Linear Models

Machine Learning



• Supervised learning: instances, concepts, and hypotheses

• Supervised learning: instances, concepts, and hypotheses



- Supervised learning: instances, concepts, and hypotheses
- Specific learners
 - Decision trees



• Supervised learning: instances, concepts, and hypotheses



- General ML ideas
 - Features as high dimensional vectors
 - Overfitting

Lecture outline

- Linear models
- What functions do linear classifiers express?

Where are we?

- Linear models
 - Introduction: Why linear classifiers and regressors?
 - Geometry of linear classifiers
 - A notational simplification
 - Learning linear classifiers: The lay of the land
- What functions do linear classifiers express?

Is learning possible at all?

- There are 2¹⁶ = 65536 possible
 Boolean functions over 4 inputs
 - Why? There are 16 possible outputs. Each way to fill these 16 slots is a different function, giving 2¹⁶ functions.
- We have seen only 7 outputs
- We *cannot* know what the rest are without seeing them
 - Think of an adversary filling in the labels every time you make a guess at the function

x_1	x_2	x_3	x_4	$\mid y$
0	0	0	0	?
0	0	0	1	?
0	0	1	0	$\rightarrow 0$
0	0	1	1	$1 \leftarrow$
0	1	0	0	$\rightarrow 0$
0	1	0	1	$\rightarrow 0$
0	1	1	0	$\rightarrow 0$
0	1	1	1	?
1	0	0	0	?
1	0	0	1	$1 \leftarrow$
1	0	1	0	?
1	0	1	1	?
1	1	0	0	$\rightarrow 0$
1	1	0	1	?
1	1	1	0	?
1	1	1	1	?

Is learning possible at all?

	 There are 2¹⁶ = 65536 possible Boolean functions over 4 inputs Why? There are 16 possible outputs. Each way to fill these 16 slots is a different function, giving 2¹⁶ functions. 	$egin{array}{c} x_1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$	$egin{array}{c} x_2 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \end{array}$	$\begin{array}{c} x_3 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{array}$	$\begin{array}{c} x_4 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ \end{array}$	$\begin{array}{c} y \\ \hline ? \\ 0 \leftarrow \\ 1 \leftarrow \\ 0 \leftarrow \end{array}$
	How could we possibly learn	any	/thi	ng?	$1 \\ 0$	$egin{array}{c} ightarrow 0 \ ightarrow 0 \ ightarrow 0 \ ightarrow 0 \ ightarrow 0$
•	We have seen only / outputs	0	1	1	1	?
			0	0	0	?
		1	0	0	1	$1 \leftarrow$
•	We <i>cannot</i> know what the rest are	1	0	1	0	?
	without seeing them	1	0	1	1	?
— Th ev fu	Think of an advorsary filling in the labels	1	1	0	0	$0 \leftarrow$
	- Think of all adversary filling in the labels	1	1	0	1	?
	every time you make a guess at the	1	1	1	0	?
	TUNCTION	1	1	1	1	?

Solution: Restrict the search space

- A *hypothesis space* is the set of possible functions we consider
 - We were looking at the space of all Boolean functions
 - Instead choose a hypothesis space that is smaller than the space of all functions
 - Only simple conjunctions (with four variables, there are only 16 conjunctions without negations)
 - Simple disjunctions
 - <u>m-of-n rules: Fix a set of n variables.</u> At least m of them must be true
 - Linear functions











Similar argument for regression



Linear regression might make smaller errors on new points

Similar argument for regression



Linear regression might make smaller errors on new points

Similar argument for regression



Linear regression might make smaller errors on new points

Recall: Regression vs. Classification

- Linear regression is about predicting real valued outputs
- Linear classification is about predicting a discrete class label
 - +1 or -1
 - SPAM or NOT-SPAM
 - Or more than two categories

Suppose we want to determine whether a robot arm is defective or not using two measurements:

- 1. The maximum distance the arm can reach *d*
- 2. The maximum angle it can rotate *a*

Suppose we want to determine whether a robot arm is defective or not using two measurements:

- 1. The maximum distance the arm can reach *d*
- 2. The maximum angle it can rotate *a*

Suppose we use a linear decision rule that predicts defective if $2d + 0.01a \ge 7$

Suppose we want to determine whether a robot arm is defective or not using two measurements:

- 1. The maximum distance the arm can reach *d*
- 2. The maximum angle it can rotate *a*

Suppose we use a linear decision rule that predicts defective if $2d + 0.01a \ge 7$

We can apply this rule if we have the two measurements For example: for a certain arm, if d = 3 and a = 200, then $2d + 0.01a = 8 \ge 7$ The arm would be labeled as not defective

Suppose we want to determine whether a robot arm is defective or not using two measurements:

- 1. The maximum distance the arm can reach *d*
- 2. The maximum angle it can rotate *a*

Suppose we use a linear decision rule that predicts defective if $2d + 0.01a \ge 7$

This rule is an example of a linear classifier Features are weighted and added up, the sum is checked against a threshold

Linear Classifiers

Inputs are d dimensional feature vectors, denoted by **x** Output is a label $y \in \{-1, 1\}$

Linear Classifiers

Inputs are d dimensional feature vectors, denoted by **x** Output is a label $y \in \{-1, 1\}$

Linear Threshold Units classify an example \mathbf{x} using parameters \mathbf{w} (a d dimensional vector) and b (a real number) according the following classification rule

Output = sign
$$(\mathbf{w}^{\mathrm{T}}\mathbf{x} + b)$$
 = sign $\left(\sum_{i} w_{i}x_{i} + b\right)$

Linear Classifiers

Inputs are d dimensional feature vectors, denoted by **x** Output is a label $y \in \{-1, 1\}$

Linear Threshold Units classify an example \mathbf{x} using parameters \mathbf{w} (a d dimensional vector) and b (a real number) according the following classification rule

Output = sign(
$$\mathbf{w}^{\mathrm{T}}\mathbf{x} + b$$
) = sign $\left(\sum_{i} w_{i}x_{i} + b\right)$

if $\mathbf{w}^{\mathrm{T}}\mathbf{x} + b \ge 0$ then predict y = +1if $\mathbf{w}^{\mathrm{T}}\mathbf{x} + b < 0$ then predict y = -1

b is called the <u>bias</u> term

The geometry of a linear classifier An illustration in two dimensions

X₂

X₁

The geometry of a linear classifier An illustration in two dimensions















We can stop writing b at each step using notational sugar:

The prediction function is $sgn(\mathbf{w}^T\mathbf{x} + b) = sgn(\sum_i w_i x_i + b)$

We can stop writing b at each step using notational sugar:

The prediction function is $\operatorname{sgn}(\mathbf{w}^T \mathbf{x} + b) = \operatorname{sgn}(\sum_i w_i x_i + b)$ Rewrite \mathbf{x} as $\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$. Call this \mathbf{x}' Rewrite \mathbf{w} as $\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$. Call this \mathbf{w}'

We can stop writing b at each step using notational sugar:

The prediction function is $\operatorname{sgn}(\mathbf{w}^T \mathbf{x} + b) = \operatorname{sgn}(\sum_i w_i x_i + b)$ Rewrite \mathbf{x} as $\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$. Call this \mathbf{x}' Rewrite \mathbf{w} as $\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$. Call this \mathbf{w}'

Note that $\mathbf{w}^{T}\mathbf{x} + b$ is the same as $\mathbf{w'}^{T}\mathbf{x'}$

We can stop writing b at each step using notational sugar:

The prediction function is $\operatorname{sgn}(\mathbf{w}^T \mathbf{x} + b) = \operatorname{sgn}(\sum_i w_i x_i + b)$ Rewrite \mathbf{x} as $\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$. Call this \mathbf{x}' Rewrite \mathbf{w} as $\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$. Call this \mathbf{w}'

Note that $\mathbf{w}^{T}\mathbf{x} + b$ is the same as $\mathbf{w'}^{T}\mathbf{x'}$

The prediction function is now $sgn({w'}^T x')$

Increases dimensionality by one

Equivalent to adding a feature that is a constant: always 1

We can stop writing b at each step using notational sugar:

The prediction function is $\operatorname{sgn}(\mathbf{w}^T\mathbf{x} + b) = \operatorname{sgn}(\sum_i w_i x_i + b)$ Rewrite \mathbf{x} as $\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$. Call this \mathbf{x}' Rewrite \mathbf{w} as $\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$. Call this \mathbf{w}'

Note that $\mathbf{w}^{T}\mathbf{x} + b$ is the same as $\mathbf{w'}^{T}\mathbf{x'}$

The prediction function is now $sgn(\mathbf{w'}^T\mathbf{x'})$

Increases dimensionality by one

Equivalent to adding a feature that is a constant: always 1

In the increased dimensional space, the vector \mathbf{w}' goes through the origin

We can stop writing b at each step using notational sugar:

The prediction function is $\operatorname{sgn}(\mathbf{w}^T \mathbf{x} + b) = \operatorname{sgn}(\sum_i w_i x_i + b)$ Rewrite \mathbf{x} as $\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$. Call this \mathbf{x}' Rewrite \mathbf{w} as $\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$. Call this \mathbf{w}'

Note that $\mathbf{w}^{T}\mathbf{x} + b$ is the same as $\mathbf{w'}^{T}\mathbf{x'}$

The prediction function is now $sgn(\mathbf{w'}^T\mathbf{x'})$

Increases dimensionality by one

Equivalent to adding a feature that is a constant: always 1

In the increased dimensional space, the vector \mathbf{w}' goes through the origin

We sometimes hide the bias b, and instead fold the bias term into the weights by adding an extra constant feature. But remember that it is there.

Questions?

Coming up (next several weeks): Linear classification

- **Perceptron**: Error driven learning, updates the hypothesis if there is an error
- **Support Vector Machines**: Define a different cost function that includes an error term *and* a term that targets future performance
- Naïve Bayes classifier: A simple linear classifier with a probabilistic interpretation
- Logistic regression: Another probabilistic linear classifier, bears similarity to support vector machines

In all cases, the prediction will be done with the same rule: $\mathbf{w}^{\mathrm{T}}\mathbf{x} + b \ge 0 \Rightarrow y = +1$ $\mathbf{w}^{\mathrm{T}}\mathbf{x} + b < 0 \Rightarrow y = -1$