

At the end of the class you should be able to:

- show how decision-tree learning works on small examples
- explain the relationship between linear and logistic regression
- explain the updates of stochastic gradient descent

Many learning algorithms can be seen as deriving from:

- decision trees
- linear (and non-linear) classifiers

Learning Decision Trees

- Representation is a decision tree.
- Bias is towards simple decision trees.
- Search through the space of decision trees, from simple decision trees to more complex ones.

A (binary) **decision tree** (for a particular target feature) is a tree where:

- each internal (non-leaf) node is labeled with a condition, a Boolean function of examples, built using input features
- each internal node has two branches, one labeled *true* and the other *false*
- each leaf of the tree is labeled with a **point estimate** of the target feature.

Decision trees are also called **classification trees** when the target is discrete, and **regression trees** when the target is real-valued.

A (binary) **decision tree** (for a particular target feature) is a tree where:

- each internal (non-leaf) node is labeled with a condition, a Boolean function of examples, built using input features
- each internal node has two branches, one labeled *true* and the other *false*
- each leaf of the tree is labeled with a **point estimate** of the target feature.

Decision trees are also called **classification trees** when the target is discrete, and **regression trees** when the target is real-valued.

- Like an if–then–else structure in a programming language.

Example Classification Data

Training Examples:

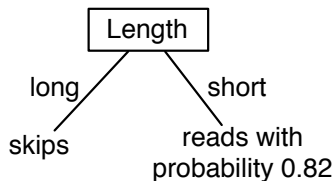
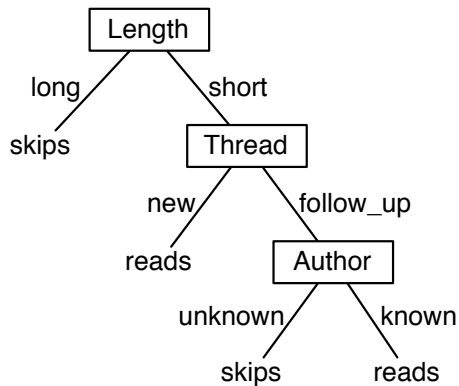
	Action	Author	Thread	Length	Where
e1	skips	known	new	long	home
e2	reads	unknown	new	short	work
e3	skips	unknown	old	long	work
e4	skips	known	old	long	home
e5	reads	known	new	short	home
e6	skips	known	old	long	work

New Examples:

e7	???	known	new	short	work
e8	???	unknown	new	short	work

We want to classify new examples on feature *Action* based on the examples' *Author*, *Thread*, *Length*, and *Where*.

Example Decision Trees



Equivalent Programs

```
define action(e):  
  if long(e): return skips  
  else if new(e): return reads  
  else if unknown(e): return skips  
  else: return reads
```


Equivalent Programs

define action(e):

if $long(e)$: return *skips*

else if $new(e)$: return *reads*

else if $unknown(e)$: return *skips*

else: return *reads*

Logic Program:

$reads(E) \leftarrow short(E) \wedge new(E)$.

$reads(E) \leftarrow short(E) \wedge follow_up(E) \wedge known(E)$.

$skips(E) \leftarrow long(E)$.

$skips(E) \leftarrow short(E) \wedge follow_up(E) \wedge unknown(E)$.

Equivalent Programs

define action(e):

if $long(e)$: return *skips*

else if $new(e)$: return *reads*

else if $unknown(e)$: return *skips*

else: return *reads*

Logic Program:

$reads(E) \leftarrow short(E) \wedge new(E)$.

$reads(E) \leftarrow short(E) \wedge follow_up(E) \wedge known(E)$.

$skips(E) \leftarrow long(E)$.

$skips(E) \leftarrow short(E) \wedge follow_up(E) \wedge unknown(E)$.

or with negation as failure:

$reads \leftarrow short \wedge new$.

$reads \leftarrow short \wedge \sim new \wedge known$.

Equivalent Programs

define action(e):

if $long(e)$: return *skips*

else if $new(e)$: return *reads*

else if $unknown(e)$: return *skips*

else: return *reads*

Logic Program:

$reads(E) \leftarrow short(E) \wedge new(E)$.

$reads(E) \leftarrow short(E) \wedge follow_up(E) \wedge known(E)$.

$skips(E) \leftarrow long(E)$.

$skips(E) \leftarrow short(E) \wedge follow_up(E) \wedge unknown(E)$.

or with negation as failure:

$reads \leftarrow short \wedge new$.

$reads \leftarrow short \wedge \sim new \wedge known$.

or as a logical formula: $reads \leftrightarrow short \wedge (new \vee known)$

- Given some training examples, which decision tree should be generated?

Issues in decision-tree learning

- Given some training examples, which decision tree should be generated?
- A decision tree can represent any discrete function of the input features.
- You need a **bias**.

Issues in decision-tree learning

- Given some training examples, which decision tree should be generated?
- A decision tree can represent any discrete function of the input features.
- You need a **bias**. Example, prefer the smallest tree.

- Given some training examples, which decision tree should be generated?
- A decision tree can represent any discrete function of the input features.
- You need a **bias**. Example, prefer the smallest tree. Least depth? Fewest nodes?

- Given some training examples, which decision tree should be generated?
- A decision tree can represent any discrete function of the input features.
- You need a **bias**. Example, prefer the smallest tree. Least depth? Fewest nodes? Which trees are the best predictors of unseen data?

- Given some training examples, which decision tree should be generated?
- A decision tree can represent any discrete function of the input features.
- You need a **bias**. Example, prefer the smallest tree. Least depth? Fewest nodes? Which trees are the best predictors of unseen data?
- How should you go about building a decision tree? The space of decision trees is too big for systematic search for the smallest decision tree.

Searching for a Good Decision Tree

- The input is a set of input features, a target feature and, a set of training examples.
- Either:
 - ▶ Stop and return a value for the target feature or a distribution over target feature values
 - ▶ Choose a condition (e.g. an input feature) to split on. build a subtree for those examples with with the condition true and the examples with the condition false.

Choices in implementing the algorithm

- When to stop:

- When to stop:
 - ▶ no more input features
 - ▶ all examples are classified the same
 - ▶ too few examples to make an informative split
 - ▶ no split give an appreciable improvement in error

Choices in implementing the algorithm

- When to stop:
 - ▶ no more input features
 - ▶ all examples are classified the same
 - ▶ too few examples to make an informative split
 - ▶ no split give an appreciable improvement in error
- Which test to split on isn't defined. Often we use **myopic** split: which single split gives smallest error?

- 1: **procedure** *DT_learner*(Cs, Y, Es, γ)
- 2: **Inputs** Cs : set of possible conditions; Y : target feature;
 Es : training examples; γ : improvement threshold
- 3: **Output** function to predict a value of Y for an example
- 4: $c := \text{select_split}(Es, Cs, \gamma)$ ▷ see next slide
- 5: **if** $c = \text{None}$ **then** ▷ stopping criterion is true
- 6: $v := \text{leaf_prediction}(Y, Es)$ ▷ Prediction on Y
- 7: **define** $T(e) = v$
- 8: **return** T
- 9: **else**
- 10: $\text{true_examples} := \{e \in Es : c(e)\}$
- 11: $t_1 := \text{DT_learner}(Cs \setminus \{c\}, Y, \text{true_examples}, \gamma)$
- 12: $\text{false_examples} := \{e \in Es : \neg c(e)\}$
- 13: $t_0 := \text{DT_learner}(Cs \setminus \{c\}, Y, \text{false_examples}, \gamma)$
- 14: **define** $T(e) = \text{if } c(e) \text{ then } t_1(e) \text{ else } t_0(e)$
- 15: **return** T

```

1: procedure select_split(Es, Y, Cs,  $\gamma$ )
2:   best_val := sum_loss(Y, Es) -  $\gamma$ 
3:   best_split := None
4:   for  $c \in Cs$  do
5:     val := sum_loss(Y, { $e \in Es \mid c(e)$ })
6:           + sum_loss(Y, { $e \in Es \mid \neg c(e)$ })
7:     if val < best_val then
8:       best_val := val
9:       best_split := c
10:  return best_split

```

For log loss: Prediction is empirical proportion of Y value

- $P = \text{leaf_prediction}(Y, Es) : v \mapsto \frac{|\{e' \in Es : Y(e')=v\}|}{|Es|}$

$$\text{sum_loss}(Y, Es) = \sum_{e \in Es} \log(P(Y(e)))$$

Example Classification Data

Training Examples:

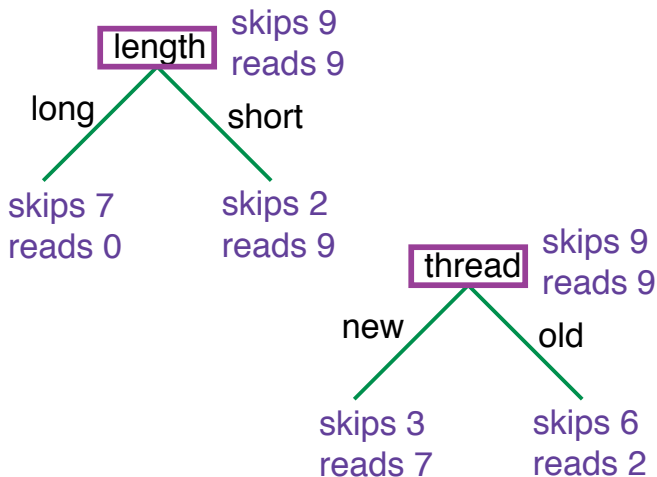
	Action	Author	Thread	Length	Where
e1	skips	known	new	long	home
e2	reads	unknown	new	short	work
e3	skips	unknown	old	long	work
e4	skips	known	old	long	home
e5	reads	known	new	short	home
e6	skips	known	old	long	work

New Examples:

e7	???	known	new	short	work
e8	???	unknown	new	short	work

Aim: classify new examples on feature *Action* based on the examples' *Author*, *Thread*, *Length*, and *Where*.

Example: possible splits



Handling Overfitting

- This algorithm can overfit the data.
This occurs when

- This algorithm can overfit the data.
This occurs when noise and correlations in the training set that are not reflected in the data as a whole.
- To handle overfitting:
 - ▶ restrict the splitting, and split only when the split is useful.
 - ▶ allow unrestricted splitting and prune the resulting tree where it makes unwarranted distinctions.
 - ▶ learn multiple trees and average them (decision forests, random forests)

Linear Function

A **linear function** of features X_1, \dots, X_n is a function of the form:

$$f^{\bar{w}}(X_1, \dots, X_n) = w_0 + w_1 X_1 + \dots + w_n X_n$$

Invent a new feature X_0 which has value 1, to make it not a special case.

$$f^{\bar{w}}(X_1, \dots, X_n) = \sum_{i=0}^n w_i X_i$$

Linear Regression

- Aim: predict feature Y from features X_1, \dots, X_n .
- A feature is a function of an example.
 $X_i(e)$ is the value of feature X_i on example e .
- **Linear regression**: predict a linear function of the input features.

$$\begin{aligned}\widehat{Y}^{\bar{w}}(e) &= w_0 + w_1 X_1(e) + \dots + w_n X_n(e) \\ &= \sum_{i=0}^n w_i X_i(e) ,\end{aligned}$$

$\widehat{Y}^{\bar{w}}(e)$ is the predicted value for Y on example e .
It depends on the weights \bar{w} .

Sum of squares error for linear regression

The sum of squares error on examples E for target Y is:

$$\begin{aligned}SSE(E, \bar{w}) &= \sum_{e \in E} (Y(e) - \hat{Y}^{\bar{w}}(e))^2 \\ &= \sum_{e \in E} \left(Y(e) - \sum_{i=0}^n w_i * X_i(e) \right)^2.\end{aligned}$$

Goal: given examples E , find weights that minimize $SSE(E, \bar{w})$.

Finding weights that minimize $Error(E, \bar{w})$

- Find the minimum analytically.
Effective when it can be done (e.g., for linear regression).

Finding weights that minimize $Error(E, \bar{w})$

- Find the minimum analytically.
Effective when it can be done (e.g., for linear regression).
- Find the minimum iteratively.
Works for larger classes of problems.
Gradient descent:

$$w_i \leftarrow w_i - \eta \frac{\partial}{\partial w_i} Error(E, \bar{w})$$

η is the gradient descent step size, the **learning rate**.

Finding weights that minimize $Error(E, \bar{w})$

- Find the minimum analytically.
Effective when it can be done (e.g., for linear regression).
- Find the minimum iteratively.
Works for larger classes of problems.
Gradient descent:

$$w_i \leftarrow w_i - \eta \frac{\partial}{\partial w_i} Error(E, \bar{w})$$

η is the gradient descent step size, the **learning rate**.

- Often update weights after each example:
 - **incremental gradient descent** updates parameters after each example
 - **stochastic gradient descent** updates parameters after a batch of (randomly selected) examplesOften much faster than updating weights after sweeping through examples, but may not converge to a local optimum

Linear Classifier

- Assume you are doing binary classification, with classes $\{0, 1\}$ (e.g., using indicator functions).

Linear Classifier

- Assume you are doing binary classification, with classes $\{0, 1\}$ (e.g., using indicator functions).
- There is no point in making a prediction of less than 0 or greater than 1.
- A **squashed linear function** is of the form:

$$f^{\bar{w}}(X_1, \dots, X_n) = f(w_0 + w_1 X_1 + \dots + w_n X_n)$$

where f is an **activation function**.

Linear Classifier

- Assume you are doing binary classification, with classes $\{0, 1\}$ (e.g., using indicator functions).
- There is no point in making a prediction of less than 0 or greater than 1.
- A **squashed linear function** is of the form:

$$f^{\bar{w}}(X_1, \dots, X_n) = f(w_0 + w_1 X_1 + \dots + w_n X_n)$$

where f is an **activation function**.

- A simple activation function is the step function:

$$f(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

Cannot be used in gradient descent because

Linear Classifier

- Assume you are doing binary classification, with classes $\{0, 1\}$ (e.g., using indicator functions).
- There is no point in making a prediction of less than 0 or greater than 1.
- A **squashed linear function** is of the form:

$$f^{\bar{w}}(X_1, \dots, X_n) = f(w_0 + w_1 X_1 + \dots + w_n X_n)$$

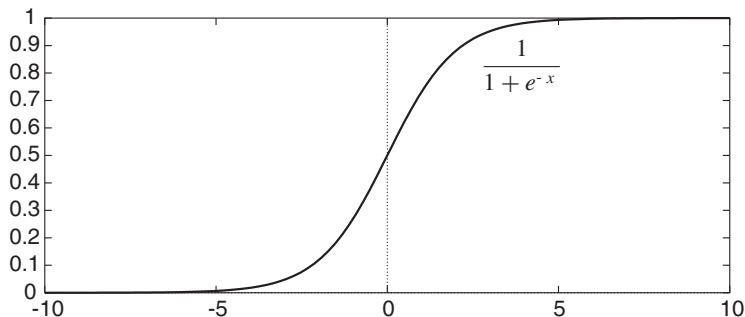
where f is an **activation function**.

- A simple activation function is the step function:

$$f(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

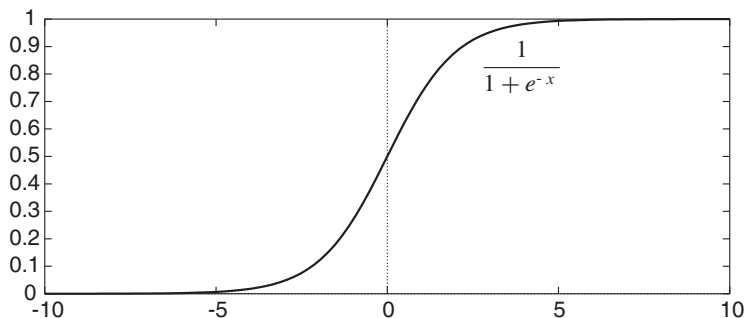
Cannot be used in gradient descent because it has a derivative of 0 almost everywhere (except at 0)

The sigmoid or logistic activation function



$$f(x) = \frac{1}{1 + e^{-x}}$$

The sigmoid or logistic activation function



$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = f(x)(1 - f(x))$$

A **logistic function** is the sigmoid of a linear function.

Logistic regression: find weights to minimize log loss of a logistic function.

Error for Squashed Linear Function

When the domain of target Y is $\{0, 1\}$:

- $\hat{Y}(e) = \text{sigmoid}(\sum_{i=0}^n w_i * X_i(e)).$
- $\delta(e) = Y(e) - \hat{Y}^w(e)$

Error for Squashed Linear Function

When the domain of target Y is $\{0, 1\}$:

- $\hat{Y}(e) = \text{sigmoid}(\sum_{i=0}^n w_i * X_i(e))$.
- $\delta(e) = Y(e) - \hat{Y}^{\bar{w}}(e)$

A natural measure for sigmoid is log loss:

$$LL(E, \bar{w}) = \sum_{e \in E} Y(e) * \log \hat{Y}(e) + (1 - Y(e)) * \log(1 - \hat{Y}(e))$$

$$\frac{\partial}{\partial w_i} LL(E, \bar{w}) =$$

Error for Squashed Linear Function

When the domain of target Y is $\{0, 1\}$:

- $\hat{Y}(e) = \text{sigmoid}(\sum_{i=0}^n w_i * X_i(e))$.
- $\delta(e) = Y(e) - \hat{Y}^{\bar{w}}(e)$

A natural measure for sigmoid is log loss:

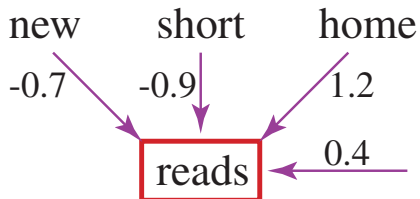
$$LL(E, \bar{w}) = \sum_{e \in E} Y(e) * \log \hat{Y}(e) + (1 - Y(e)) * \log(1 - \hat{Y}(e))$$

$$\frac{\partial}{\partial w_i} LL(E, \bar{w}) = \sum_{e \in E} \delta(e) * X_i(e)$$

Linear Learner with Stochastic Gradient Descent

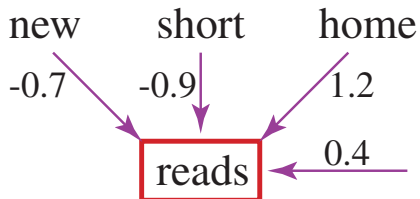
```
1: procedure Linear_learner( $Xs, Y, Es, \eta, b$ )
2:   • Input features:  $Xs = \{X_1, \dots, X_n\}$ . Target feature:  $Y$ .
   Examples:  $Es$ . Learning rate:  $\eta$ . Batch size:  $b$ 
3:   initialize  $w_0, \dots, w_n$  randomly
4:   define  $pred(e) = \phi(\sum_i w_i * X_i(e))$ 
5:   repeat
6:     for each  $i \in [0, n]$  do  $d[i] := 0$ 
7:     select batch  $B \subseteq Es$  of size  $b$ 
8:     for each example  $e$  in  $B$  do
9:        $error := pred(e) - Y(e)$ 
10:      for each  $i \in [0, n]$  do
11:         $d_i := d_i + error * X_i(e)$ 
12:      for each  $i \in [0, n]$  do
13:         $w_i := w_i - \eta * d_i / b$ 
14:   until termination
15:   return  $pred$ 
```

Simple Example



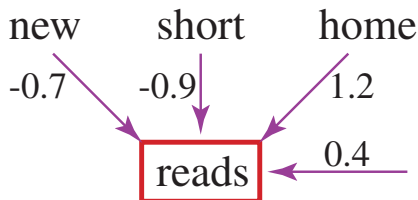
Ex	new	short	home	reads		δ	SSE
				Predicted	Obs		
e1	0	0	0	$f(0.4) \approx 0.6$	0	-0.6	0.36
e2	1	1	0		0		
e3	1	0	1		1		

Simple Example



Ex	new	short	home	reads		δ	SSE
				Predicted	Obs		
e1	0	0	0	$f(0.4) \approx 0.6$	0	-0.6	0.36
e2	1	1	0	$f(-1.2) \approx 0.23$	0		
e3	1	0	1	$f(0.9) \approx 0.71$	1		

Simple Example



Ex	new	short	home	reads		δ	SSE
				Predicted	Obs		
e1	0	0	0	$f(0.4) \approx 0.6$	0	-0.6	0.36
e2	1	1	0	$f(-1.2) \approx 0.23$	0	-0.23	0.053
e3	1	0	1	$f(0.9) \approx 0.71$	1	0.29	0.084

Linearly Separable

- A classification is **linearly separable** if there is a hyperplane where the classification is *true* on one side of the hyperplane and *false* on the other side.
- For the sigmoid function, the hyperplane is when:

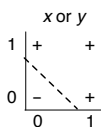
Linearly Separable

- A classification is **linearly separable** if there is a hyperplane where the classification is *true* on one side of the hyperplane and *false* on the other side.
- For the sigmoid function, the hyperplane is when:

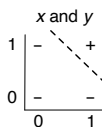
$$w_0 + w_1 * X_1 + \dots + w_n * X_n = 0$$

This separates the predictions > 0.5 and < 0.5 .

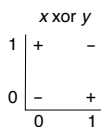
- linearly separable implies the error can be arbitrarily small



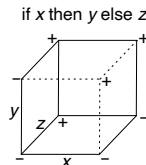
(a)



(b)



(c)



(d)

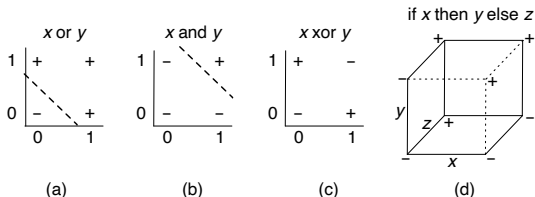
Linearly Separable

- A classification is **linearly separable** if there is a hyperplane where the classification is *true* on one side of the hyperplane and *false* on the other side.
- For the sigmoid function, the hyperplane is when:

$$w_0 + w_1 * X_1 + \dots + w_n * X_n = 0$$

This separates the predictions > 0.5 and < 0.5 .

- linearly separable implies the error can be arbitrarily small



Kernel Trick: use functions of input features (e.g., product)

Which linear separator to use can result in various algorithms:

- Perceptron
- Logistic Regression
- Support Vector Machines (SVMs)
- ...

- It's easy for a logistic function to represent “at least two of X_1, \dots, X_k are true”:

$$\frac{w_0 \quad w_1 \quad \dots \quad w_k}{\quad}$$

- It's easy for a logistic function to represent "at least two of X_1, \dots, X_k are true":

$$\frac{w_0 \quad w_1 \quad \cdots \quad w_k}{-15 \quad 10 \quad \cdots \quad 10}$$

This concept forms a large decision tree.

- Consider representing a conditional: "If X_7 then X_2 else X_3 ":
 - ▶ Simple in a decision tree.
 - ▶ For a linear separator

- It's easy for a logistic function to represent "at least two of X_1, \dots, X_k are true":

$$\frac{w_0 \quad w_1 \quad \cdots \quad w_k}{-15 \quad 10 \quad \cdots \quad 10}$$

This concept forms a large decision tree.

- Consider representing a conditional: "If X_7 then X_2 else X_3 ":
 - ▶ Simple in a decision tree.
 - ▶ For a linear separator it is impossible to represent as it is not linearly separable