

# More about Naïve Bayes

Instructor: Wei Xu

Some slides adapted from Dan Jurfasky and Brendan O'connor

# Market Research

(due 11:59pm, Friday, Feb 3rd)

- When was the company started?
- Who were the founders?
- What kind of organization is it? (publicly traded company, privately held company, non-profit organization, other)
- What is company's main business model?
- How does the company generate revenue?
- What is the finance situation of the company? (stock price, annual report, news)

# Market Research

(due 11:59pm, Friday, Feb 3rd)

- Why is the company interested in speech or NLP technologies or both?
- What are specific areas or applications of speech/ NLP the company is interested in?
- What products of the company use speech or NLP technologies?
- What the main users of their speech or NLP technologies?
- Does the company hold any patent using speech or NLP technologies?
- Does the company publish any papers on speech or NLP technologies?

# Market Research

(due 11:59pm, Friday, Feb 3rd)

- Is the company recently hiring in NLP? interns? PhD?
- What specific expertise within speech or NLP the company is looking to hire?
- How is the press coverage of the company?
- How many employees do the company have?
- An estimation of how many speech/NLP experts currently in the company?
- Any notable speech/NLP researcher or recent hires in the company?
- Which city is the company's speech/NLP research office located?

# Text Classification

?



parser  
language  
label  
translation  
...

Machine Learning	NLP	Garbage Collection	Planning	GUI
learning	<u>parser</u>	garbage	planning	...
<u>training</u>	tag	collection	temporal	
algorithm	training	memory	reasoning	
shrinkage	<u>translation</u>	optimization	plan	
network...	<u>language...</u>	region...	<u>language...</u>	

# Classification Methods: Supervised Machine Learning

- *Input:*
  - a document  $d$
  - a fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$
  - A training set of  $m$  hand-labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$
- *Output:*
  - a learned classifier  $\gamma: d \rightarrow c$

# Naïve Bayes Classifier

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

MAP is “maximum a posteriori”  
= most likely class

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

Dropping the  
denominator

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

Document  $d$   
represented as  
features  $x_1, \dots, x_n$

# Multinomial Naïve Bayes: Assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities  $P(x_i | c_j)$  are independent given the class  $c$ .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$



# Multinomial Naïve Bayes: Assumptions

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$= \operatorname{argmax}_{c \in C} \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

estimate from data

# Multinomial Naïve Bayes: Learning

- maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

use unique word  $w_i$  as  
features (in place of  $x_i$ )

# Multinomial Naïve Bayes: Learning

- Calculate  $P(c_j)$  terms
  - For each  $c_j$  in  $C$  do

$docs_j \leftarrow$  all docs with class =  $c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

# Multinomial Naïve Bayes: Learning

- maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

use unique word  $w_i$  as  
features (in place of  $x_i$ )

# Multinomial Naïve Bayes: Learning

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

fraction of times word  $w_i$  appears  
among all words in documents of  
topic  $c_j$

- Create mega-document for topic  $j$  by concatenating all docs in this topic
  - Use frequency of  $w$  in mega-document

# Zero Probabilities Problem

- What if we have seen no training documents with the word *fantastic* and classified in the topic **positive**?

$$\hat{P}(\text{"fantastic"} \mid \text{positive}) = \frac{\text{count}(\text{"fantastic"}, \text{positive})}{\sum_{w \in V} \text{count}(w, \text{positive})} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$c_{MAP} = \operatorname{argmax}_{c \in C} \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

# Problem with Maximum Likelihood

- What if we have seen no training documents with the word *fantastic* and classified in the topic **positive**?

re <https://moviepilot.com/posts/3522674>



START WRITING

SUPERHEROES HORROR YOUNG ADULT TELEVISION VIDEO



#Thefantasticfour

## THE FANTASTIC FOUR ' FAR FROM FANTASTIC

September 4, 2015 at 14:49PM



# Laplace (add-1) smoothing for Naïve Bayes

$$\begin{aligned}\hat{P}(w_i | c) &= \frac{\text{count}(w_i, c)}{\sum_{w \in V} (\text{count}(w, c))} \\ &= \frac{\text{count}(w_i, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V|}\end{aligned}$$

smoothing to avoid  
zero probabilities

- For *unknown* words (which completely doesn't occur in training set), we can ignore them.



# Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate  $P(w_k | c_j)$  terms
  - $Text_j \leftarrow$  single doc containing all *docs*<sub>j</sub>
  - For each word  $w_k$  in *Vocabulary*  
 $n_k \leftarrow$  # of occurrences of  $w_k$  in  $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha | \text{Vocabulary} |}$$

smoothing to avoid  
zero probabilities  
(often use  $\alpha = 1$ )

# Exercise

# Naïve Bayes: Practical Issues

$$\begin{aligned}c_{MAP} &= \operatorname{argmax}_c P(c|x_1, \dots, x_n) \\&= \operatorname{argmax}_c P(x_1, \dots, x_n|c)P(c) \\&= \operatorname{argmax}_c P(c) \prod_{i=1}^n P(x_i|c)\end{aligned}$$

- Multiplying together lots of probabilities
- Probabilities are numbers between 0 and 1

Q: What could go wrong here?

# Underflow



## underflow

**Contents** [\[hide\]](#)

1 English

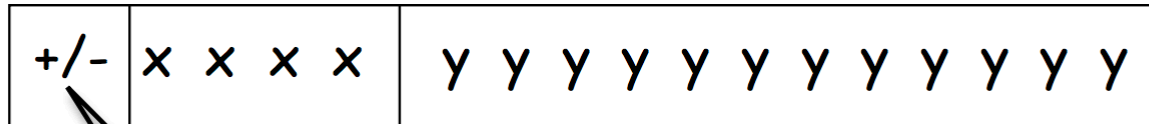
**Noun** [\[edit\]](#)

**underflow** (*plural* **underflows**)

1. (*computing*) A condition in which the **value** of a computed **quantity** is smaller than the smallest non-zero value that can be physically stored; usually treated as an **error** condition
2. (*computing*) The error condition that results from an attempt to retrieve an item from an empty **stack**

# Floating Point Numbers

**Exponent E**      significand F (also called ***mantissa***)



Sign bit

In **decimal** it means (+/-) **1**. yyyyyyyyyyyyyy  $\times 10^{\text{xxxx}}$

In **binary**, it means (+/-) **1**. yyyyyyyyyyyyyy  $\times 2^{\text{xxxx}}$

(The 1 is implied)

# Floating Point Numbers

## IEEE 754 double precision (64 bits)

S	exponent	significand
---	----------	-------------

1

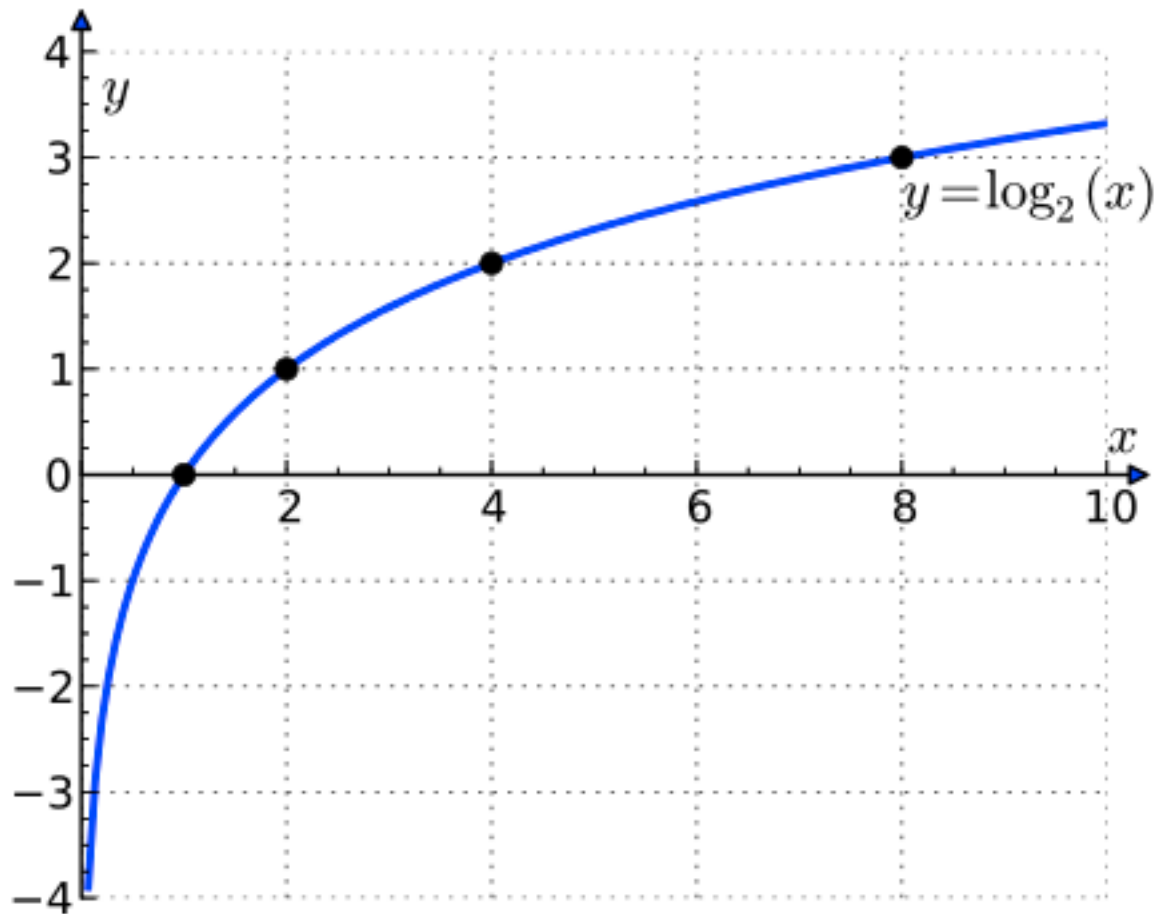
11 bits

52 bits

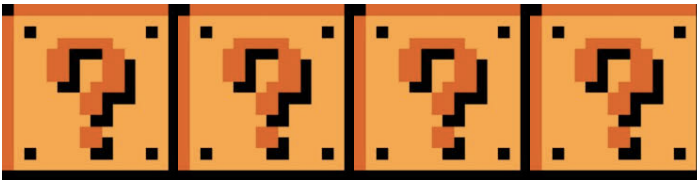
Largest =  $1.111\dots \times 2^{+1023}$

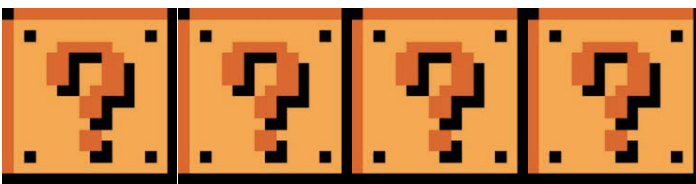
Smallest =  $1.000\dots \times 2^{-1024}$

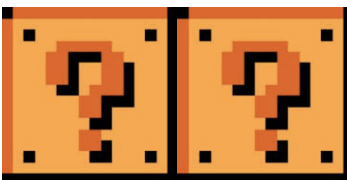
# Working with Probabilities in Log Space



# Log Identities (review)

$$\log(a \times b) =$$
A horizontal row of four orange Super Mario Bros. ? blocks, each with a black question mark and a small black dot in the center.

$$\log\left(\frac{a}{b}\right) =$$
A horizontal row of four orange Super Mario Bros. ? blocks, each with a black question mark and a small black dot in the center.

$$\log(a^n) =$$
A horizontal row of two orange Super Mario Bros. ? blocks, each with a black question mark and a small black dot in the center.



# Naïve Bayes with Log Probabilities

$$c_{MAP} = \operatorname{argmax}_c P(c|x_1, \dots, x_n)$$

$$= \operatorname{argmax}_c P(c) \prod_{i=1}^n P(x_i|c)$$

$$= \operatorname{argmax}_c \log \left( P(c) \prod_{i=1}^n P(x_i|c) \right)$$

because log is  
monotonic increasing

$$= \operatorname{argmax}_c \log P(c) + \sum_{i=1}^n \log P(x_i|c)$$

# Naïve Bayes with Log Probabilities

$$c_{MAP} = \operatorname{argmax}_c \log P(c) + \sum_{i=1}^n \log P(x_i|c)$$

Q: Why don't we have to worry about floating point underflow anymore?

# Working with Probabilities in Log Space

x	log(x)
0.000000000001	-11
0.00001	-5
0.0001	-4
0.001	-3
0.01	-2
0.1	-1

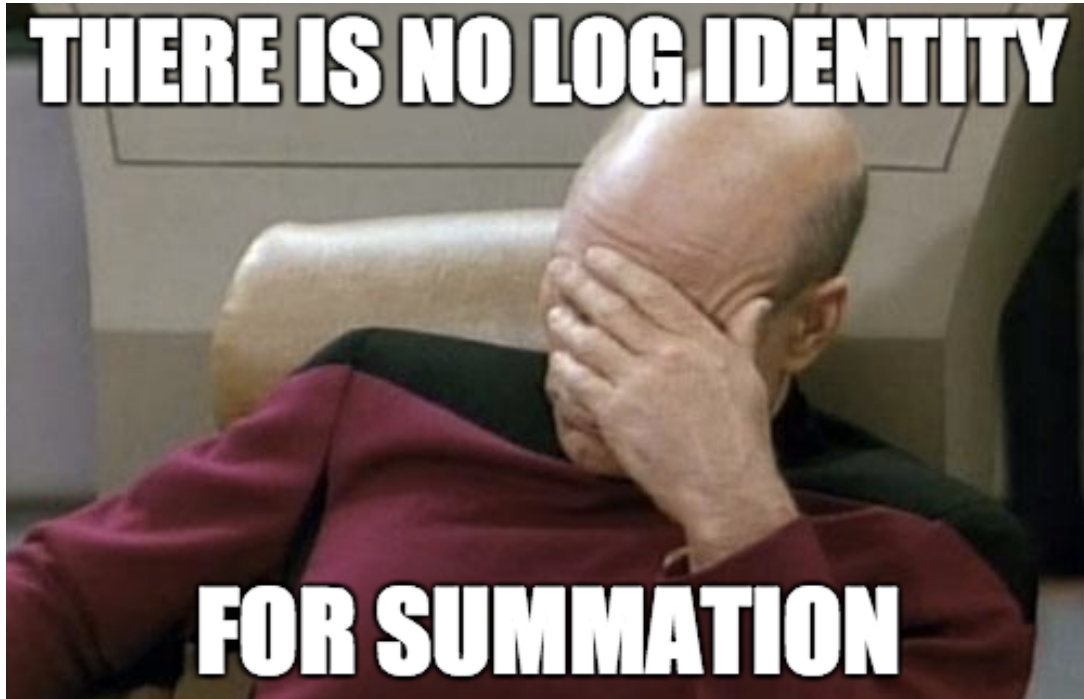
What if we want to calculate posterior log-probabilities?

$$P(c|x_1, \dots, x_n) = \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$\log P(c|x_1, \dots, x_n) = \log \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$= \log P(c) + \sum_{i=1}^n \log P(x_i|c) - \log \left[ \sum_{c'} P(c') \prod_{i=1}^n P(x_i|c') \right]$$

What if we want to calculate posterior log-probabilities?



$$P(c|x_1, \dots, x_n)$$

$$\frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$\log P(c|x_1, \dots, x_n)$$

$$\frac{\log P(c) + \sum_{i=1}^n \log P(x_i|c)}{\log \left[ \sum_{c'} P(c') \prod_{i=1}^n P(x_i|c') \right]}$$

$$= \log P(c) + \sum_{i=1}^n \log P(x_i|c) - \log \left[ \sum_{c'} P(c') \prod_{i=1}^n P(x_i|c') \right]$$

# Log Exp Sum Trick

- We have: a bunch of log probabilities:  
 $\log(p_1), \log(p_2), \log(p_3), \dots \log(p_n)$
- We want:  $\log(p_1 + p_2 + p_3 + \dots p_n)$
- We could convert back from log space, sum then take the log.

Q: Is this a good idea?

# Log Exp Sum Trick

- We have: a bunch of log probabilities:  
 $\log(p_1), \log(p_2), \log(p_3), \dots \log(p_n)$
- We want:  $\log(p_1 + p_2 + p_3 + \dots p_n)$
- We could convert back from log space, sum then take the log.

If the probabilities are very small, this will result in floating point underflow

# Log Exp Sum Trick

$$\log\left[\sum_i \exp(x_i)\right] = x_{max} + \log\left[\sum_i \exp(x_i - x_{max})\right]$$



# Another issue: Smoothing

$$\hat{P}(w_i|c) = \frac{\text{count}(w, c) + 1}{\sum_{w' \in V} \text{count}(w', c) + |V|}$$

# Another issue: Smoothing

Can think of alpha as a “pseudo-count”.  
Imaginary number of times this word has been seen.

$$\hat{P}(w_i|c) = \frac{\text{count}(w, c) + \alpha}{\sum_{w' \in V} \text{count}(w', c) + \alpha|V|}$$

# Another issue: Smoothing

Alpha doesn't necessarily need to be 1 (hyperparameter)

$$\hat{P}(w_i|c) = \frac{\text{count}(w, c) + \alpha}{\sum_{w' \in V} \text{count}(w', c) + \alpha|V|}$$

Hyperparameters are parameters that cannot be directly learned from the standard training process, and need to be predefined.

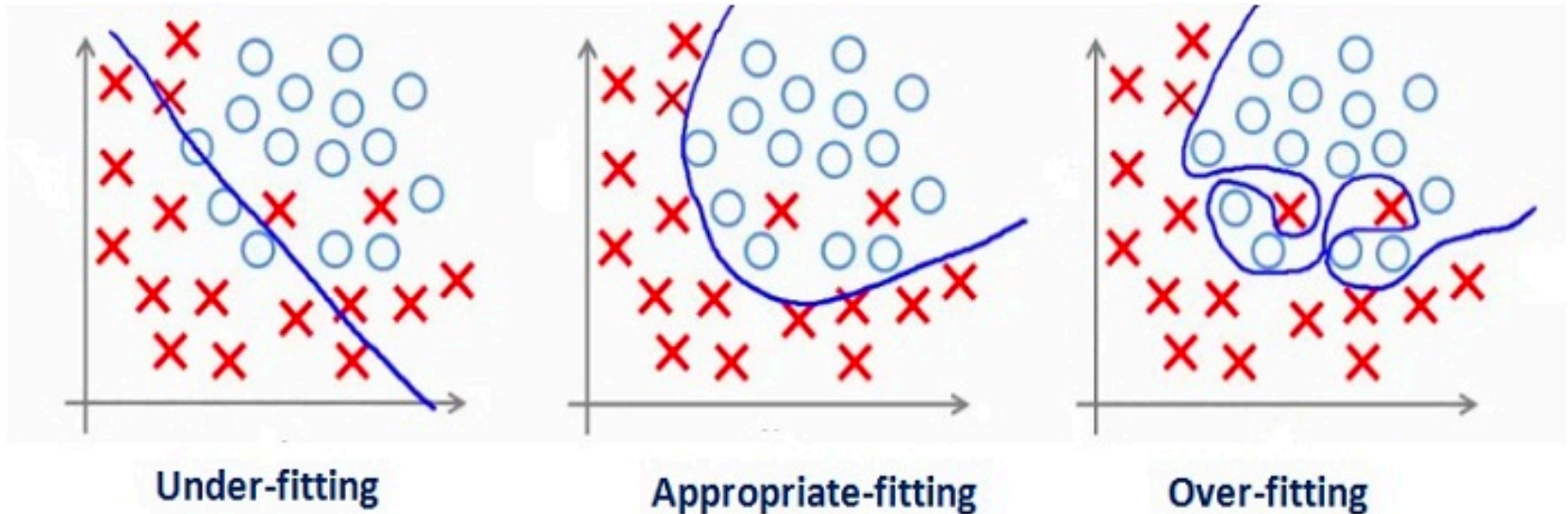
# Another issue: Smoothing

$$\hat{P}(w_i|c) = \frac{\text{count}(w, c) + \alpha}{\sum_{w' \in V} \text{count}(w', c) + \alpha|V|}$$

- What if alpha = 0?
- What if alpha = 0.000001?
- What happens as alpha gets very large?

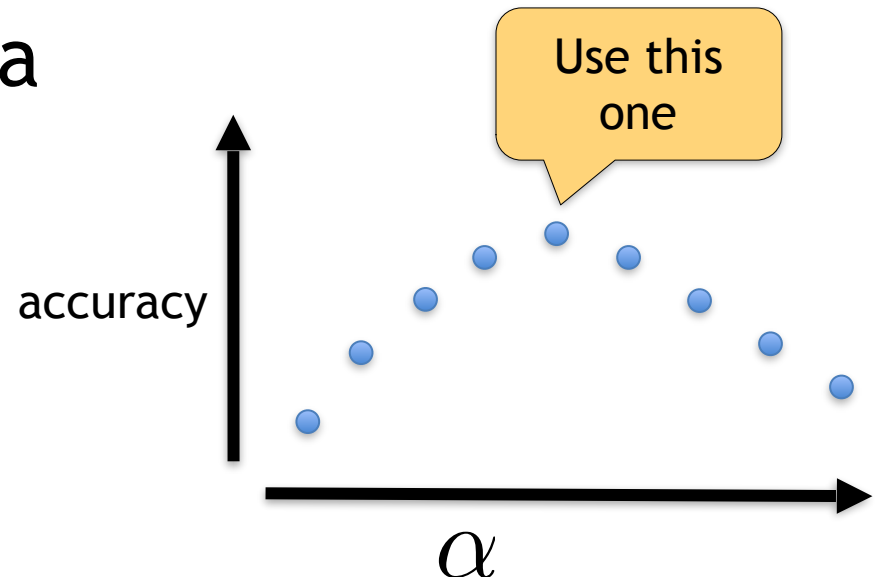
# Overfitting

- Model parameters fits the training data well, but generalize poorly to test data
- How to check for overfitting?
  - Training vs. test accuracy
- Pseudo-count parameter combats overfitting



# How to pick Alpha?

- Split train vs. test (dev)
- Try a bunch of different values
- Pick the value of alpha that performs best
- What values to try?  
Grid search
  - $(10^{-2}, 10^{-1}, \dots, 10^2)$



# Data Splitting

- Train vs. Test
- Better:
  - Train (used for fitting model **parameters**)
  - Dev (used for tuning **hyperparameters**)
  - Test (reserve for final evaluation)
- Cross-validation

# Feature Engineering

- What is your word / feature representation
  - Tokenization rules: splitting on whitespace?
  - Uppercase is the same as lowercase?
  - Numbers?
  - Punctuation?
  - Stemming?