

NAÏVE BAYES

David Kauchak
CS159 Fall 2020

1

Admin


Assignment 6a

Assignment 6b

2

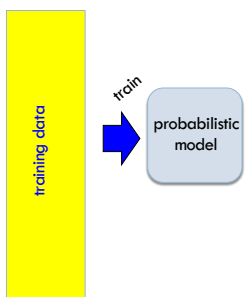
Machine Learning is...

Machine learning is about predicting the future based on the past.
-- Hal Daume III



3

Probabilistic Modeling



training data

train

probabilistic model

Model the data with a probabilistic model

specifically, learn $p(\text{features}, \text{label})$

$p(\text{features}, \text{label})$ tells us how likely these features and this example are

4

An example: classifying fruit

Training data

examples	label
red, round, leaf, 3oz, ...	apple
green, round, no leaf, 4oz, ...	apple
yellow, curved, no leaf, 4oz, ...	banana
green, curved, no leaf, 5oz, ...	banana

train → probabilistic model: $p(\text{features}, \text{label})$

5

Probabilistic models

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(\text{features}, \text{label})$ → 0.004

6

Probabilistic model vs. classifier

Probabilistic model:

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(\text{features}, \text{label})$ → 0.004

Classifier:

yellow, curved, no leaf, 6oz → probabilistic model: $p(\text{features}, \text{label})$ → banana

7

Probabilistic models: classification

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(\text{features}, \text{label})$ → 0.004

Given an unlabeled example: yellow, curved, no leaf, 6oz predict the label

How do we use a probabilistic model for classification/prediction?

8

Probabilistic models

Probabilistic models define a *probability distribution* over features and labels:



For each label, ask for the probability under the model
Pick the label with the highest probability

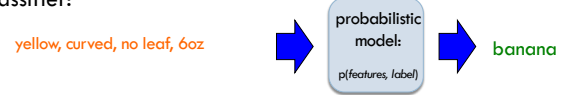
9

Probabilistic model vs. classifier

Probabilistic model:



Classifier:



Why probabilistic models?

10

Probabilistic models

Probabilities are nice to work with

- range between 0 and 1
- can combine them in a well understood way
- lots of mathematical background/theory

Provide a strong, well-founded groundwork

- Allow us to make clear decisions about things like smoothing
- Tend to be much less "heuristic"
- Models have very clear meanings

11

Probabilistic models: big questions

1. Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?
2. How do train the model, i.e. how do we estimate the probabilities for the model?
3. How do we deal with overfitting (i.e. smoothing)?

12

Basic steps for probabilistic modeling

Step 1: pick a model	<p>Probabilistic models</p> <p>Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?</p> <p>How do we train the model, i.e. how do we estimate the probabilities for the model?</p> <p>How do we deal with overfitting?</p>
Step 2: figure out how to estimate the probabilities for the model	
Step 3 (optional): deal with overfitting	

13

What was the data generating distribution?

The diagram illustrates a data generating distribution represented by a blue oval containing several fruits: two apples (one red, one green), three bananas, and two lemons. Two green arrows point upwards from this oval to two separate groups of fruits. The group on the left is labeled 'Training data' and contains a mix of these fruits. The group on the right is labeled 'Test set' and also contains a mix of these fruits, representing a sample drawn from the same underlying distribution.

14

Step 1: picking a model

What we're really trying to do is model the data generating distribution, that is how likely the feature/label combinations are

The diagram shows a blue oval containing a data generating distribution of fruits: one green apple, one red apple, three bananas, and two lemons.

15

Some math

$$p(\text{features}, \text{label}) = p(x_1, x_2, \dots, x_m, y)$$

$$= p(y)p(x_1, x_2, \dots, x_m | y)$$

What rule?

16

Some math

$$\begin{aligned}
 p(\text{features}, \text{label}) &= p(x_1, x_2, \dots, x_m, y) \\
 &= p(y) p(x_1, x_2, \dots, x_m | y) \\
 &= p(y) p(x_1 | y) p(x_2, \dots, x_m | y, x_1) \\
 &= p(y) p(x_1 | y) p(x_2 | y, x_1) p(x_3, \dots, x_m | y, x_1, x_2) \\
 &= p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})
 \end{aligned}$$

17

Step 1: pick a model

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

So, far we have made NO assumptions about the data

$$p(x_m | y, x_1, x_2, \dots, x_{m-1})$$

How many entries would the probability distribution table have if we tried to represent all possible values and we had 7000 binary features?

18

Full distribution tables

x_1	x_2	x_3	...	y	$p()$
0	0	0	...	0	*
0	0	0	...	1	*
1	0	0	...	0	*
1	0	0	...	1	*
0	1	0	...	0	*
0	1	0	...	1	*
			...		

All possible combination of features!

Table size: $2^{7000} = ?$

19

2⁷⁰⁰⁰

```

1621696755662202026466665085478377095191112430363743256235982084151527023162702352987080237879
4460004651996019099530984538652557892546513204107022110253564658647431585227076599373340842842
72242001228187826007293108261704319448426639207784125099996860169436066600112098175792966787
81962523770653294757256478055809293844657218640216108626600816097132874749264352087401101862
690842327501724665231129395523505905454421455477250950906507889478094683592939574112569473438
6191215296848474344406741204174020887540371869421701550220735398381224299258743537536161041593
43594557666561701790904172597025336526626820218084938928126997095285708906637557541434487608
82483699419938024151975145101251270438290872809195384763028578118540240995889964192277601255
360491156240349994714416090573084242931396211953679373012944795600248333570738998392029910322
34659803853069042980174009801732521069130797124201696339723021835300758978451952584853710885
819563173700074380516741189134617501484521767984296782842287373127422122022517597535994839257
02987790706355347902449354353866605125910795672914312162977887848185522928196541766009803989
979916814047493842157435158026038115106828460678973048382922034604277576550737656754730702714
466223487685709621261074762705203049488907208978593690470634285483166866563271744606581835
60966484958001276175461457216176955575199211750751406775104496728590822558547771447242334900
76402632176089211355256124119453870268029904001838585057671936968975936612135688883680023840
932567380777501891470304962150996983853975207154939633923720287592041517294930790977853625108
3200928396048072795488706954662168804465211249307629009199071774235503913511744153297374799300
89955830518881133347964641136800499940373745460035288112326328218661131046550728992296946
915601858083982074170460683212438815202609584696588161375826382921029547343888832163627122302
921229795384868355483537106034077891774170263636562027269554375177807413134551018100094688094
0781122057380335371124632958916237089580476224595091825301636909236240671411644331656159828058
3720783439888562390892028440902553829376
    
```

Any problems with this?

20

Full distribution tables

x_1	x_2	x_3	...	y	$p(\cdot)$
0	0	0	...	0	*
0	0	0	...	1	*
1	0	0	...	0	*
1	0	0	...	1	*
0	1	0	...	0	*
0	1	0	...	1	*
			...		

- Storing a table of that size is impossible!
- How are we supposed to learn/estimate each entry in the table?

21

Step 1: pick a model

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

So, far we have made NO assumptions about the data

Model selection involves making assumptions about the data

We've done this before, n-gram language model, parsing, etc.

These assumptions allow us to represent the data more compactly and to estimate the parameters of the model

22

Naïve Bayes assumption

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

$$p(x_i | y, x_1, x_2, \dots, x_{i-1}) = p(x_i | y)$$

What does this assume?

23

Naïve Bayes assumption

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

$$p(x_i | y, x_1, x_2, \dots, x_{i-1}) = p(x_i | y)$$

Assumes feature i is independent of the other features given the label

Is this true for text, say, with unigram features?

24

Naïve Bayes assumption

$$p(x_i | y, x_1, x_2, \dots, x_{i-1}) = p(x_i | y)$$

For most applications, this is not true!

For example, the fact that “San” occurs will probably make it *more likely* that “Francisco” occurs

However, this is often a reasonable approximation:

$$p(x_i | y, x_1, x_2, \dots, x_{i-1}) \approx p(x_i | y)$$

25

Naïve Bayes model

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

$$= p(y) \prod_{j=1}^m p(x_j | y) \quad \text{naïve Bayes assumption}$$

$p(x_i | y)$ is the probability of a particular feature value given the label

How do we model this?

- for binary features (e.g., “banana” occurs in the text)
- for discrete features (e.g., “banana” occurs x_i times)
- for real valued features (e.g, the text contains x_i proportion of verbs)

26

$p(x | y)$

Binary features (aka, Bernoulli Naïve Bayes) :

$$p(x_j | y) = \begin{cases} \theta_j & \text{if } x_j = 1 \\ 1 - \theta_j & \text{otherwise} \end{cases} \quad \text{biased coin toss!}$$

27

Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

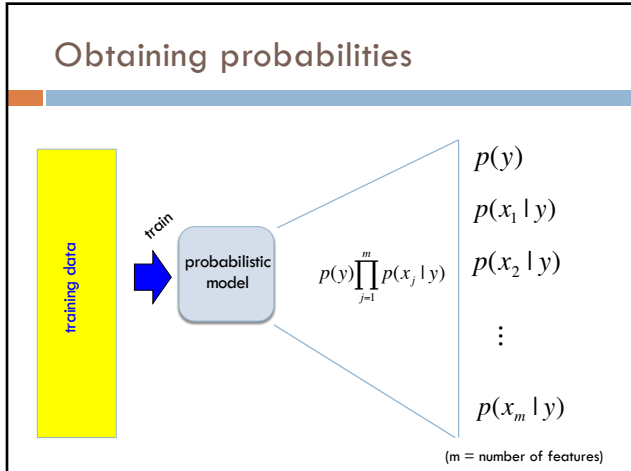
Probabilistic models

Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?

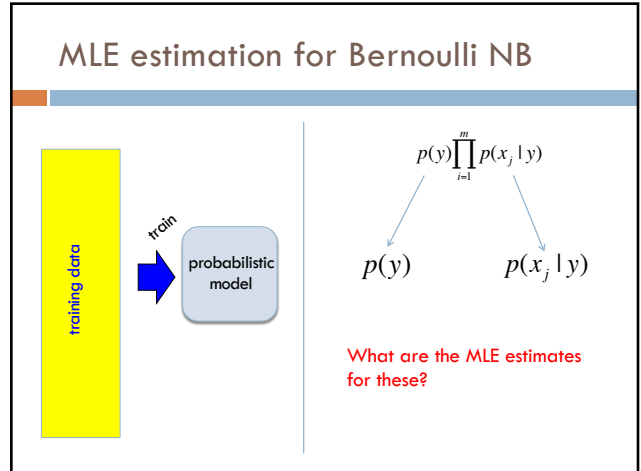
How do train the model, i.e. how to we we **estimate the probabilities** for the model?

How do we deal with overfitting?

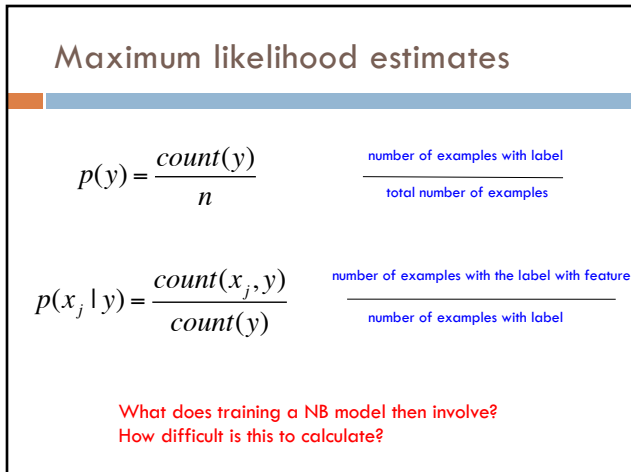
28



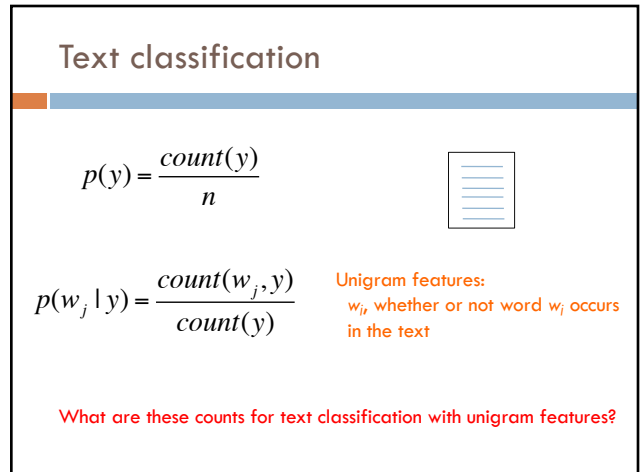
29



30



31



32

Text classification

$$p(y) = \frac{\text{count}(y)}{n}$$

number of texts with label
total number of texts

$$p(w_j | y) = \frac{\text{count}(w_j, y)}{\text{count}(y)}$$

number of texts with the label with word w_j
number of texts with label

33

Naïve Bayes classification

yellow, curved, no leaf, 6oz, banana → NB Model $p(\text{features}, \text{label})$ → 0.004

$$p(y) \prod_{j=1}^m p(x_j | y)$$

Given an unlabeled example: yellow, curved, no leaf, 6oz predict the label

How do we use a probabilistic model for classification/prediction?

34

NB classification

probabilistic model: $p(\text{features}, \text{label})$

yellow, curved, no leaf, 6oz, banana → $p(y=1) \prod_{j=1}^m p(x_j | y=1)$ → pick largest

yellow, curved, no leaf, 6oz, apple → $p(y=2) \prod_{j=1}^m p(x_j | y=2)$ →

$$\text{label} = \operatorname{argmax}_{y \in \text{labels}} p(y) \prod_{j=1}^m p(x_j | y)$$

35

NB classification

probabilistic model: $p(\text{features}, \text{label})$

yellow, curved, no leaf, 6oz, banana → $p(y=1) \prod_{j=1}^m p(x_j | y=1)$ → pick largest

yellow, curved, no leaf, 6oz, apple → $p(y=2) \prod_{j=1}^m p(x_j | y=2)$ →

Notice that each label has its own separate set of parameters, i.e. $p(x_j | y)$

36

Bernoulli NB for text classification

probabilistic model: $p(\text{features}, \text{label})$

$p(y=1) \prod_{j=1}^m p(w_j | y=1)$

$p(y=2) \prod_{j=1}^m p(w_j | y=2)$

pick largest

How good is this model for text classification?

37

Bernoulli NB for text classification

$(1, 1, 1, 0, 0, 1, 0, 0, \dots)$

$w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8$

$p(y=1) \prod_{j=1}^m p(w_j | y=1)$

$p(y=2) \prod_{j=1}^m p(w_j | y=2)$

pick largest

For text classification, what is this computation?
Does it make sense?

38

Bernoulli NB for text classification

$(1, 1, 1, 0, 0, 1, 0, 0, \dots)$

$w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8$

$p(y=1) \prod_{j=1}^m p(w_j | y=1)$

$p(y=2) \prod_{j=1}^m p(w_j | y=2)$

pick largest

Each word that occurs, contributes $p(w_j | y)$
Each word that does NOT occur, contributes $1 - p(w_j | y)$

39

Generative Story

To classify with a model, we're given an example and we obtain the probability

We can also ask how a given model would generate an example

This is the "generative story" for a model

Looking at the generative story can help understand the model

We also can use generative stories to help develop a model

40

Bernoulli NB generative story 

$$p(y) \prod_{j=1}^m p(x_j | y)$$

What is the generative story for the NB model?

41

Bernoulli NB generative story 

$$p(y) \prod_{j=1}^m p(x_j | y)$$

1. Pick a label according to $p(y)$
 - roll a biased, num_labels-sided die
2. For each feature:
 - Flip a *biased* coin:
 - if heads, include the feature
 - if tails, don't include the feature

What does this mean for text classification, assuming unigram features?

42

Bernoulli NB generative story 

$$p(y) \prod_{j=1}^m p(w_j | y)$$

1. Pick a label according to $p(y)$
 - roll a biased, num_labels-sided die
2. For each word in your vocabulary:
 - Flip a *biased* coin:
 - if heads, include the word in the text
 - if tails, don't include the word

43

Bernoulli NB

$$p(y) \prod_{j=1}^m p(x_j | y)$$

Pros/cons?

44

Bernoulli NB

Pros

- Easy to implement
- Fast!
- Can be done on large data sets

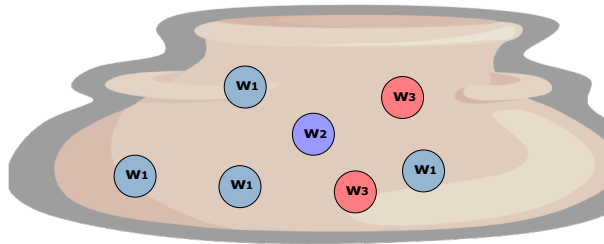
Cons

- Naïve Bayes assumption is generally not true
- Performance isn't as good as other models
- For text classification (and other sparse feature domains) the $p(x_i=0|y)$ can be problematic

45

Another generative story

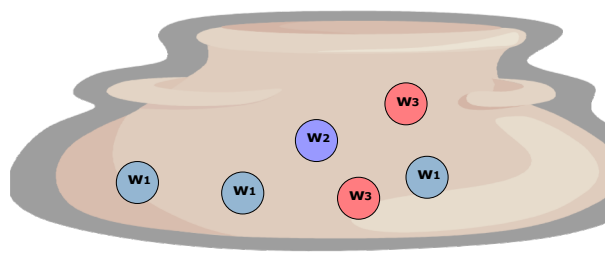
Randomly draw words from a "bag of words" until document length is reached



46

Draw words from a fixed distribution

Selected: w_1

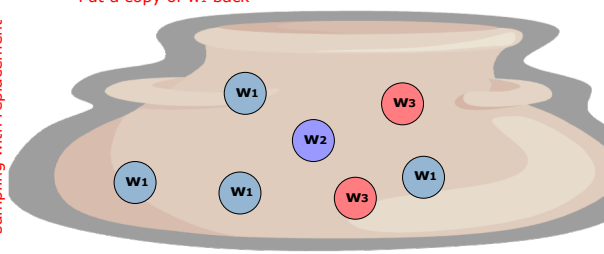


47

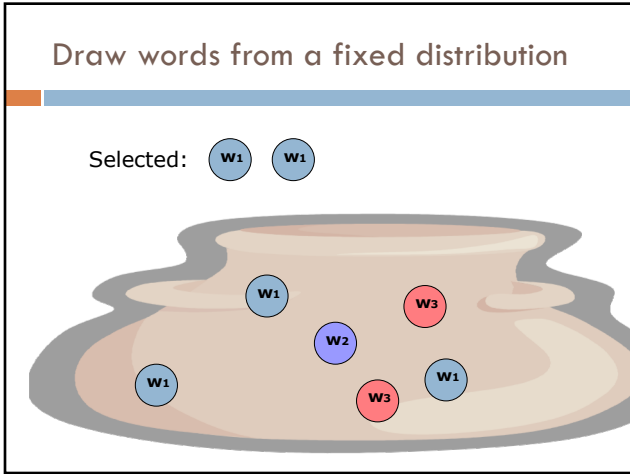
Draw words from a fixed distribution

Selected: w_1
Put a copy of w_1 back

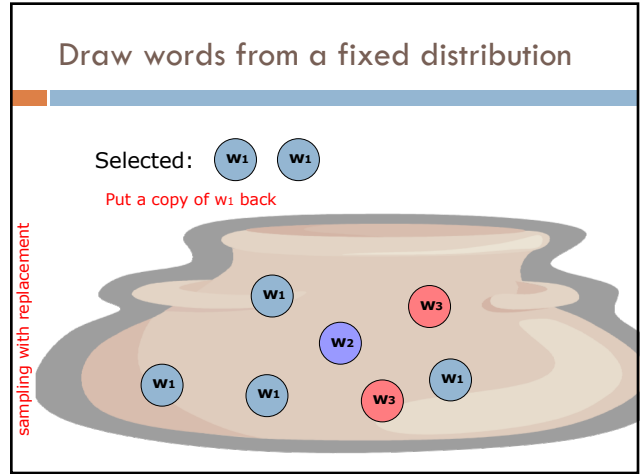
sampling with replacement



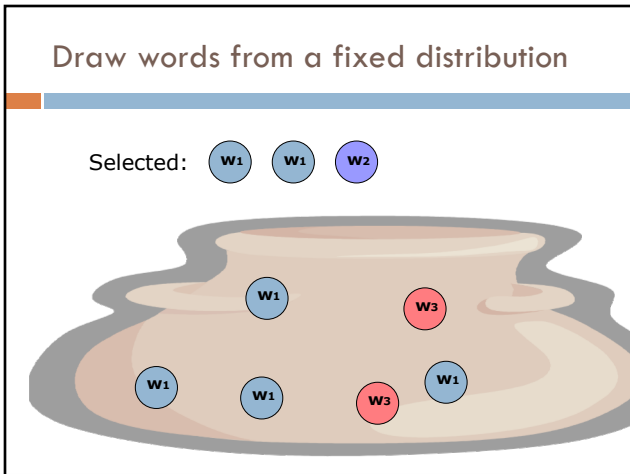
48



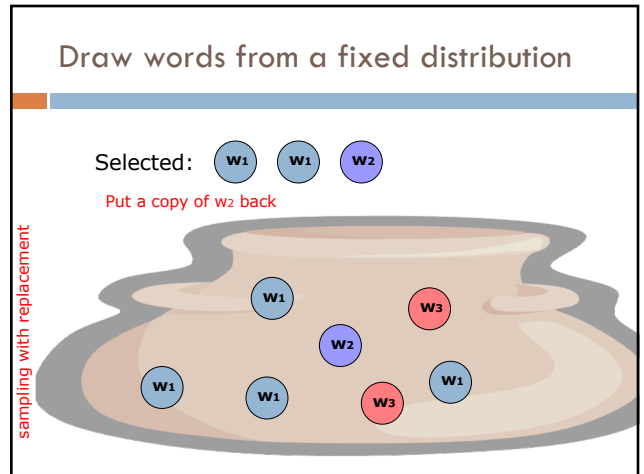
49



50



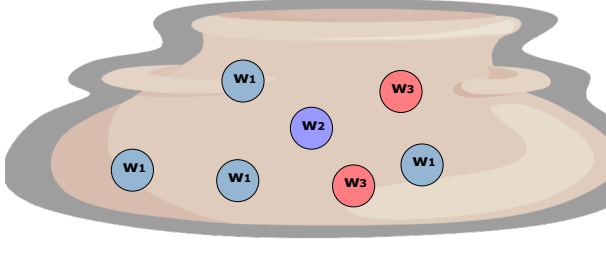
51



52

Draw words from a fixed distribution

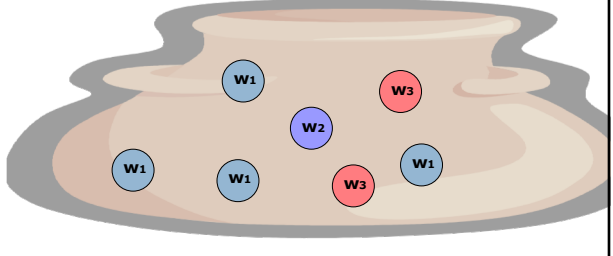
Selected: w_1 w_1 w_2 ...



53

Draw words from a fixed distribution

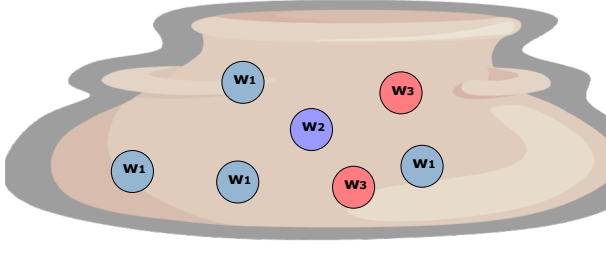
Is this a NB model, i.e. does it assume each individual word occurrence is independent?



54

Draw words from a fixed distribution

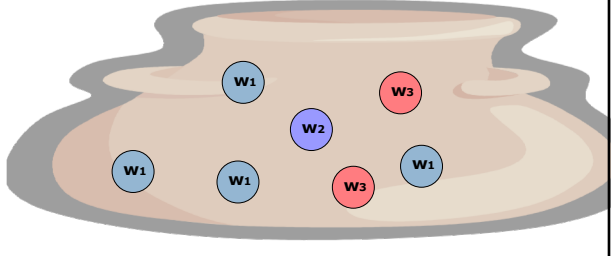
Yes! Doesn't matter what words were drawn previously, still the same probability of getting any particular word



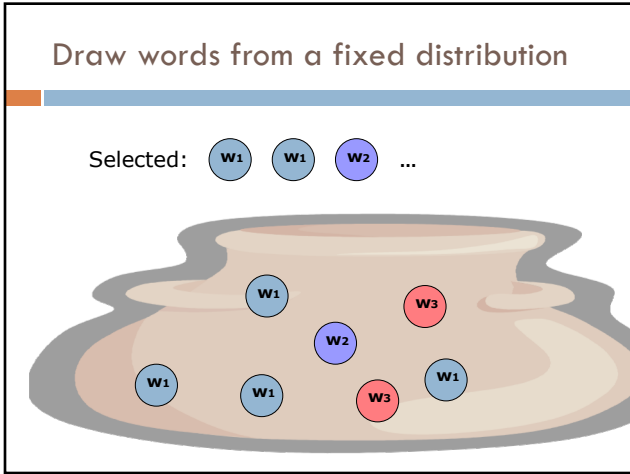
55

Draw words from a fixed distribution

Does this model handle multiple word occurrences?



56



57