# GEOMETRIC VIEW OF DATA

David Kauchak
CS 158 – Fall 2025

1

## Admin

Assignment 2 out and due on Sunday

Assignment 1 solution posted under the "File" tab on canvas (use them to debug!)

Assignment 1 back soon

Keep reading

Mentor hours: Wednesday, 6-8pm (Alan, Edmunds downstairs)

Office hours:
- Mon, Tue, Thu: 10-11am
- Thu: 4-5pm

2

# Proper Experimentation



u13007351 fotosearch.com

3

## Experimental setup

**REAL WORLD USE OF ML ALGORITHMS**

past

Training Data

*learn*

future

Testing Data

*predict*

(data with labels)

(data without labels)

How do we tell how well we're doing?

4

## Real-world classification

Google has labeled training data, for example from people clicking the "spam" button, but when new messages come in, they're not labeled

| | | | |
|---|---|---|---|
| fmcory | (no subject) - I am in the military unit here in Afghanistan,we some some amount of funds that we war | 7:18 am |
| corowamotorinn | (no subject) - plz revert for the deal | 6:51 am |
| perfectemail1 | nnnnnnnnnnnnnnnnnnn - nnnnnnnnnnnnnnnnnnnnnnnn | 2:56 am |
| DRESURI | SOBETE | COLAN. Pregateste-ta de frig! Alege din 1000 modele de ciorapi, cumpara acum la cel mai bun pret! - Per | Sep 15 |
| Soroush Madjzoob | Stop burning money; get the most out of your investment! - Unsubscribe To remove yourself from | Sep 14 |
| Oihane Irazoki Sanchez | (no subject) - The BRITISH JUMBO COMPANY has Award your Id with the sum of 3000000.00. Send | Sep 14 |
| Long, Bruce [NS] | (no subject) - The JUMBO COMPANY has Picked you for a lump sum payout of 3000000.00. To clair | Sep 14 |
| h_044 | EEIC2013–EI–Submission: Sept 20th - 2013 3rd International Conference on Electric and Electroni | Sep 13 |
| Soroush Madjzoob | Did you know the wrong technology can cost you money? - Dear David, Technology has become t | Sep 13 |
| SantechUSA.com | Pimp Up Your Network and Save Money Doing It! - Call for consulting! 888.923.1000 FREE Our mis | Sep 13 |
| Soroush Madjzoob | When is the last time you checked your backups? - Unsubscribe To remove yourself from this ema | Sep 13 |
| Soroush Madjzoob | Is your data at risk? Get Simple, Secure & Scalable Cloud-based Backup in 3 steps! - Saccount_r | Sep 13 |
| Eden Newsletter | Get Your Free Gifts - Up To 50% Savings + Free Shipping Having trouble reading this email? view in | Sep 12 |
| AcademicPub | Meet the cutting edge in customized course materials - AcademicPub: Your Book - Your Way Acad | Sep 12 |
| Mail Administrator | Your e-mail quota has been reached! (Action Required) - Attention User, MAILBOX QUOTA EXCEI | Sep 12 |
| Wells Fargo Online | New message from Wells Fargo Online - You have 1 new message . Please Login to your account ( | Sep 12 |
| Carter, Susan | System Administrator. - Your Mailbox Is Almost Full "CLICK HERE" Update Your Mail Box And Incri | Sep 12 |

5

## Classification evaluation

Data   Label

Labeled data

0
0
1
1
0
1
0

Use the labeled data we have already to create a test set with known labels!

Why can we do this?

We assume there's an underlying distribution that generates both the training and test examples
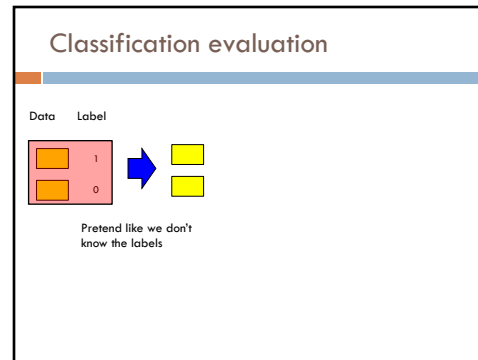
6

## Classification evaluation

Data   Label

Labeled data

0
0
1
1
0

Training data

1
0

Testing data

7

## Classification evaluation

Data   Label

Labeled data

0
0
1
1
0

Training data

train a classifier

classifier

1
0

Testing data

8

## Classification evaluation

Data    Label

1

0

Pretend like we don't
know the labels

9

## Classification evaluation

Data    Label

1

0

classifier

1

1

Pretend like we don't
know the labels

Classify

10

## Classification evaluation

Data    Label

1

0

classifier

1

1

Pretend like we don't
know the labels

Classify

How could we score
these for classification?

Compare predicted labels
to actual labels

11

## Test accuracy

To evaluate the model, compare the predicted
labels to the actual labels

prediction    Label

**Accuracy**: the proportion of
examples where we correctly
predicted the label

12

## Proper testing

Training Data

*learn*

Test Data

Evaluate model

One way to do algorithm development:
- try out an algorithm
- evaluate on test data
- repeat until happy with results

**Is this ok?**

No. Although we're not explicitly looking at the examples, we're still "cheating" by biasing our algorithm to the test data

13

## Proper testing

Test Data

Evaluate model

Once you look at/use test data **it is no longer test data!**

So, how can we evaluate our algorithm during development?

14

## Development set

Labeled Data

(data with labels)

All Training Data

Training Data

Development Data

Test Data

PEEKING

15

## Proper testing

Training Data

*learn*

Development Data

Evaluate model

Using the **development data**:
- try out an algorithm
- evaluate on development data
- repeat until happy with results

**When satisfied, evaluate on test data**

16

4

## Proper testing

Training
Data

learn

Using the **development data:**
- try out an algorithm
- evaluate on development data
- repeat until happy with results

Development
Data

Evaluate model

Any problems with this?

17

## Overfitting to development data

Be careful not to overfit to the development data!

All
Training
Data

Training
Data

Development
Data

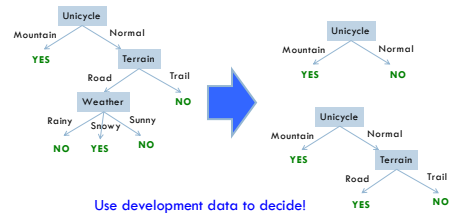Often we'll split off development data multiple times (in fact, on the fly)... you can still overfit, but this helps avoid it

18

## Pruning revisited

Unicycle

Mountain    Normal

YES         Terrain

Road    Trail

Weather    NO

Rainy  Snowy  Sunny

NO    YES    NO

Unicycle

Mountain    Normal

YES         NO

Unicycle

Mountain    Normal

YES         Terrain

Road    Trail

YES     NO

Which should we pick?

19

## Pruning revisited

Unicycle

Mountain    Normal

YES         Terrain

Road    Trail

Weather    NO

Rainy  Snowy  Sunny

NO    YES    NO

Unicycle

Mountain    Normal

YES         NO

Unicycle

Mountain    Normal

YES         Terrain

Road    Trail

YES     NO

Use development data to decide!

20

## Machine Learning: A Geometric View



21

## Apples vs. Bananas

| Weight | Color | Label |
|---|---|---|
| 4 | Red | Apple |
| 5 | Yellow | Apple |
| 6 | Yellow | Banana |
| 3 | Red | Apple |
| 7 | Yellow | Banana |
| 8 | Yellow | Banana |
| 6 | Yellow | Apple |

Can we visualize this data?

22

## Apples vs. Bananas

**Turn features into numerical values**
(read the book for a more detailed discussion of this)

| Weight | Color | Label |
|---|---|---|
| 4 | 0 | Apple |
| 5 | 1 | Apple |
| 6 | 1 | Banana |
| 3 | 0 | Apple |
| 7 | 1 | Banana |
| 8 | 1 | Banana |
| 6 | 1 | Apple |



We can view examples as points in an *n*-dimensional space where *n* is the number of features

23

## Examples in a feature space



feature$_2$

feature$_1$

- ● label 1
- ○ label 2
- ○ label 3

24

6

## Test example: what class?



feature2

feature1

- ● label 1
- ○ label 2
- ○ label 3

25

## Test example: what class?



feature2

closest to red

feature1

- ● label 1
- ○ label 2
- ● label 3

26

## Another classification algorithm?

To classify an example *d*:

Label *d* with the label of the closest example to *d* in the training set



27

## What about this example?



feature2

feature1

- ● label 1
- ○ label 2
- ○ label 3

28

## What about this example?



feature2

closest to red, but…

feature1

- ● label 1
- ○ label 2
- ○ label 3

29

## What about this example?



feature2

Most of the next closest are blue

feature1

- ● label 1
- ○ label 2
- ○ label 3

30

## k-Nearest Neighbor (k-NN)

To classify an example *d*:
- Find *k* nearest neighbors of *d*
- Choose as the label the majority label within the *k* nearest neighbors

31

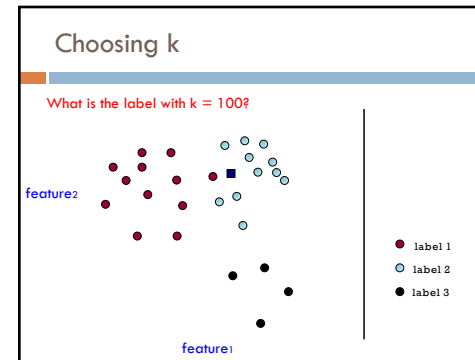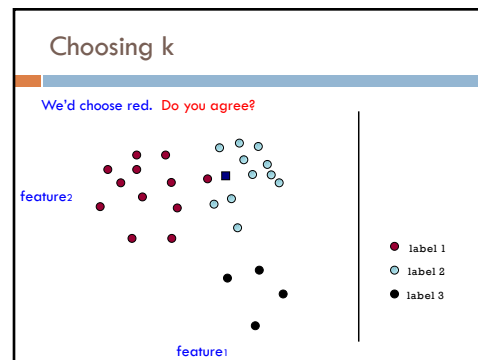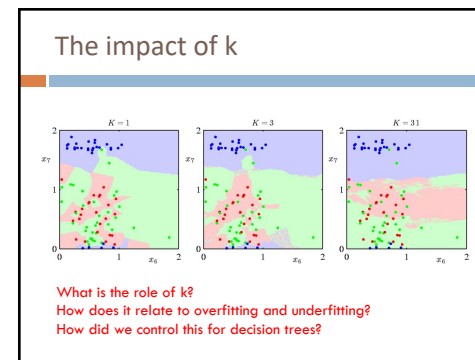## k-Nearest Neighbor (k-NN)

To classify an example *d*:
- Find *k* *nearest* neighbors of *d*
- Choose as the label the majority label within the *k* nearest neighbors

How do we measure "nearest"?

32

## Euclidean distance

In two dimensions, how do we compute the distance?

$(b_1, b_2)$

$(a_1, a_2)$

$$D(a,b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

33

## Euclidean distance

In n-dimensions, how do we compute the distance?

$(b_1, b_2, ..., b_n)$

$(a_1, a_2, ..., a_n)$

$$D(a,b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + ... + (a_n - b_n)^2}$$

34

## Euclidean distance

In n-dimensions, how do we compute the distance?

$(b_1, b_2, ..., b_n)$

$(a_1, a_2, ..., a_n)$

Measuring distance/similarity is a domain-specific problem and there are many, many different variations!

35

## Decision boundaries

The **decision boundaries** are places in the features space where the classification of a point/example changes

● label 1
○ label 2
● label 3

Where are the decision boundaries for k-NN?

36

9

## k-NN decision boundaries



label 1
label 2
label 3

k-NN gives locally defined decision boundaries between classes

37

## Choosing k

What is the label with k = 1?

feature2



label 1
label 2
label 3

feature1

38

## Choosing k

We'd choose red. Do you agree?

feature2



label 1
label 2
label 3

feature1

39

## Choosing k

What is the label with k = 3?

feature2



label 1
label 2
label 3

feature1

40

## Choosing k

We'd choose blue.  Do you agree?

feature2

feature1

● label 1
○ label 2
● label 3

41

## Choosing k

What is the label with k = 100?

feature2

feature1

● label 1
○ label 2
● label 3

42

## Choosing k

We'd choose red.  Do you agree?

feature2

feature1

● label 1
○ label 2
● label 3

43

## The impact of k

$K = 1$        $K = 3$        $K = 31$

What is the role of k?
How does it relate to overfitting and underfitting?
How did we control this for decision trees?

44

## k-Nearest Neighbor (k-NN)

To classify an example **d**:
- Find **k** nearest neighbors of **d**
- Choose as the class the majority class within the **k** nearest neighbors

How do we choose **k**?

45

## How to pick k

Common heuristics:
- often 3, 5, 7
- choose an odd number to avoid ties

Use development data

46

## k-NN variants

To classify an example **d**:
- Find **k** nearest neighbors of **d**
- Choose as the class the majority class within the **k** nearest neighbors

Any variation ideas?

47

## k-NN variations

Instead of **k** nearest neighbors, count majority from all examples within a fixed distance

Weighted **k**-NN:
- Right now, all examples are treated equally
- weight the "vote" of the examples, so that closer examples have more vote/weight
- often use some sort of exponential decay

48

## Decision boundaries for decision trees



label 1
label 2
label 3

What do the decision boundaries for decision trees like?

49

## Decision boundaries for decision trees



label 1
label 2
label 3

Axis-aligned splits/cuts of the data

50

## Decision boundaries for decision trees



label 1
label 2
label 3

What types of data sets will DT work poorly on?

51

## Problems for DT



52

## Decision trees vs. *k*-NN

Which is faster to train?

Which is faster to classify?

Do they use the features in the same way to label the examples?

53

## Decision trees vs. *k*-NN

Which is faster to train?
*k*-NN doesn't require any training!

Which is faster to classify?
For most data sets, decision trees

Do they use the features in the same way to label the examples?
k-NN treats all features equally!  Decision trees "select" important features

54

## Machine learning models

Some machine learning approaches make strong assumptions about the data
- If the assumptions are true it can often lead to better performance
- If the assumptions aren't true, the approach can fail miserably

Other approaches don't make many assumptions about the data
- This can allow us to learn from more varied data
- But, they are more prone to overfitting
- and generally require more training data

55

## Data generating distribution

We are going to use the *probabilistic model* of learning

There is some probability distribution over example/label pairs called the *data generating distribution*

**Both** the training data **and** the test set are generated based on this distribution

What is a probability distribution?

56

14

## Probability distribution

Describes how likely (i.e. probable) certain events are



- Describes probabilities for all possible events
- Probabilities are between 0 and 1 (inclusive)
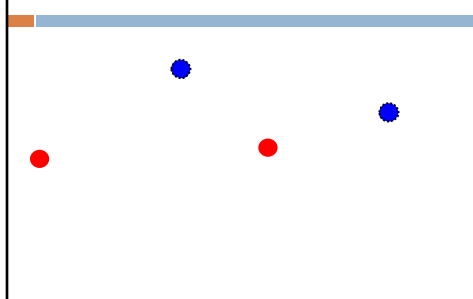- Sum of probabilities over all events is 1

57

## data generating distribution
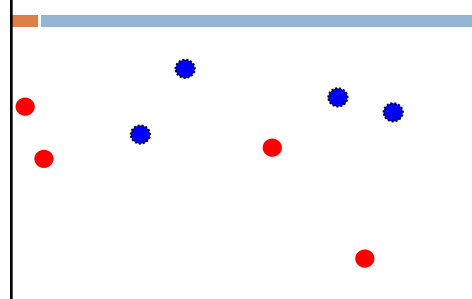


Training data          Test set

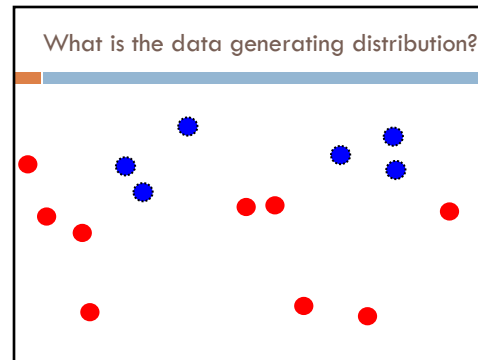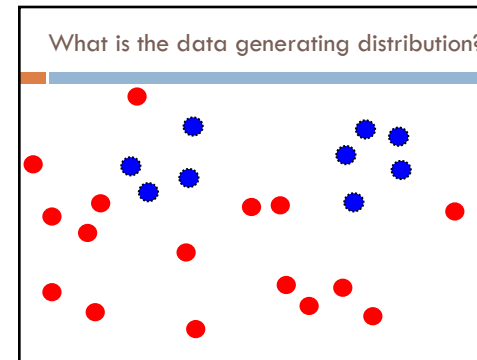data generating distribution

58

## What is the data generating distribution?
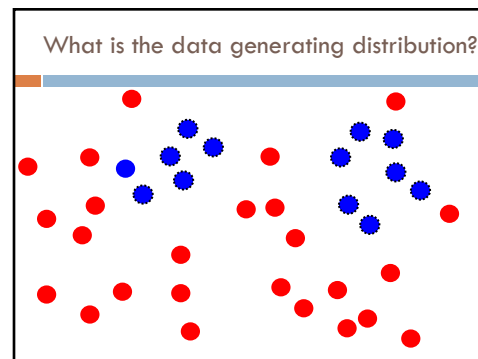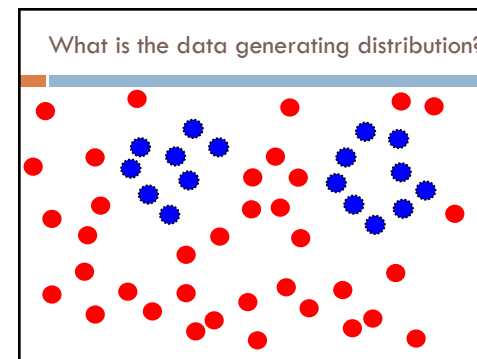


59

## What is the data generating distribution?



60

What is the data generating distribution?

61

What is the data generating distribution?

62

What is the data generating distribution?

63

What is the data generating distribution?

64

## Actual model



65

## Model assumptions

If you don't have strong assumptions about the model, it can take you a longer to learn

Assume now that our model of the blue class is two circles

66

## What is the data generating distribution?



67

## What is the data generating distribution?



68

What is the data generating distribution?

69

What is the data generating distribution?

70

What is the data generating distribution?

71

Actual model

72

## What is the data generating distribution?



Knowing the model beforehand can drastically improve the learning and the number of examples required

73

## What is the data generating distribution?



74

## Make sure your assumption is correct, though!



75

## Machine learning models

What are the *model* assumptions (if any) that *k*-NN and decision trees make about the data?

Are there data sets that could never be learned correctly by either?

76

## k-NN model



No model assumptions. Assumes that proximity relates to class

77

## Decision tree model



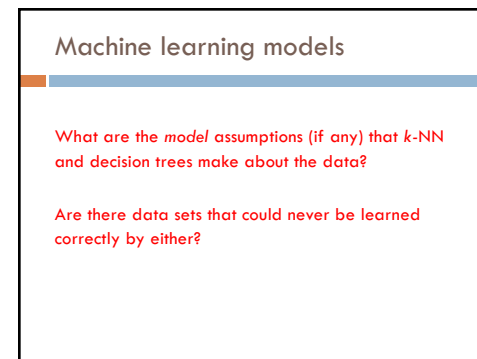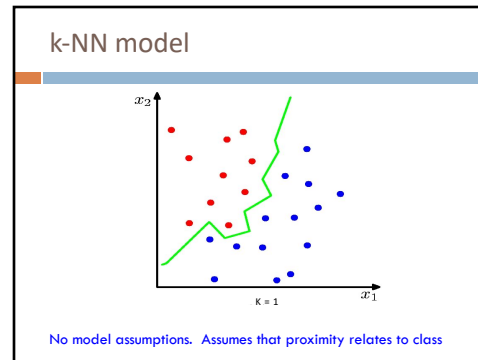- label 1
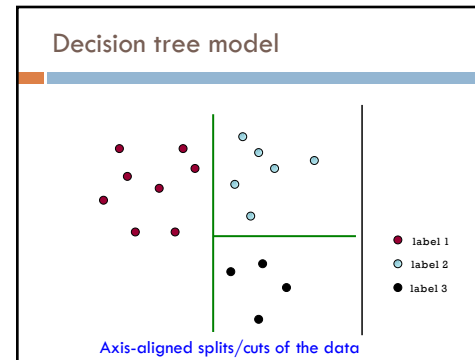- label 2
- label 3

Axis-aligned splits/cuts of the data

78

## Bias

The "bias" of a model is how strong the model assumptions are.

low-bias classifiers make minimal assumptions about the data (*k*-NN and DT are generally considered low bias)

high-bias classifiers make strong assumptions about the data

79

## Linear models

A strong high-bias assumption is *linear separability*:
- in 2 dimensions, can separate classes by a line
- in higher dimensions, need hyperplanes

A *linear model* is a model that assumes the data is linearly separable



80

## An aside: a thought experiment

What is a 100,000-dimensional space like?

You're a 1-D creature, and you decide to buy a 2-unit apartment

2 rooms (very, skinny rooms)

81

## Another thought experiment

What is a 100,000-dimensional space like?

Your job's going well and you're making good money. You upgrade to a 2-D apartment with 2-units per dimension
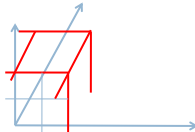
4 rooms (very, flat rooms)

82

## Another thought experiment

What is a 100,000-dimensional space like?

You get promoted again and start having kids and decide to upgrade to another dimension.

Each time you add a dimension, the amount of space you have to work with goes up exponentially

8 rooms (very, normal rooms)

83

## Another thought experiment

What is a 100,000-dimensional space like?

Sundar Pichai steps down as CEO of Google (Alphabet) and they ask you if you'd like the job. You decide to upgrade to a 100,000 dimensional apartment.

How much room do you have?
Can you have a big party?

$2^{100,000}$ rooms (it's very quiet and lonely...) = ~$10^{30}$ rooms per person if you invited everyone on the planet

84

## The challenge

Our intuitions about space/distance don't scale with dimensions!

85