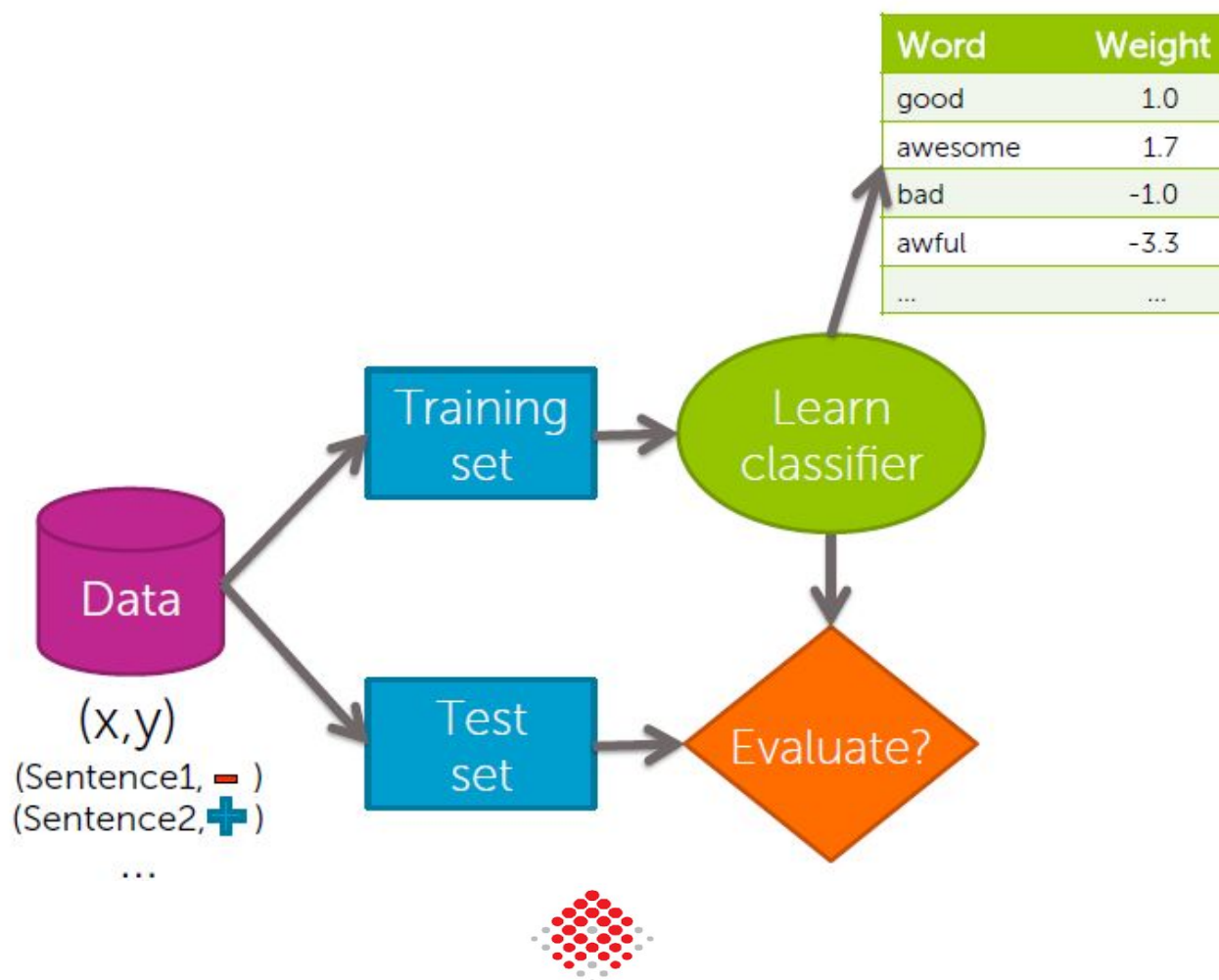
For linear classifiers:

- When 2 weights are non-zero
⇒ **line**
- When 3 weights are non-zero
⇒ **plane**
- When many weights are non-zero
⇒ **hyperplane**

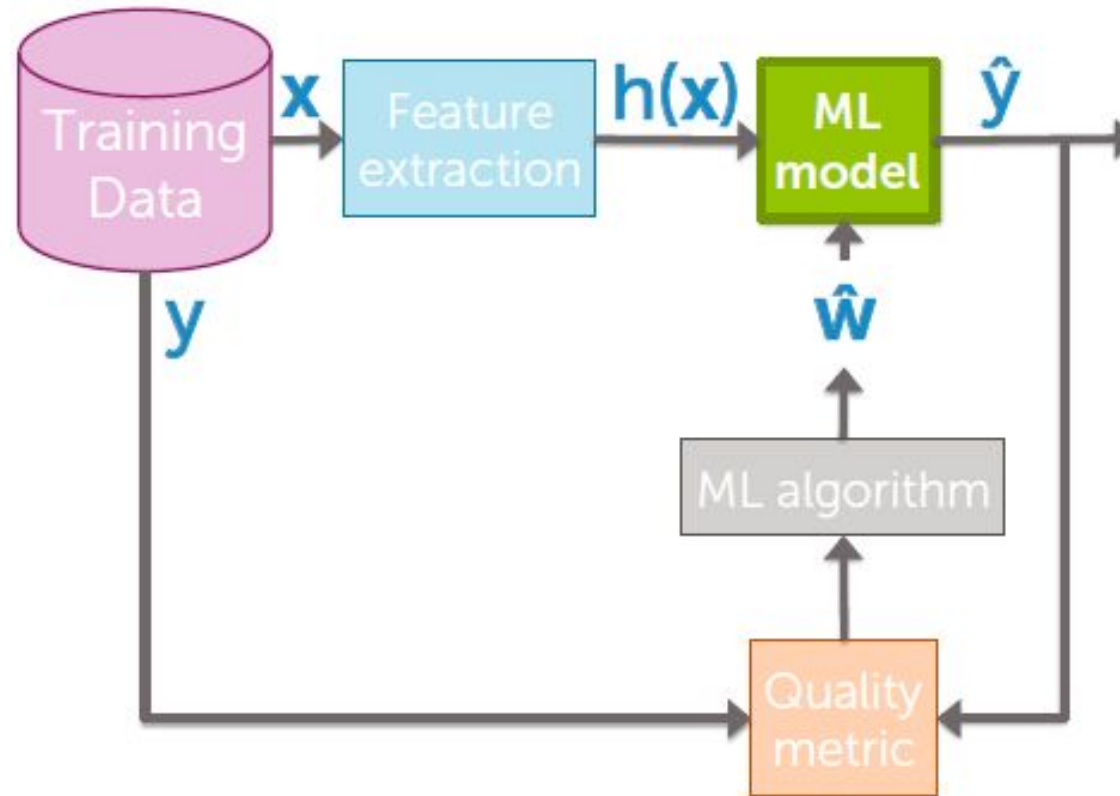
For more general classifiers
⇒ more complicated shapes



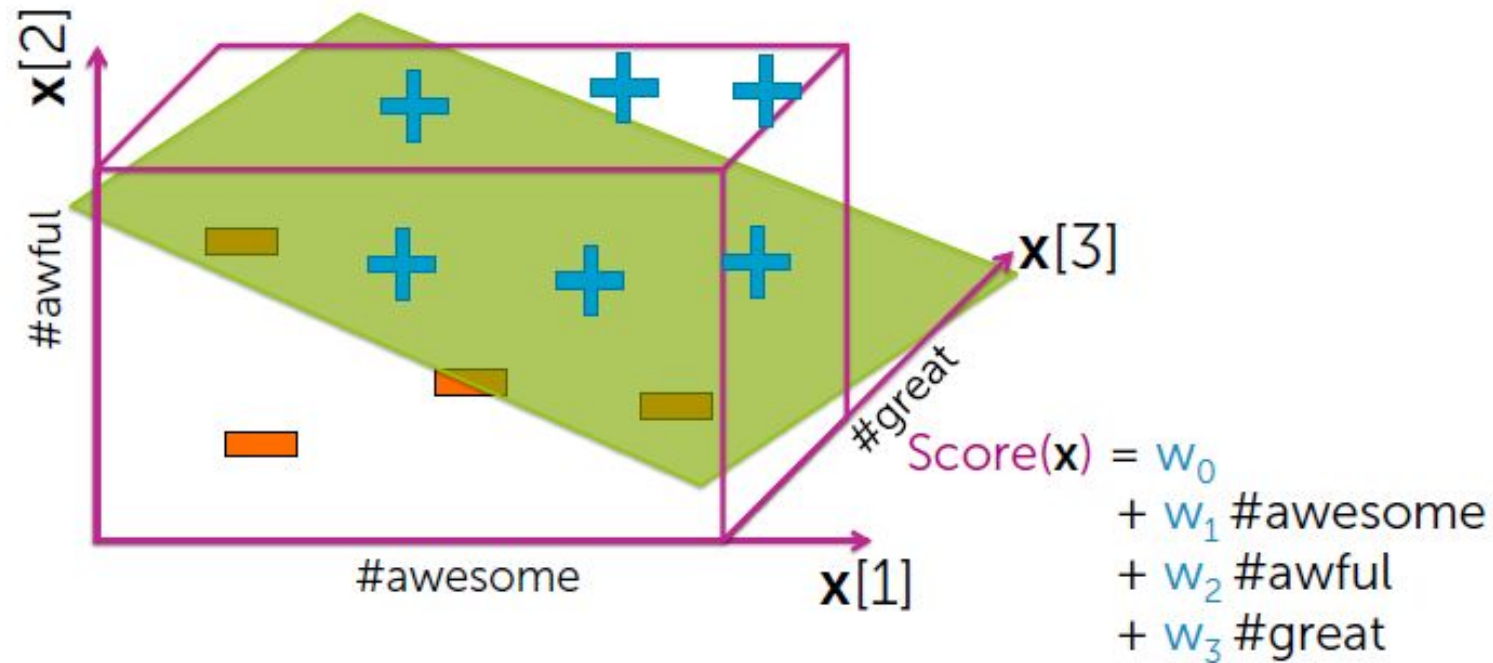
Training a classifier → Learning weights



Machine learning system



Coefficients of a classifier



General notation

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

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

\nwarrow d -dim vector

Notational conventions:

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

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

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

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



Simple hyperplane

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

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

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

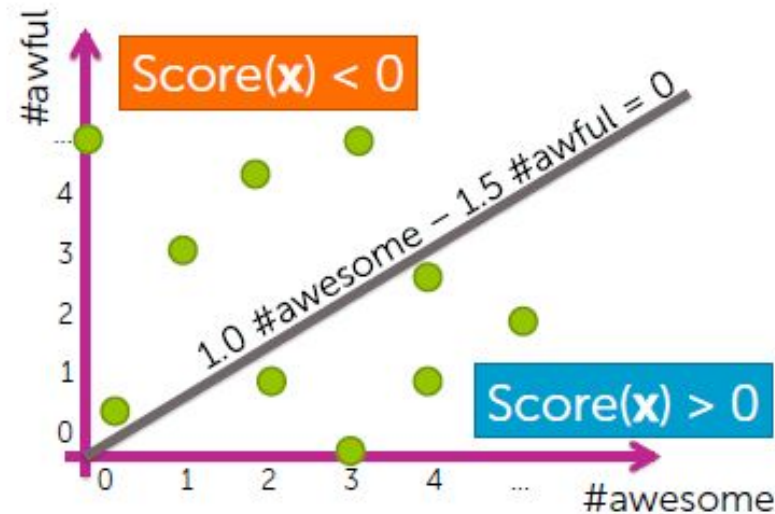
How did it
arrive at this?



Decision boundary: Effect of changing coefficients

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

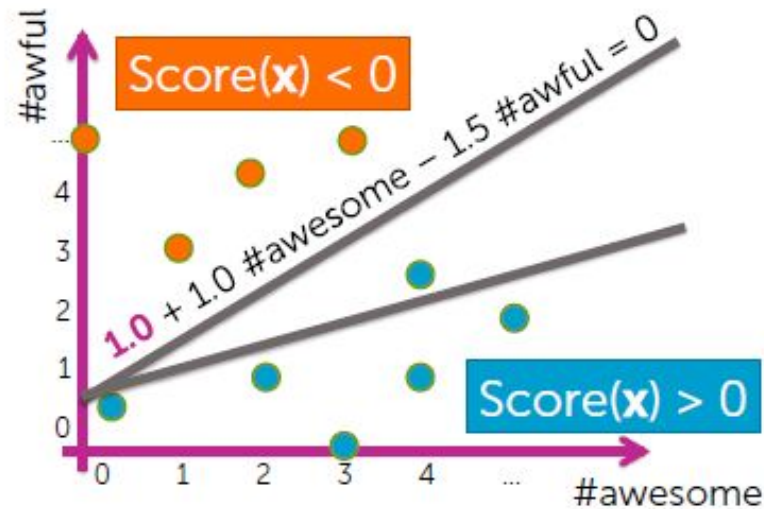
→ $\text{Score}(x) = 1.0 \text{ \#awesome} - 1.5 \text{ \#awful}$



Decision boundary: Effect of changing coefficients

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

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



More generic features...

D-dimensional hyperplane

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

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

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

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

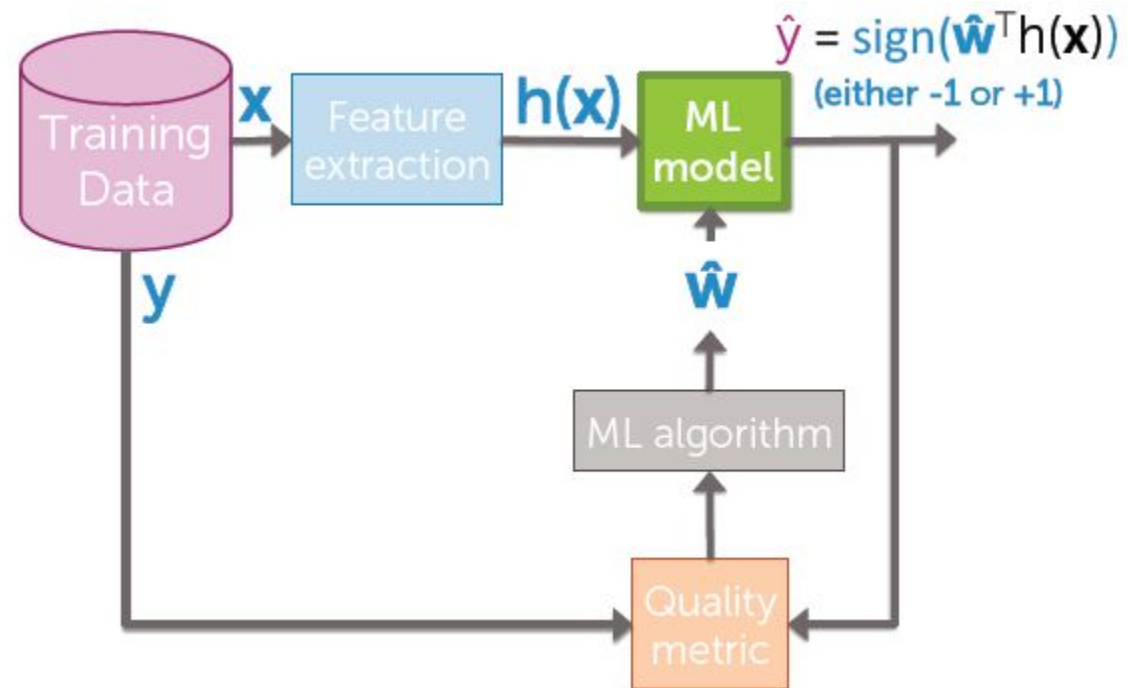
or, $\text{tf-idf}(\text{"awful"})$

...

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



ML model

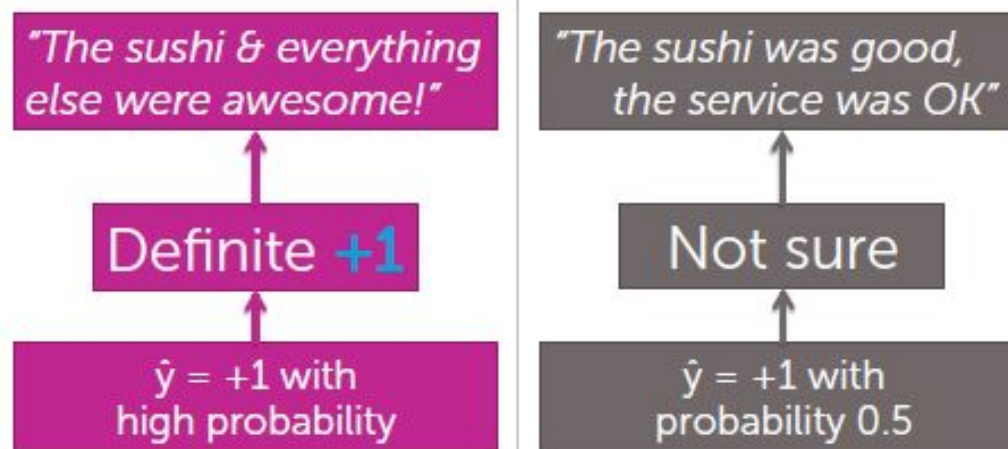


Class probability



How confident is your prediction?

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



An intuition on probability

Probability a review is positive is 0.7

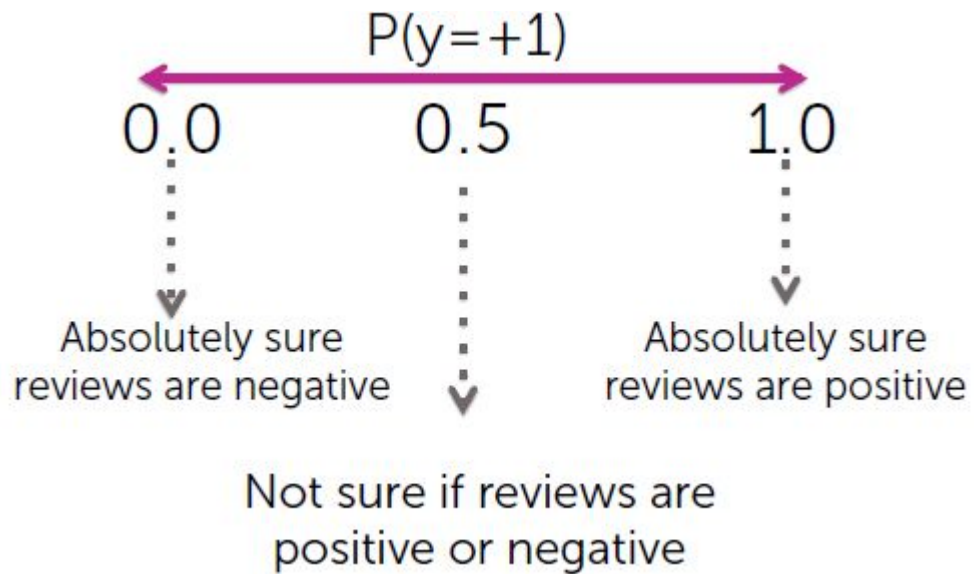


x = review text	y = sentiment
All the sushi was delicious! Easily best sushi in Seattle.	+1
The sushi & everything else were awesome!	+1
My wife tried their ramen, it was pretty forgettable.	-1
The sushi was good, the service was OK	+1
...	...

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



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		
Probabilities sum up to 1		



Conditional probability

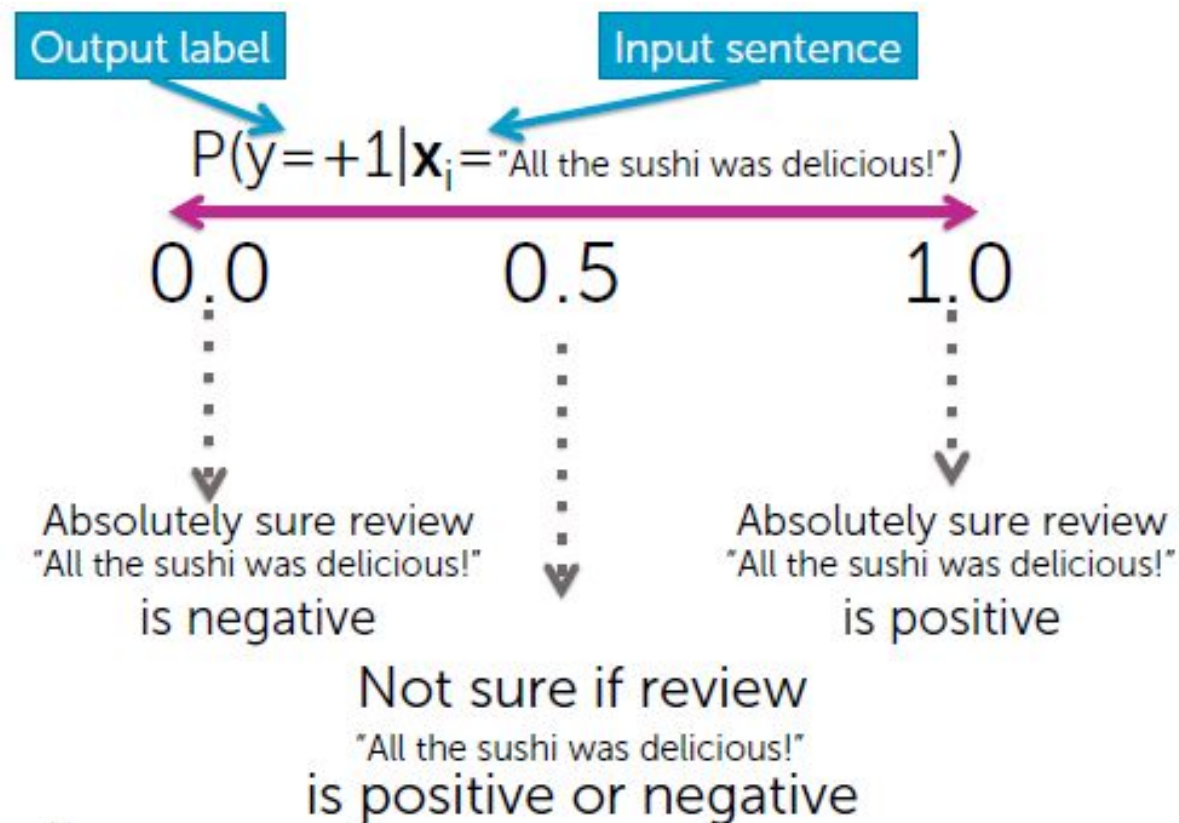
Probability a review with 3 "awesome" and 1 "awful" is positive is 0.9

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

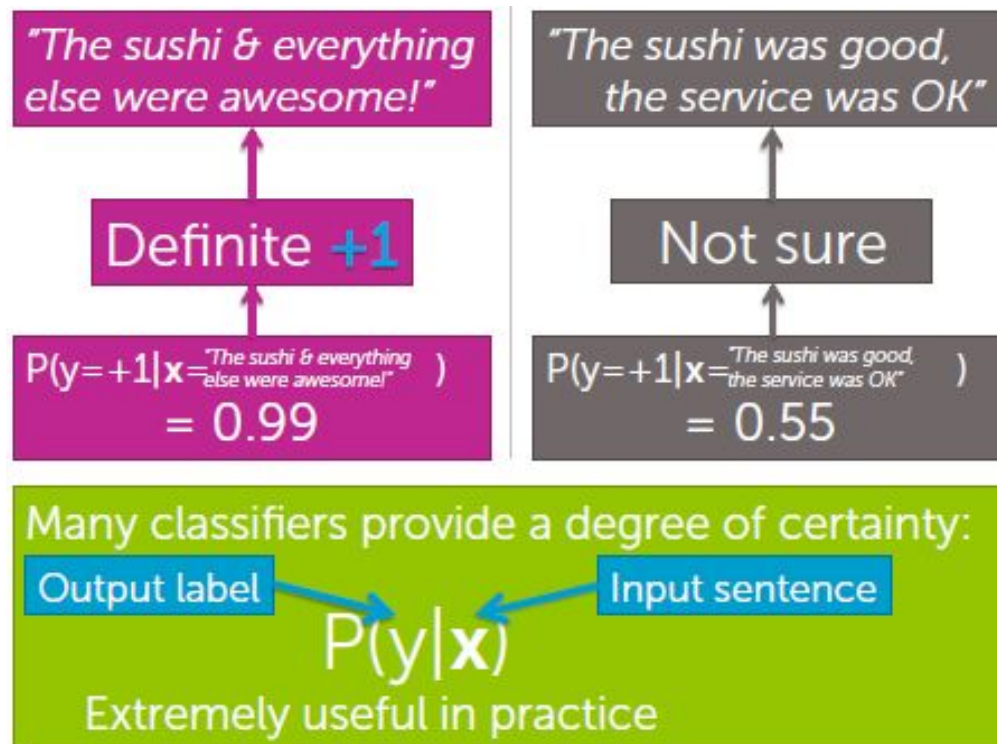
I expect 90% of rows with reviews containing 3 "awesome" & 1 "awful" to have y = +1 (Exact number will vary for each specific dataset)



Interpreting conditional probabilities



How confident is your prediction?



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	2	-1
3	3	-1
4	1	+1
...

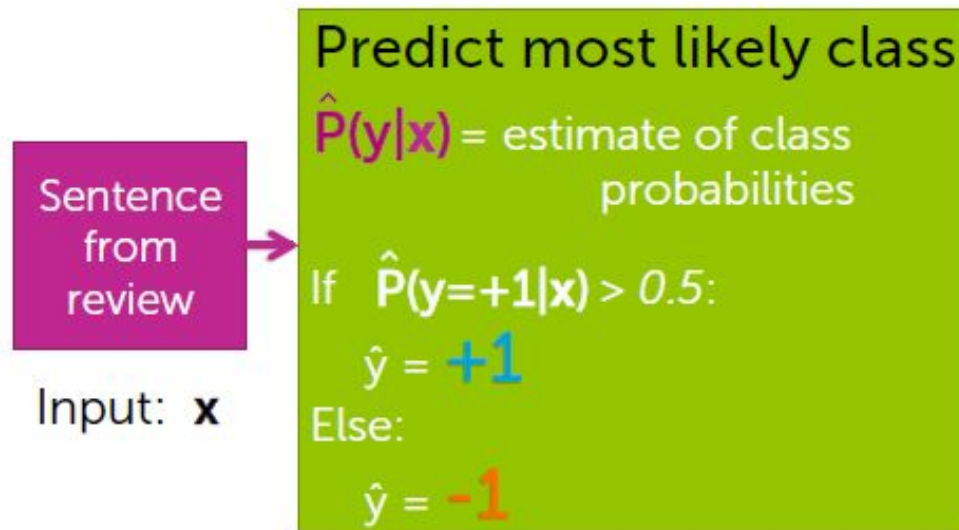
Optimize **quality metric**
on training data

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

Useful for
predicting \hat{y}



Estimate class probabilities

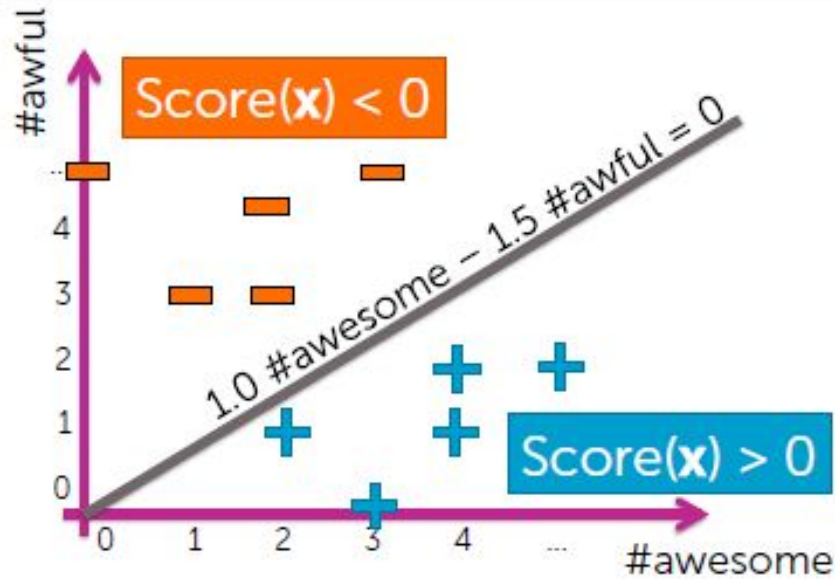


- Estimating $\hat{P}(\mathbf{y}|\mathbf{x})$ improves **interpretability**:
 - Predict $\hat{y} = +1$ **and** tell me how sure you are

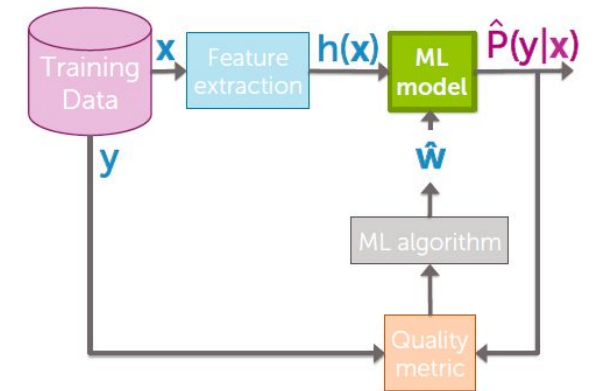


Revisit the “Score”

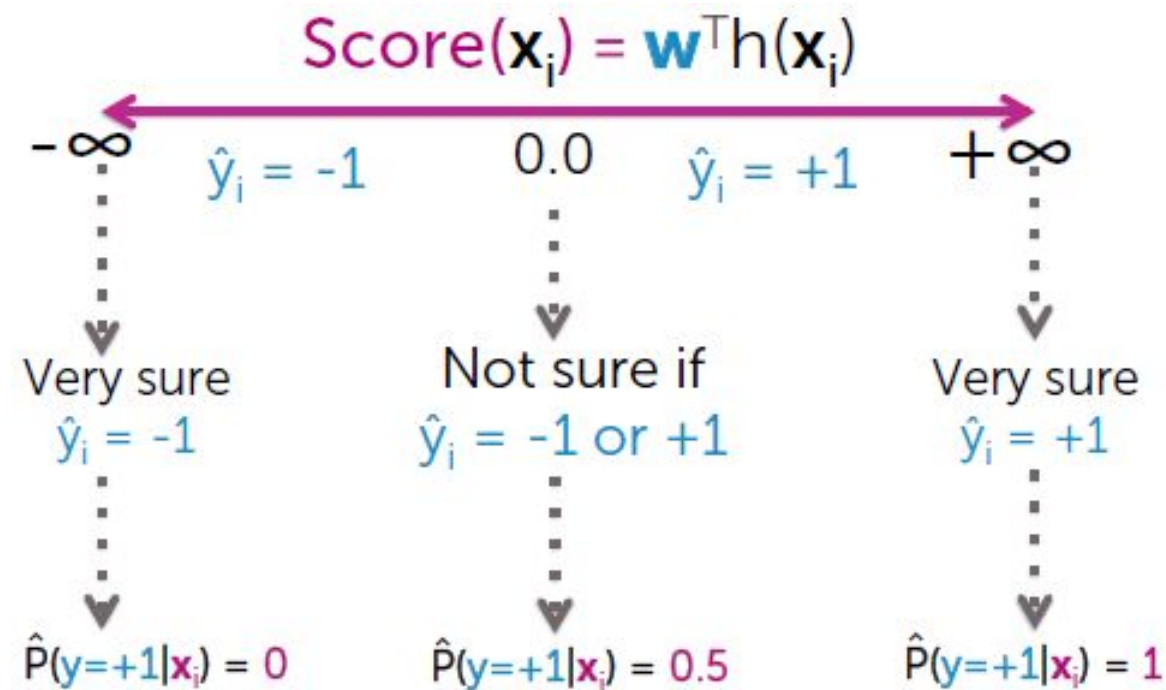
$$\begin{aligned}\text{Score}(\mathbf{x}_i) &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) \\ &= \mathbf{w}^T \mathbf{h}(\mathbf{x}_i)\end{aligned}$$



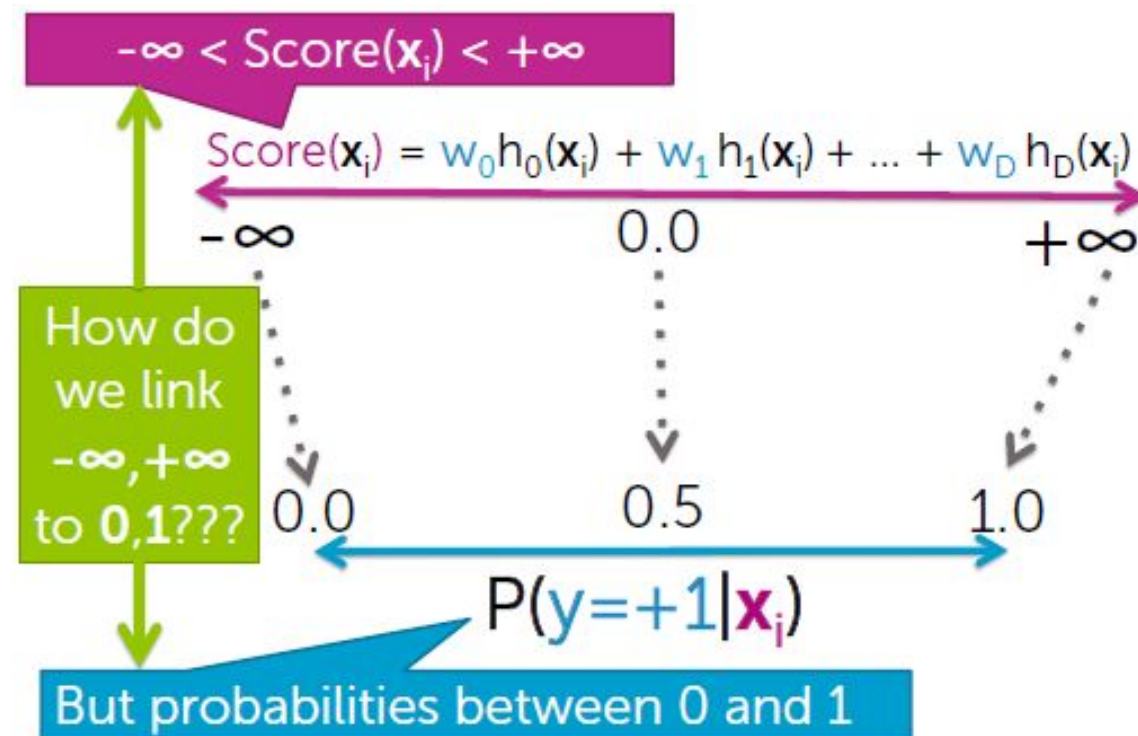
Relate
Score(\mathbf{x}_i) to
 $\hat{P}(y=+1|\mathbf{x}, \hat{\mathbf{w}})$?



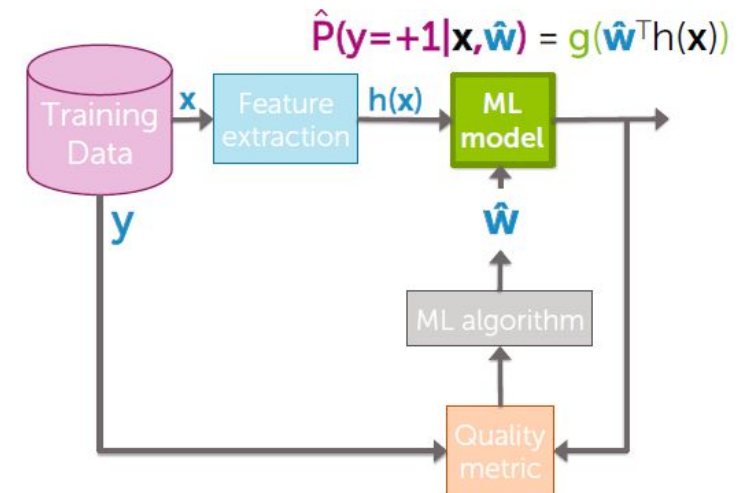
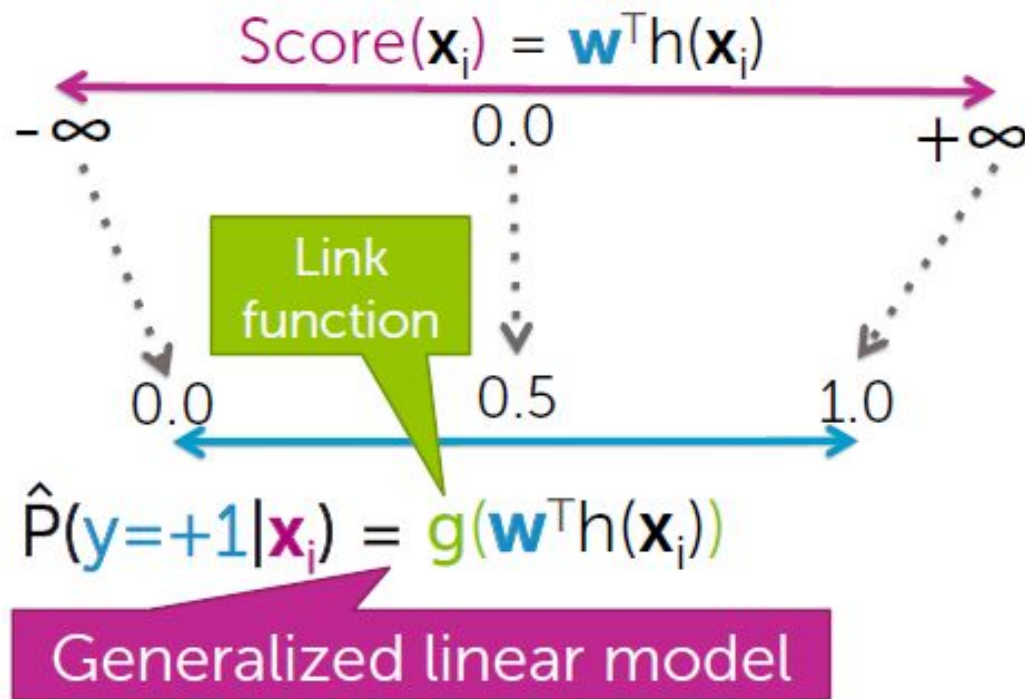
Interpreting “Score(x)”



Why not we use regression?



Link Function: Squeeze to $[0, 1]$



Logistic Regression Classifier

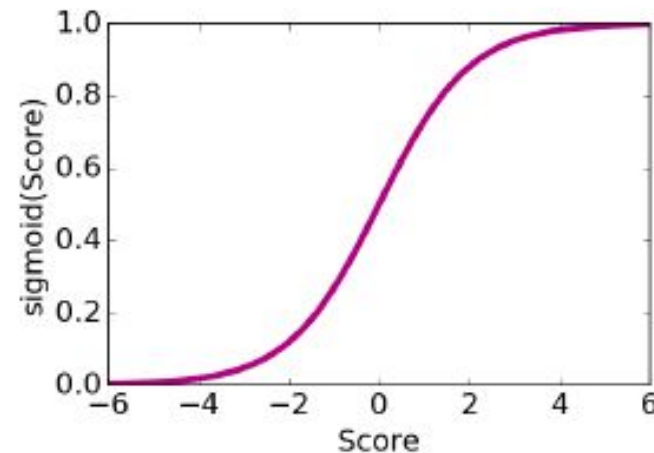
“Linear score with logistic link function”



Logistic function (sigmoid, logit)

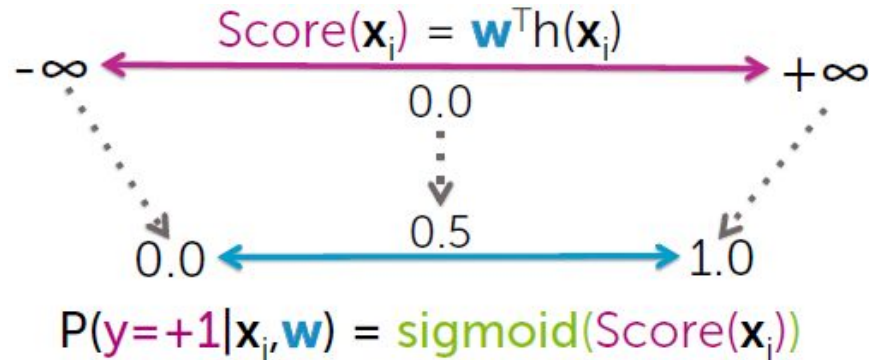
$$\text{sigmoid}(\text{Score}) = \frac{1}{1 + e^{-\text{Score}}}$$

Score	$-\infty$	-2	0.0	+2	$+\infty$
sigmoid(Score)					

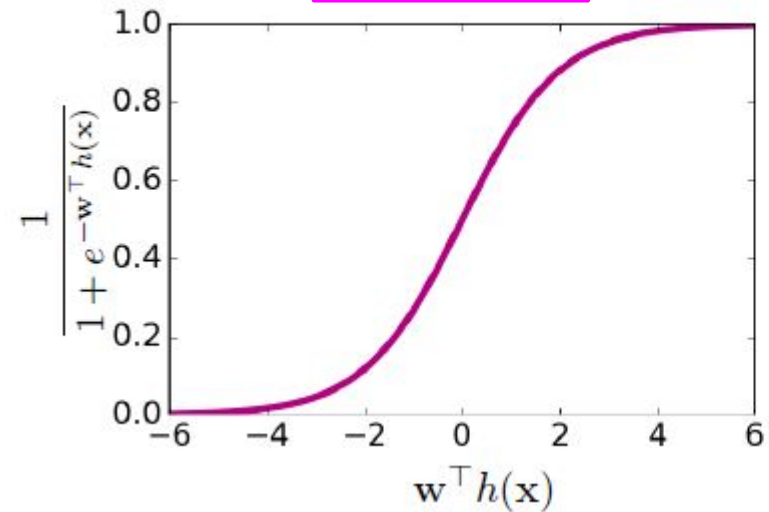


Logistic regression model

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



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



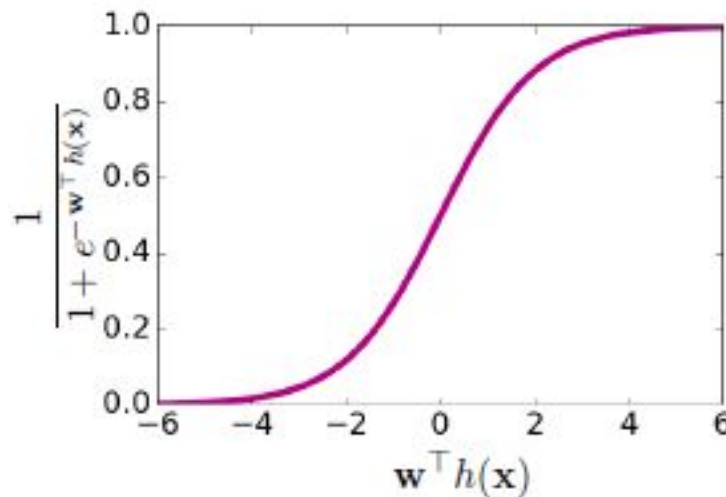
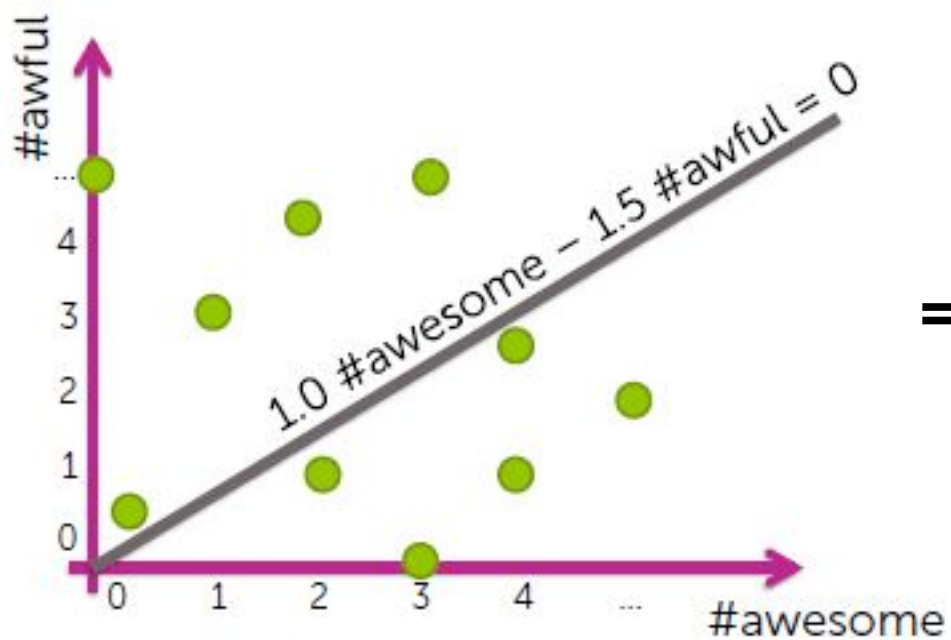
Score

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



The two steps

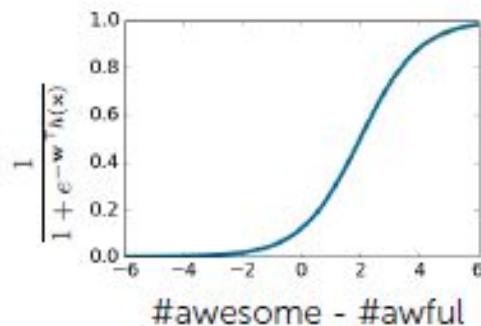
Linear decision boundary \Rightarrow Logistic regression model



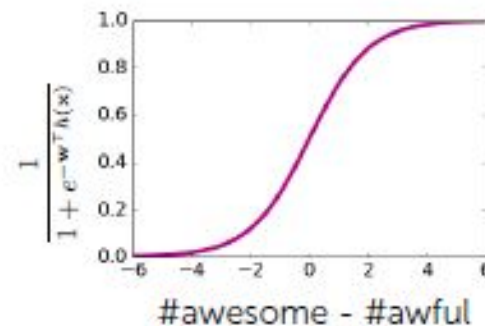
Effects of coefficients on logistic regression model

Linear decision boundary \Rightarrow Logistic regression model

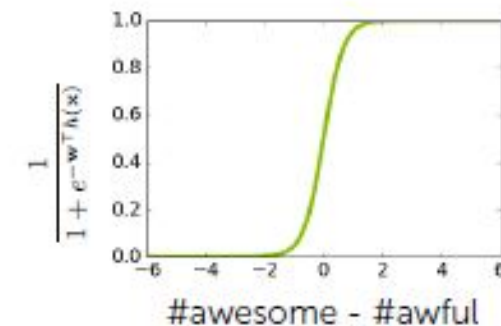
w_0	-2
$w_{\text{\#awesome}}$	+1
$w_{\text{\#awful}}$	-1



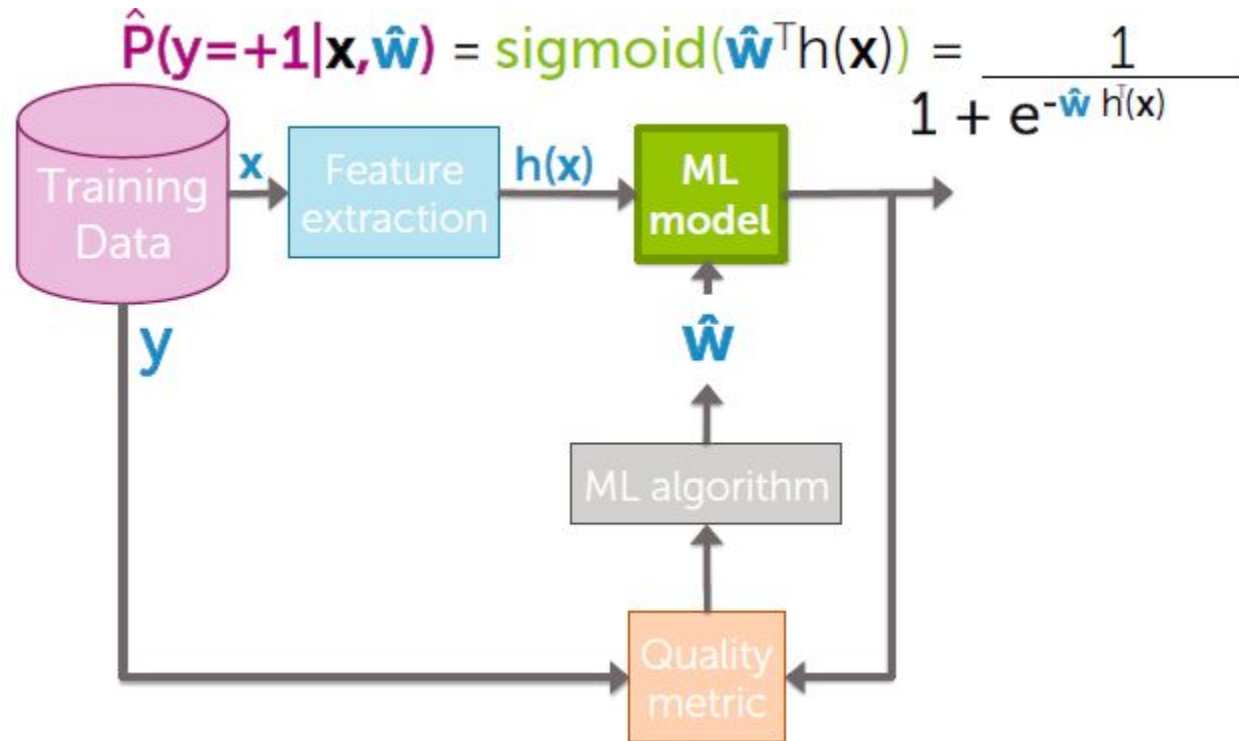
w_0	0
$w_{\text{\#awesome}}$	+1
$w_{\text{\#awful}}$	-1



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



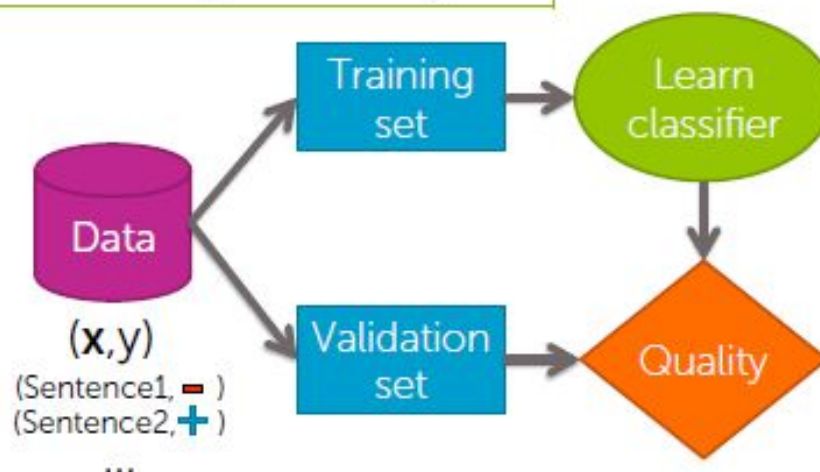
Logistic regression model



Training a classifier → Learning coefficients

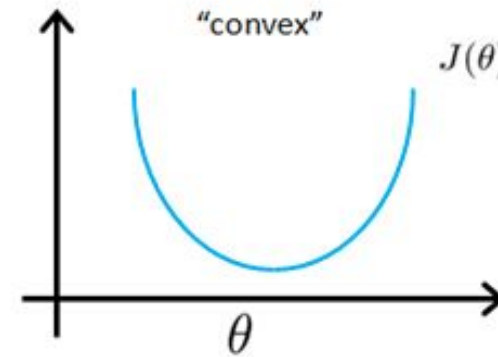
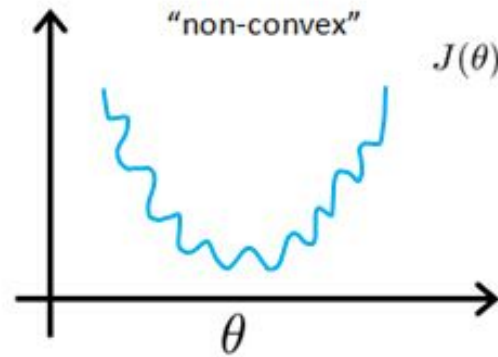
Word	Coefficient	Value
	\tilde{w}_0	-2.0
good	\tilde{w}_1	1.0
awesome	\tilde{w}_2	1.7
bad	\tilde{w}_3	-1.0
awful	\tilde{w}_4	-3.3
...

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



Find “best” classifier

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



This nonlinear function now causes the cost to be **non-convex**

That means it can easily get stuck in many local minimas

Recall \rightarrow $\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{h}(\mathbf{x}_i)^T \mathbf{w})^2$

which is the sum of cost of all samples for linear regression




Find “best” classifier

Logistic function have to be redefined so that it is **convex**

Let \hat{y} or $H_w(x)$ be the function of the decision boundary,
the convex “version” of the logistic function is approximated as

$$J = \text{Cost}(H_w(x), y) = -y \log(H_w(x)) - (1-y) \log(1-H_w(x))$$

Next, compute the gradient of the cost function, which is part of the update procedure:

$$w^{t+1} := w^t - \underbrace{\sum (H_w(x) - y) x}_{\Delta J}$$


Encoding categorical inputs

Numerical inputs:

- #awesome, age, salary, ...
- Intuitive when multiplied by coefficient, e.g. 1.5 #awesome

Categorical inputs:

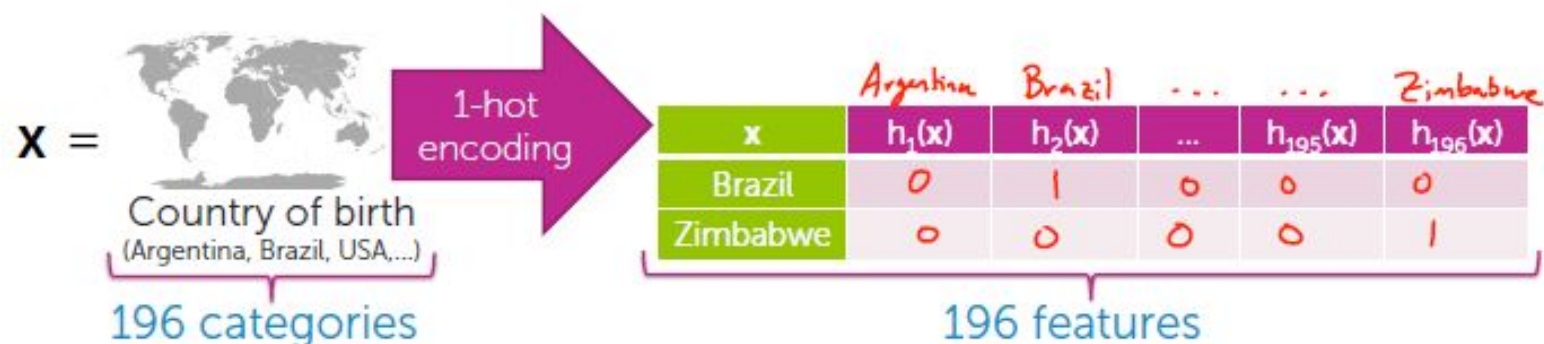


How do we multiply category by coefficient???

Must convert categorical inputs into numeric features



Categories to numerics

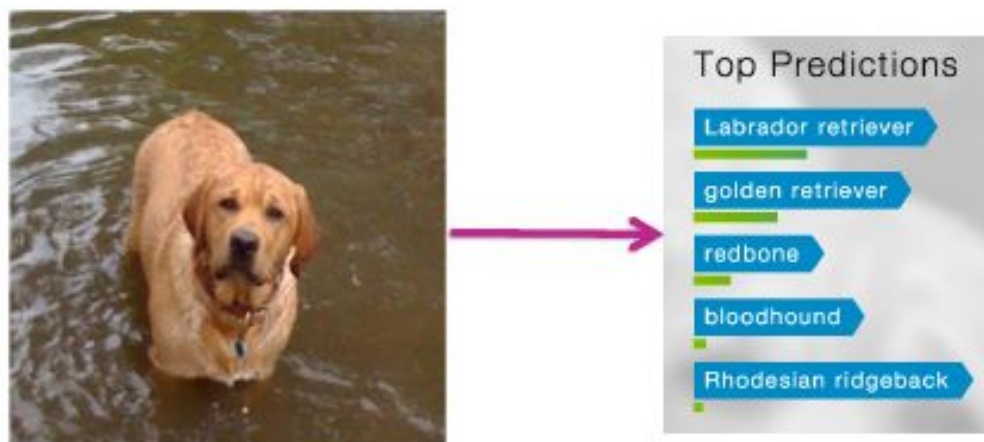


Multi-class Classification

Using 1-versus-All



Multiclass classification



Input: x
Image pixels

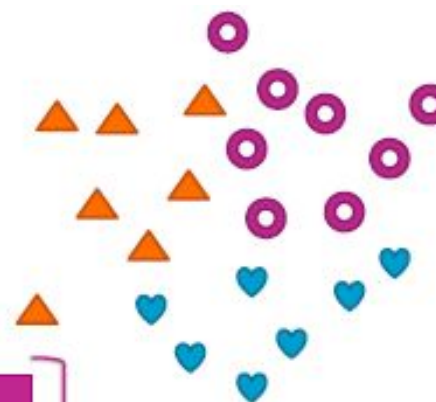
Output: y
Object in image



Multiclass classification formulation

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

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

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

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



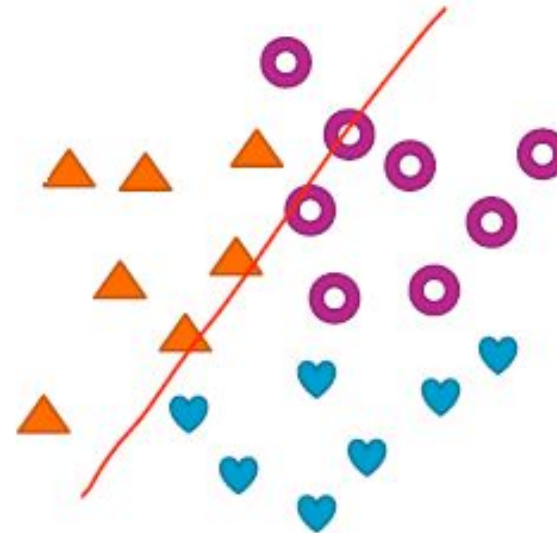
1-versus-all

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

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

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

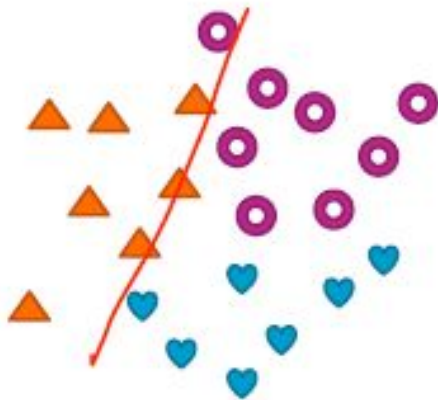
Predict: $\hat{P}(y=\triangle | \mathbf{x}_i) = \hat{P}_{\triangle}(y=+1 | \mathbf{x}_i)$



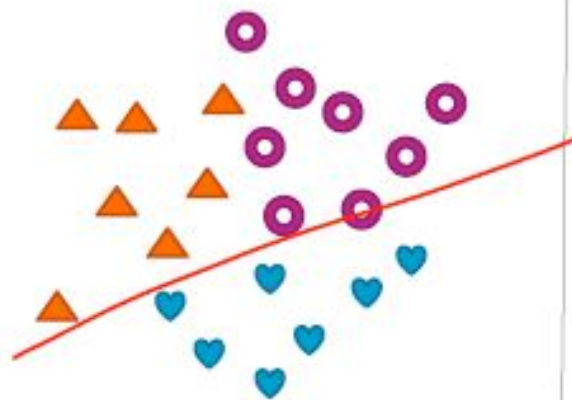
1-versus-all

Simple multiclass classification using C 2-class models

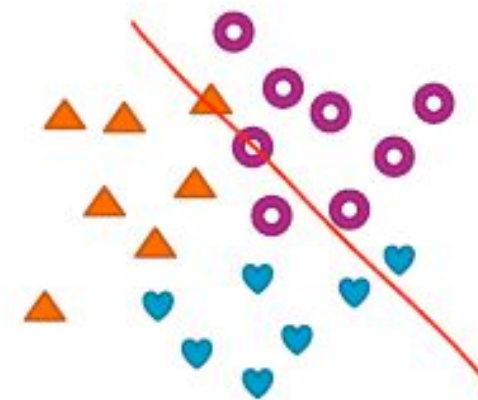
$$\hat{P}(y=\triangle | \mathbf{x}_i) =$$



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



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



Multiclass training



Input: \mathbf{x}_i

Multiclass training

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

Predict most likely class

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

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

If $\hat{P}_c(y=+1|\mathbf{x}_i) > \text{max_prob}$:

$\hat{y} = c$

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



Evaluating Classifiers



Classification error & Accuracy

Error measures fraction of mistakes made in classification:

$$\text{error} = \frac{\# \text{ of mistakes}}{\text{Total \# of sentences}}$$

(Best possible value: 0.0)

Accuracy measures the fraction of correct predictions:

$$\text{accuracy} = \frac{\# \text{ of correct}}{\text{Total \# of sentences}}$$

(Best possible value: 1.0)



What's a good accuracy?

- For binary classification:
 - Half the time, you'll get it right! (on average)
→ accuracy = 0.5
- For k classes, accuracy = $1/k$
 - 0.333 for 3 classes, 0.25 for 4 classes,...

At the very, very, very least,
you should healthily beat random...
Otherwise, it's (usually) pointless...



Is 90% accuracy good? Depends...

2010 data shows:
"90% emails sent are spam!"

Predicting every email is spam
gets you 90% accuracy!!!

Majority class prediction

Amazing performance when
there is class imbalance

(but silly approach)

- One class is more common than others
- Beats random (if you know the majority class)







Ask the hard questions...

- Is there class imbalance?
- How does it compare to a simple, baseline approach?
 - Random guessing
 - Majority class
 - ...
- Most importantly:
what accuracy does my application need?
 - What is good enough for my user's experience?
 - What is the impact of the mistakes we make?



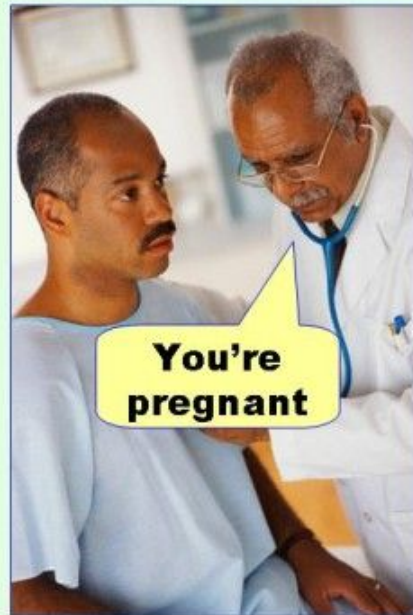
Types of Errors

		Predicted label	
			
True label		True Positive	False Negative (FN)
		False Positive (FP)	True Negative

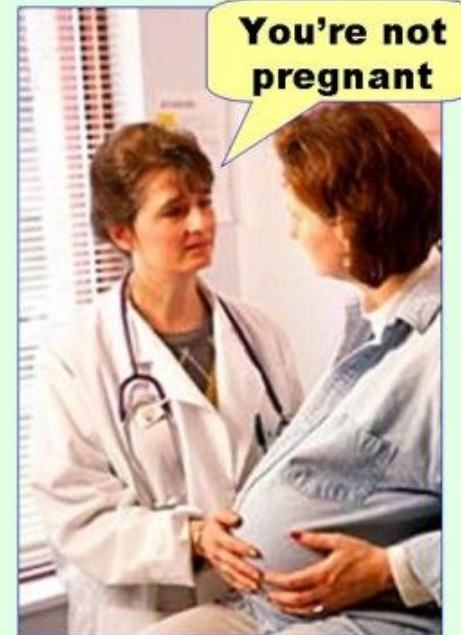


Types of Errors

Type I error
(false positive)



Type II error
(false negative)



Cost of different types of mistakes





Different (and high) in some applications!

	Spam filtering	Medical diagnosis
False negative	Annoying	Disease not treated
False positive	Email lost	Wasteful treatment



Confusion matrix – binary

Test set: 100 samples

		Predicted label	
			
True label		50	10
		5	35

Accuracy = ?



Confusion matrix – multiclass

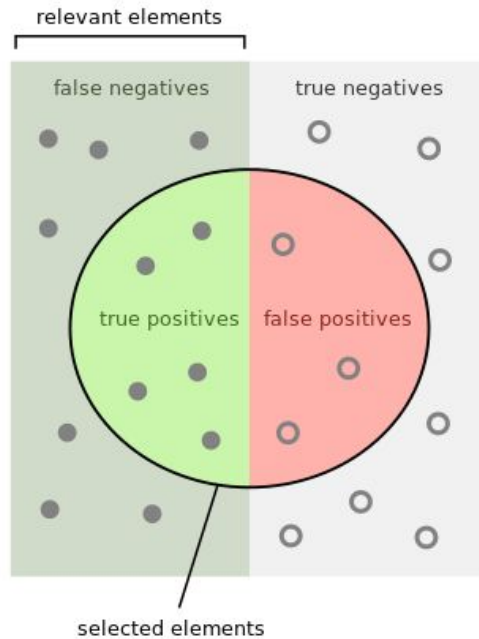
Test set: 100 samples

		Predicted label		
		Healthy	Cold	Flu
True label	Healthy	60	8	2
	Cold	4	12	4
	Flu	0	4	8

Accuracy = ?



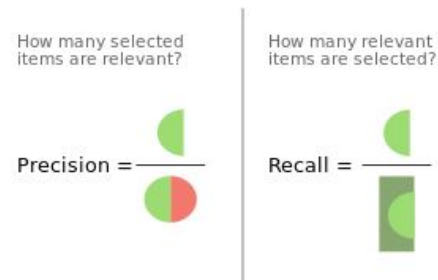
Precision, Recall



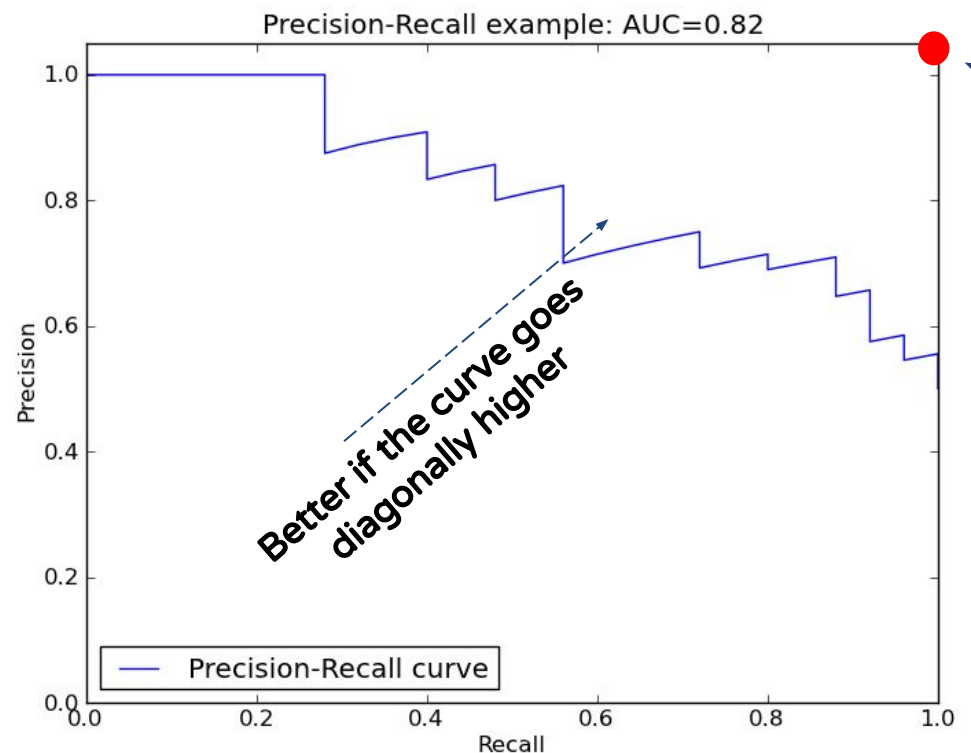
		Predicted label	
		+	-
True label	+	True Positive	False Negative (FN)
	-	False Positive (FP)	True Negative

$$\text{Recall} = \frac{TP}{TP + FN}$$

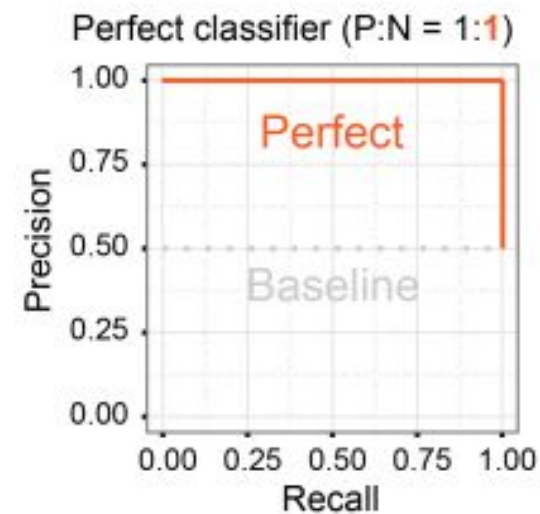
$$\text{Precision} = \frac{TP}{TP + FP}$$



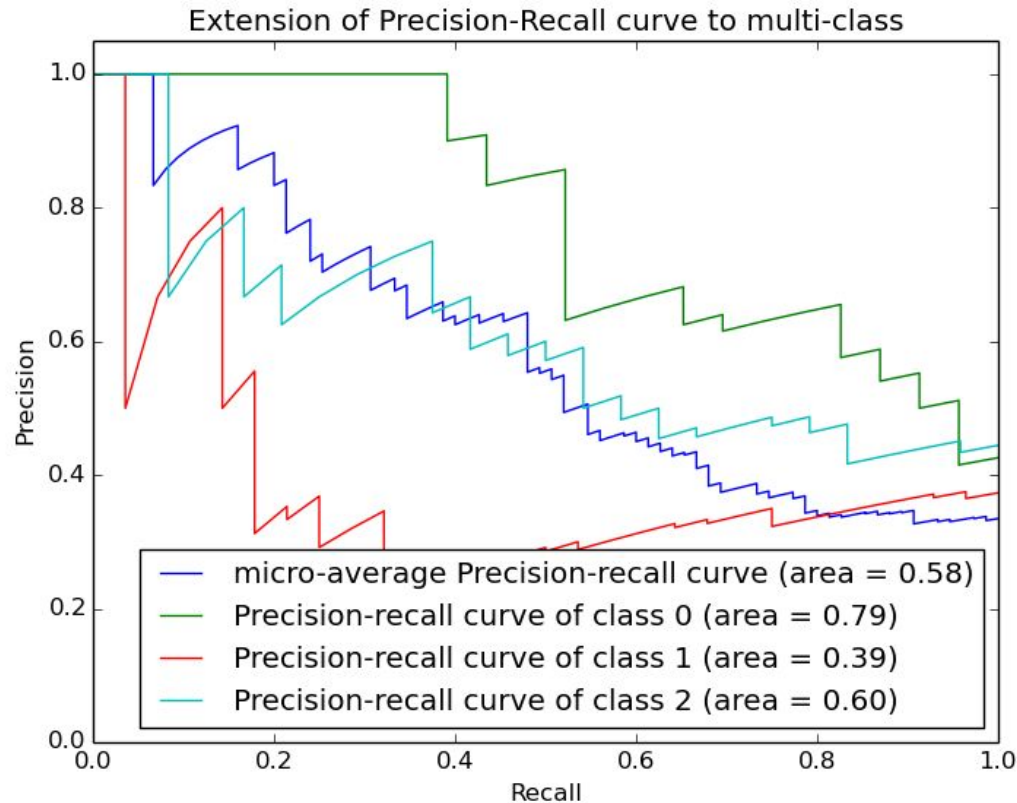
Precision-Recall Curve



This is where you want your performance to be!



Precision-Recall Curve



Can be extended to multiclass case

⇒ You get separate P-R curves for each class!



Go To Exercises

Learning Curves

How much data do I need?

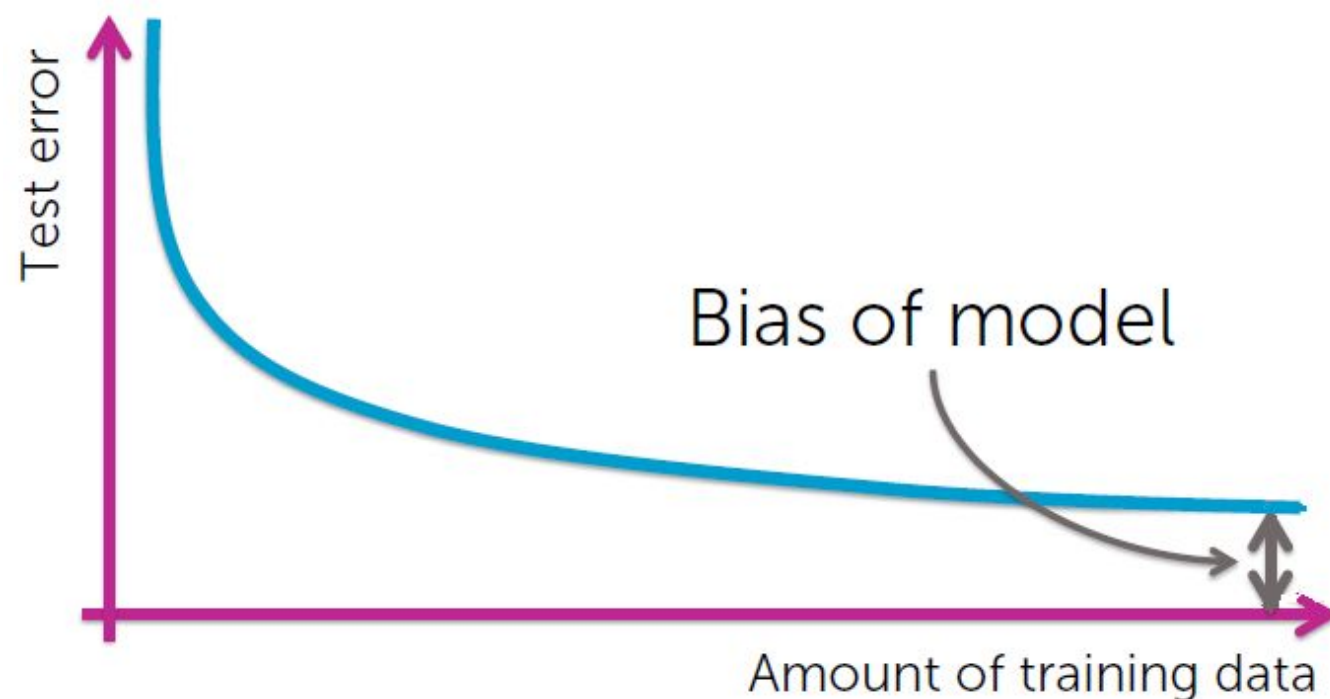


How much data needed for learning?

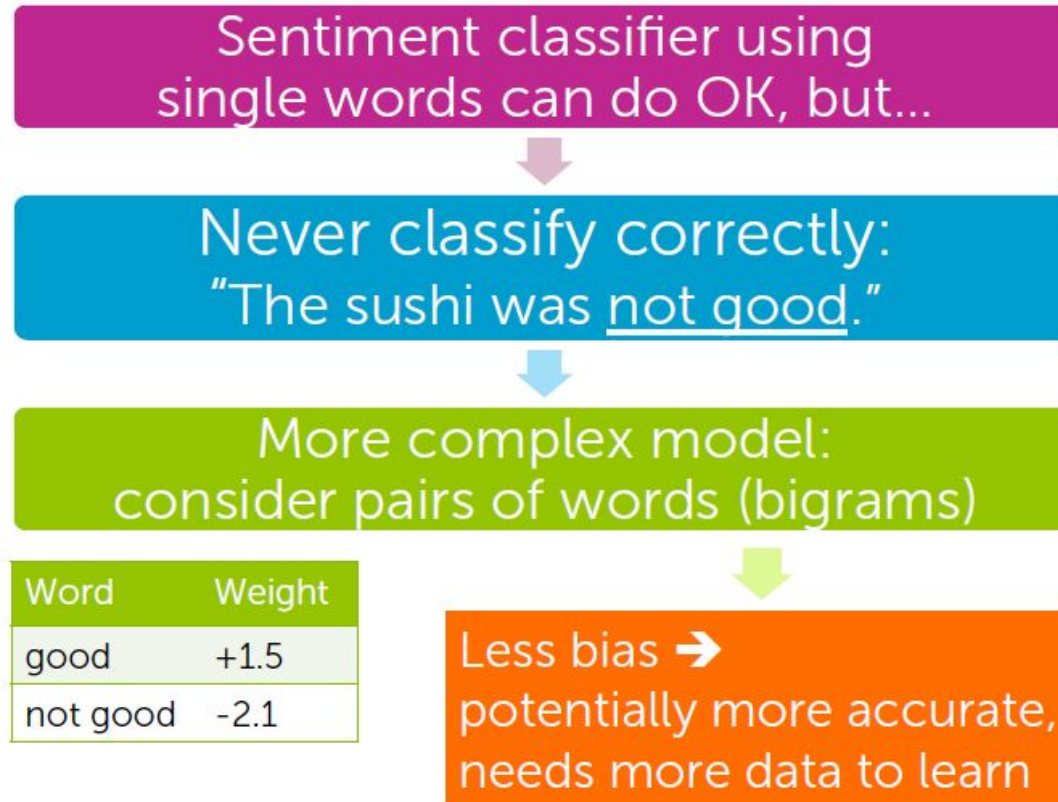
- The more the merrier 😊
 - But data quality is most important factor
- Theoretical techniques sometimes can bound how much data is needed
 - Typically too loose for practical application
 - But provide guidance
- In practice:
 - More complex models require more data
 - Empirical analysis can provide guidance



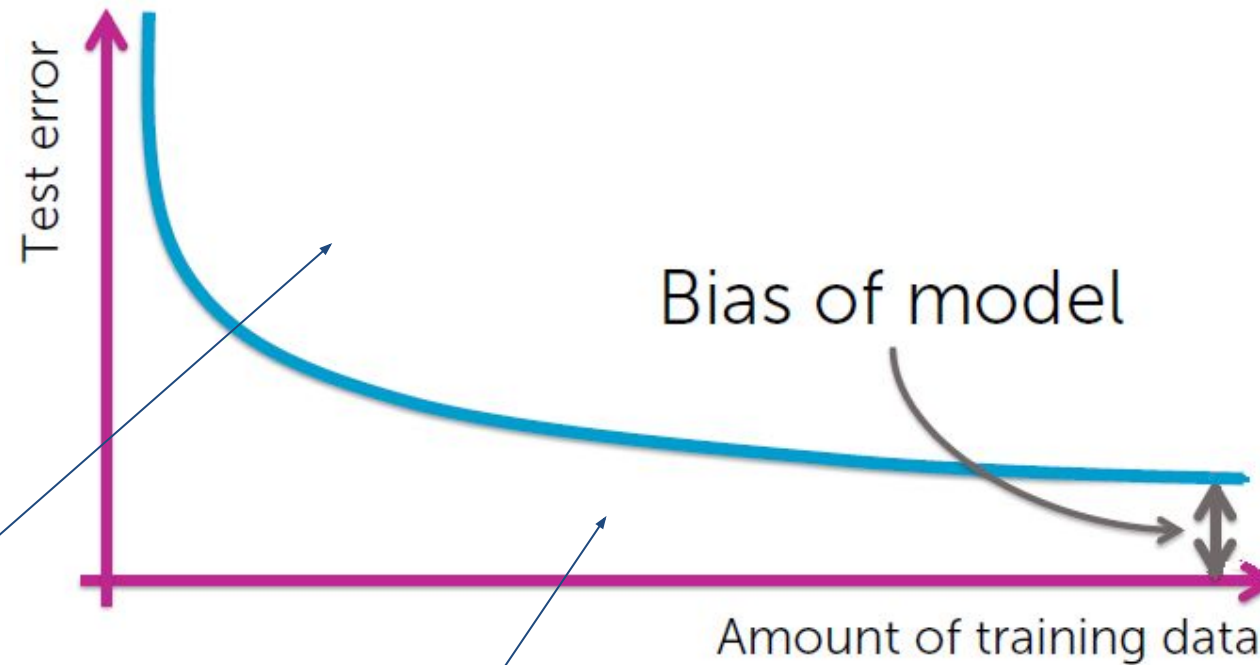
Is there a limit? Yes, for most models



More complex models have less bias



Models with less bias



**Need more data to learn well
But with sufficient data, bias is greatly reduced!**



Try Kaggle

Breast Cancer Wisconsin (Diagnostic) Data
Set

<https://www.kaggle.com/uciml/breast-cancer-wisconsin-data>

