# CSE/STAT 416

## Decision Trees

**Vinitra Swamy**
**University of Washington**
**July 13, 2020**

# Logistics

- If you have a question, there is a high chance somebody else in the class the same question too

- Homework 3 –
  - Extension until Friday
  - Concept question #11 has been removed

Today:

- Naïve Bayes

- Decision Trees

# Probability Classifier

**Idea**: Estimate probabilities $\hat{P}(y|x)$ and use those for prediction
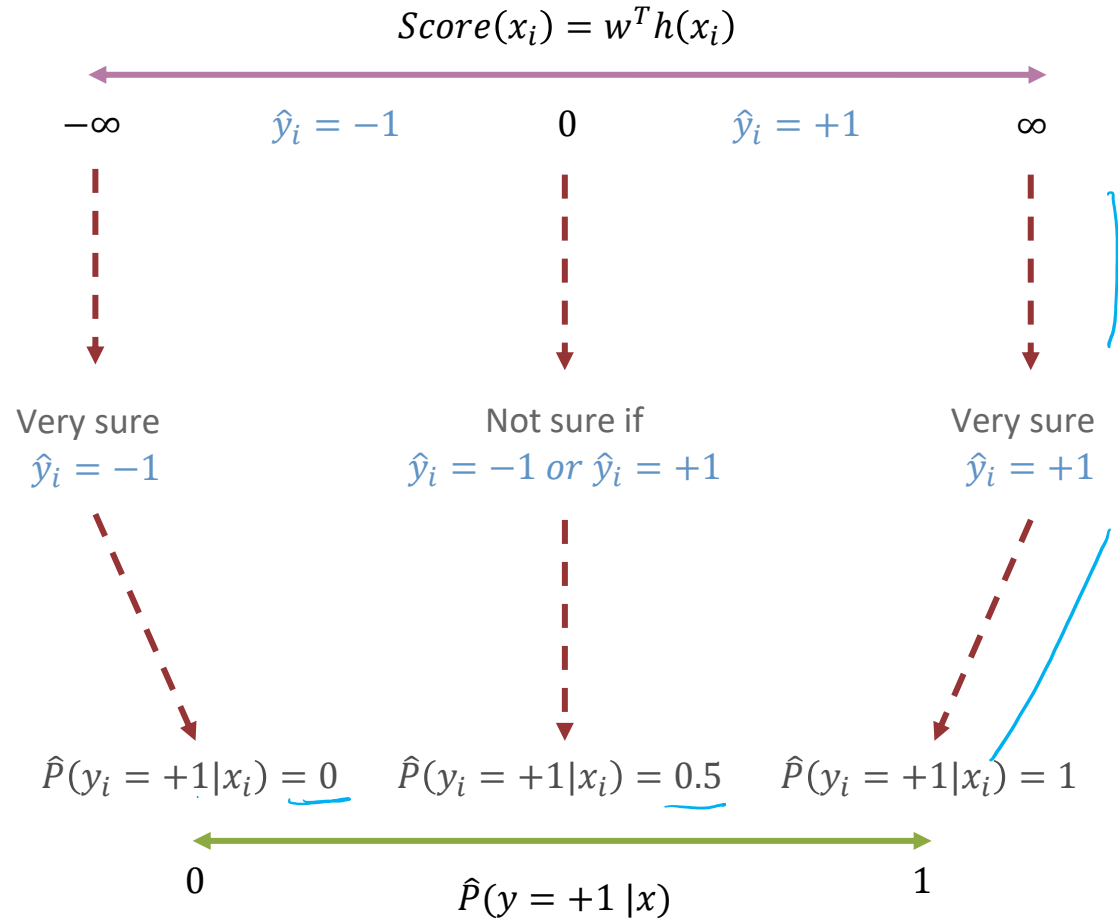
**Probability Classifier**

Input $x$: Sentence from review

- Estimate class probability $\hat{P}(y = +1|x)$

- If $\hat{P}(y = +1|x) > 0.5$:
  - $\hat{y} = +1$
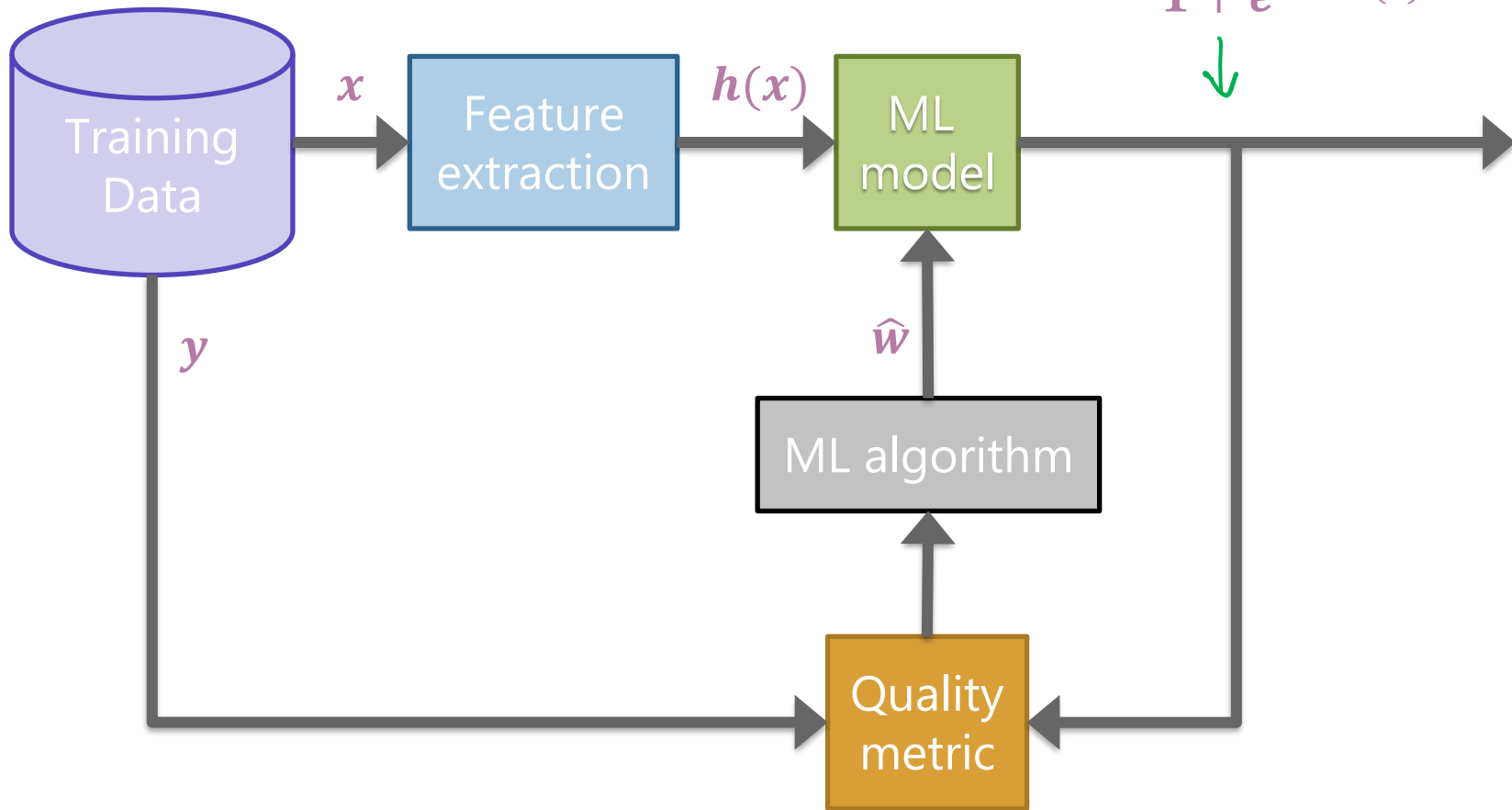
- Else:
  - $\hat{y} = -1$

**Notes**:

- Estimating the probability improves **interpretability**

# Interpreting Score

$$Score(x_i) = w^T h(x_i)$$

$-\infty$     $\hat{y}_i = -1$     $0$     $\hat{y}_i = +1$     $\infty$

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

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

Very sure
$\hat{y}_i = +1$

$\hat{P}(y_i = +1|x_i) = 0$     $\hat{P}(y_i = +1|x_i) = 0.5$     $\hat{P}(y_i = +1|x_i) = 1$

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

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

# Naïve Bayes

# Idea: Naïve Bayes

$x = $ *"The sushi & everything else was awesome!"*

$P\left(y = +1 \mid x = \text{"The sushi & everything else was awesome!"}\right)?$

$P\left(y = -1 \mid x = \text{"The sushi & everything else was awesome!"}\right)?$

**Idea**: Select the class with the highest probability!

Bayes Rule: $P(y = +1 \mid x) = \dfrac{P(x \mid y = +1) P(y = +1)}{P(x)}$

$$\dfrac{P\left(\text{"The sushi & everything else was awesome!"} \mid +1\right) \boxed{P(+1)}}{P(\text{"The sushi & everything else was awesome!"})}$$

Since we're just trying to find out which class has the greater probability, we can discard the divisor.

# Problem

**Idea**: Select the class with the highest probability!

**Problem**: We have not seen the sentence before.

**Assumption**: Words are independent from each other.

# + reviews

# reviews

$x = $ *"The sushi & everything else was awesome!"*

$$\frac{P(\text{"The sushi \& everything else was awesome!"}|+1)\ P(+1)}{P(\text{"The sushi \& everything else was awesome!"})}$$

$P(\text{"The sushi \&everything else was awesome!"} \mid +1)$
$= P(\text{The} \mid +1) * P(sushi \mid +1) * P(\& \mid +1) * P(everything \mid +1)$
$* P(else \mid +1) * P(was \mid +1) * P(awesome \mid +1)$

$$P(\text{"awesome"} \mid +1)?$$

$$= \frac{\text{\# of times "awesome" appears in + reviews}}{\text{\# of words in + reviews}}$$

# Zeros

If a feature is missing in a class everything becomes zero.

$$P(\text{"The sushi \&everything else was awesome!"} \mid + 1)$$

$$= P(\text{The} \mid +1) * P(sushi \mid + 1) * P(\& \mid + 1) * P(everything \mid + 1)$$

$$* P(else \mid + 1) * P(was \mid + 1) * P(awesome \mid + 1)$$

$$= 0$$

Solutions?

- Take the log (product becomes a sum: linear classifier)
- Laplacian Smoothing (adding a constant to avoid multiplying by zero)

# Naïve Bayes vs Logistic Regression

Logistic Regression:

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

Naïve Bayes:

$$P(y|\mathsf{x}_1, \mathsf{x}_2, \ldots, \mathsf{x}_d) = \prod_{j=1}^{d} P(x_j|y)\ P(y)$$

# Naïve Bayes vs Logistic Regression

**Generative vs Discriminative Classifiers** $P(x|y) \ P(y)$

**Generative:** defines a model for generating x (e.g. Naïve Bayes)

**Discriminative:** only cares about defining and optimizing a decision boundary (e.g. Logistic Regression) $P(y|x) \longrightarrow W_s$

# Properties

- Linear Classifier for discrete values

- Continuous Variables - Gaussian Naïve Bayes

- **Gaussian Naïve Bayes is equivalent to a Logistic Regression!**

- Naïve Bayes very efficient for discrete data: only counts

- Naïve Bayes works well for big datasets

# Multiclass Classification

- Everything works with multiple classes!



Input: **x**
**Image pixels**

Output: y
Object in image

**Take max of:**
*P(Labrador retriever|x), P(golden retriever|x), P(redbone|x), P(bloodhound|x), P(Rhodesian ridgeback|x)*

# Decision Trees

# How do we make decisions?



**COVID-19 Self-Guide Decision Tree**

Are you having symptoms? (Fever, cough, difficulty breathing)

NO → Testing is not needed. If you have been in close contact with someone who has COVID-19, stay home and monitor symptoms for 14 days.

YES → Are your symptoms severe?

YES → Call your doctor or call 911.

NO → Are you over 60 or do you have an underlying medical condition like diabetes, cancer or heart disease?

NO → Testing is not needed. Stay home for 7 days from symptom onset and 72 hours after fever is gone and symptoms improve (whichever is longer) to avoid getting anyone else sick.

YES → Call your doctor to determine if testing is needed.

There is no treatment or cure for COVID-19. For most people, the illness is generally mild and can be safely managed at home. Testing is only indicated for individuals who are at risk of serious illness, like people over 60 or with underlying medical conditions. It's important that we follow these testing guidelines to protect healthcare workers and avoid spreading the virus in our communities. Everyone, regardless of whether they have symptoms, should practice social distancing and good hand hygiene.

https://www.holzer.org/coronavirus-covid-19-updates/

15

# XOR
(Exclusive Or)

- A line might not always support our decisions.

# Credit history explained

Did I pay previous loans on time?

**Example:** excellent, good, or fair

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

# Income

What's my income?

**Example:**
$80K per year

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

# Loan terms

How soon do I need to pay the loan?

**Example:** 3 years, 5 years,…

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

# Personal information

Age, reason for the loan, marital status,...

**Example:** Home loan for a married couple

| Credit History |
| :---: |
| ★★★★ |

| Income |
| :---: |
| ★★★ |

| Term |
| :---: |
| ★★★★★ |

| Personal Info |
| :---: |
| ★★★ |

# Intelligent application

# Classifier review



Loan Application

Input: $\mathbf{x}_i$

Classifier MODEL

Output: $\hat{y}$ Predicted class

$\hat{y}_i = +1$

Safe

Risky

$\hat{y}_i = -1$

# Setup

N = 9

Data (N observations, 3 features)

| Credit | Term | Income | y |
|---------|-------|--------|-------|
| excellent | 3 yrs | high | safe |
| fair | 5 yrs | low | risky |
| fair | 3 yrs | high | safe |
| poor | 5 yrs | high | risky |
| excellent | 3 yrs | low | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

Evaluation: classification error

Many possible decisions: number of trees grows exponentially!

# Decision Trees



- **internal node:** testing a feature
- **branch:** splits into possible values of a feature
- **leaf:** final decision (the class value)

# Growing Trees

- Grow the trees using a greedy approach
- What do we need?

- feature importance to cut off possibilities early
  → weight
  → order

- stopping

## Decision stump: 1 level

| Credit | Term | Income | y |
|--------|------|--------|------|
| excellent | 3 yrs | high | safe |
| fair | 5 yrs | low | risky |
| fair | 3 yrs | high | safe |
| poor | 5 yrs | high | risky |
| excellent | 3 yrs | low | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

**Loan status:**
Safe Risky

**Root**
6 3

Split on Credit

Credit?

excellent
2  0

fair
3  1

poor
1  2

Subset of data with
Credit = excellent

Subset of data with
Credit = fair

Subset of data with
Credit = poor

# Making predictions

For each intermediate node,
set $\hat{y}$ = majority value



**Loan status:**
Safe  Risky

Root
6  3

credit?

excellent
2  0

fair
3  1

poor
1  2

Safe

Safe

Risky

# How do we select the best feature?

.

\* Select the split with lowest classification error

## Choice 1: Split on Credit

**Loan status:**
Safe  Risky

Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

## Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years

5 years

# Calculate the node values.

| Credit | Term | Income | y |
|--------|------|--------|---|
| excellent | 3 yrs | high | safe |
| fair | 5 yrs | low | risky |
| fair | 3 yrs | high | safe |
| poor | 5 yrs | high | risky |
| excellent | 3 yrs | low | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

## Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

# How do we select the best feature?

* Select the split with lowest classification error

### Choice 1: Split on Credit

**Loan status:**
Safe  Risky

Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

### Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

☕ Brain Break



9:45

# How do we measure effectiveness of a split?

**Loan status:**
Safe  Risky

Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

Idea: Calculate classification error
of this decision stump

$$\text{Error} = \frac{\text{\# mistakes}}{\text{\# data points}}$$

# Calculating classification error

**Step 1:** ŷ = class of majority of data in node

**Step 2:** Calculate classification error of predicting ŷ for this data

**Loan status:**
Safe  Risky

Root
6  3

6 correct

Safe

3 mistakes

ŷ = **majority class**

Error = $\dfrac{3}{9}$

= 0.333

| Tree | Classification error |
|------|----------------------|
| (root) | 0.33 |

# Choice 1: Split on Credit history?

Does a split on Credit reduce classification error below 0.33?

## Choice 1: Split on Credit

**Loan status:**
Safe  Risky



Root
6 3

Credit?

| excellent | fair | poor |
|-----------|------|------|
| 2  0 | 3  1 | 1  2 |

# Split on Credit: Classification error

## Choice 1: Split on Credit

**Loan status:**
Safe  Risky



Root
6  3

Credit?

| excellent | fair | poor |
|-----------|------|------|
| 2    0 | 3    1 | 1    2 |

Safe  Safe  Risky

**0 mistakes**  **1 mistake**  **1 mistake**

$$\text{Error} = \frac{2}{9}$$

$$= 0.22$$

| Tree | Classification error |
|------|----------------------|
| (root) | 0.33 |
| Split on credit | 0.22 |

## Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

Safe

Risky

# Evaluating the split on Term

## Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

Safe

Risky

1 mistake

2 mistakes

Error = $\dfrac{3}{9}$

$= 0.333$

| Tree | Classification error |
|---|---|
| (root) | 0.33 |
| Split on credit | 0.22 |
| Split on term | 0.33 |

# Choice 1 vs Choice 2: Comparing split on credit vs term

| Tree | Classification error |
|------|------|
| (root) | 0.33 |
| split on credit | 0.22 |
| split on loan term | 0.33 |

## Choice 1: Split on Credit

**Loan status:**
Safe  Risky

Root
6  3

Credit?

excellent
2  0

**WINNER**

poor
1  2

## Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

# Split Selection Summary

- Given a subset of data M (a node in a tree)

- For each remaining feature $h_i(x)$:

  1. Split data of M according to feature $h_i(x)$
  2. Compute classification error of split

- Chose feature $h^*(x)$ with lowest classification error

# Greedy Algorithm

- If split is perfect (classification error = 0) or out of features:
  - Stop

- Else:
  - repeat split selection with next stump

# Decision stump:
## 1 level

**Loan status:**
Safe Risky

Root
6 3

Split on Credit

Credit?

excellent
2 0

fair
3 1

poor
1 2

Subset of data with
Credit = excellent

Subset of data with
Credit = fair

Subset of data with
Credit = poor

# Stopping

**Loan status:**
Safe Risky

Root
6 3

Credit?

excellent
2 0

fair
3 1

poor
2 1

Safe

All data points are Safe
nothing else to do with this
subset of data

Leaf node

# Second level

# Next Step Tree

*Real valued features*

| Income | Credit | Term | y |
|--------|--------|------|---|
| $105 K | excellent | 3 yrs | Safe |
| $112 K | good | 5 yrs | Risky |
| $73 K | fair | 3 yrs | Safe |
| $69 K | excellent | 5 yrs | Safe |
| $217 K | excellent | 3 yrs | Risky |
| $120 K | good | 5 yrs | Safe |
| $64 K | fair | 3 yrs | Risky |
| $340 K | excellent | 5 yrs | Safe |
| $60 K | good | 3 yrs | Risky |

# Threshold split

**Loan status:**
Safe Risky

Root
22   18

Split on Income

Income?

< $60K
8   13

>= $60K
14   5

Subset of data with
Income >= $60K

# Best threshold?

**Infinite possible values of t**

Income = t*

Income <  t*

Income >= t*

Safe
Risky

Income

$10K

$120K

# Threshold between points

Same classification error for any threshold split between $v_A$ and $v_B$

**Income**

$v_A$  $v_B$

Safe

Risky

$10K

$120K

# Threshold split selection algorithm

- Step 1: Sort the values of a feature $h_j(x)$ :

  Let $\{v_1, v_2, v_3, \ldots v_N\}$ denote sorted values

- Step 2:
  - For i = 1 … N-1
    - Consider split $t_i = (v_i + v_{i+1}) / 2$
    - Compute classification error for threshold split $h_j(x) >= t_i$
  - Chose the $t^*$ with the lowest classification error

$$\frac{\#mistakes}{\#pts}$$

# Split on Age >= 38

# Each split partitions the 2-D space

# Depth 1: Split on x[1]



**y values**

−　+

Root
18　13

x[1]

x[1] < -0.07
13　3

x[1] >= -0.07
4　11

# Depth 2

# Threshold split caveat

y values
- +

Root
18  13

x[1]

x[1] < -0.07
13  3

x[1] >= -0.07
4  11

For threshold splits, same feature can be used multiple times

x[1]

x[2]

x[1] < -1.66
7  0

x[1] >= -1.66
6  3

x[2] < 1.55
1  11

x[2] >= 1.55
3  0

# Decision boundaries

- Decision boundaries can be complex!



**Depth 1**

**Depth 2**

**Depth 10**

# Comparing decision boundaries



**Decision Tree**

STAT/CSE 416: Intro to Machine Learning

# Overfitting

*original stopping criteria:*

*1. 0 classification error*

*2. no more features left*

- Deep decision trees are prone to overfitting
  - Decision boundaries are interpretable but not stable
  - Small change in the dataset leads to big difference in the outcome

- Overcoming Overfitting:
  - Early stopping
    - Fixed length depth
    - Stop if error does not considerably decrease
  - Pruning
    - Grow full length trees
    - Prune nodes to balance a complexity penalty

*K-D trees*

# Predicting probabilities



**Loan status:**
Safe Risky

Root
18 12

Credit?

excellent
9 2

fair
6 9

poor
3 1

Safe

Risky

Safe

$$P(y = \text{Safe} \mid x)$$

$$= \frac{3}{3 + 1} = 0.75$$

# Depth 1 probabilities

# Depth 2 probabilities

STAT/CSE 416: Intro to Machine Learning
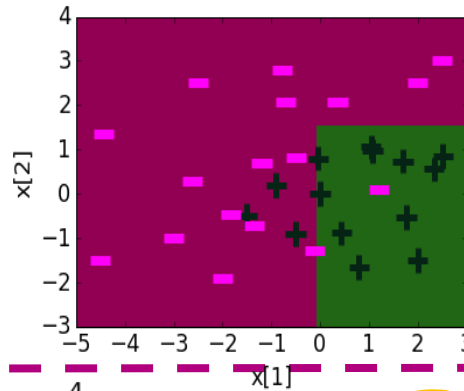
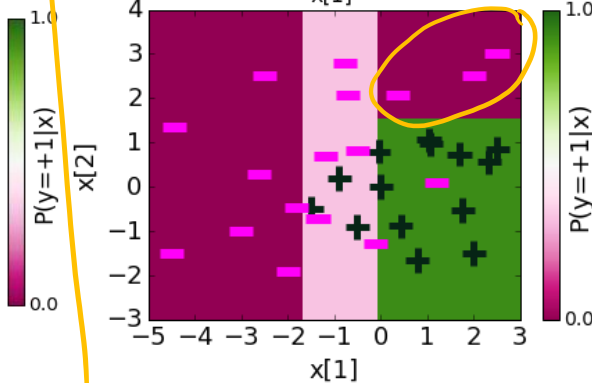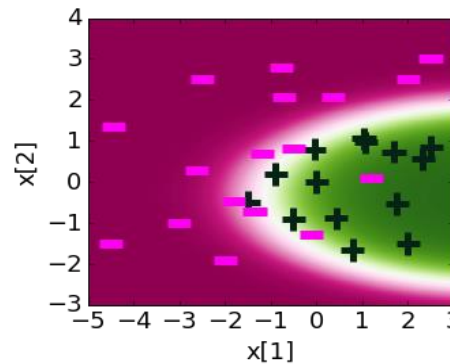# Comparison with logistic regression



**Logistic Regression (Degree 2)**

**Decision Trees (Depth 2)**

**Class**

**Probability**

# Recap

What you can do now:

- Define a decision tree classifier

- Interpret the output of a decision trees

- Learn a decision tree classifier using greedy algorithm

- Traverse a decision tree to make predictions
  - Majority class predictions