



TEXT CLASSIFICATION

Sandeep Soni

09/25/2023

CLASS LOGISTICS

- Hw6 due

HW6 REVIEW

What are the similarities and differences between KMeans clustering and topic modeling?

HW6 REVIEW

KMeans

LDA

Similarities

Word lists are likely to be very similar

Differences

KMeans surfaces words that are more frequent

LDA surfaces words that are less frequent and more cohesive to a topic

QUESTION FOR THE DAY

“How to predict a label for given text?”

AGENDA

- Classification problem
- Naive Bayes
- Training/test setup

CLASSIFICATION

- Input: document (e.g., email)
- Output: label (e.g., spam/ham)



FORMAL TASK



FORMAL TASK

- x is input (e.g., an email)
- $x \in \mathcal{X}$ (e.g., set of emails)



FORMAL TASK

- x is input (e.g., an email)
 - $x \in \mathcal{X}$ (e.g., set of emails)
- y is output (e.g., spam)
 - $y \in \mathcal{Y}$ (e.g., {spam, ham})



FORMAL TASK

- x is input (e.g., an email)
 - $x \in \mathcal{X}$ (e.g., set of emails)
- y is output (e.g., spam)
 - $y \in \mathcal{Y}$ (e.g., {spam, ham})
- $y = h(x)$
 - h maps instances to labels



CLASSIFICATION

- The true mapping function h is not known to us so we want to find \hat{h} that's a **closest** approximation



RULE BASED CLASSIFICATION

- \hat{h} (“I wish to discuss personal investment business matters with you so as to be able to learn of the available investment opportunities in your region or country.”)

\hat{h}

if email contains phrase
“investment opportunities” then
spam



SUPERVISED LEARNING

- Learn \hat{h} from training data given in the form of $\langle x, y \rangle$ pairs



CLASSIFICATION PROBLEMS

Task	\mathcal{X}	\mathcal{Y}
Language ID	text	{english, mandarin, hindi, ...}
spam classification	email	{spam, ham}
party affiliation	speech	{republican, democrat}
sentiment analysis	text	{positive, negative, mixed, neutral}
music genre	lyric	{rock, pop, jazz, rap,...}

$$\hat{h}(x)$$

- We still need to resolve:
 - How to learn this mapping function? (e.g., which method to use, how to optimize, etc)
 - How to represent the input to this function?

NAIVE BAYES

- One simple yet quite effective classification method is Naive Bayes
 - Similar to LDA, it's a generative model
 - We'll represent input text as a bag of words vector

REFRESHER: CHAIN RULE

- If x and y are random variables, the joint probability can be factorized using chain rule

$$P(x, y) = P(y)P(x | y)$$

$$P(x, y) = P(x)P(y | x)$$

REFRESHER: CHAIN RULE

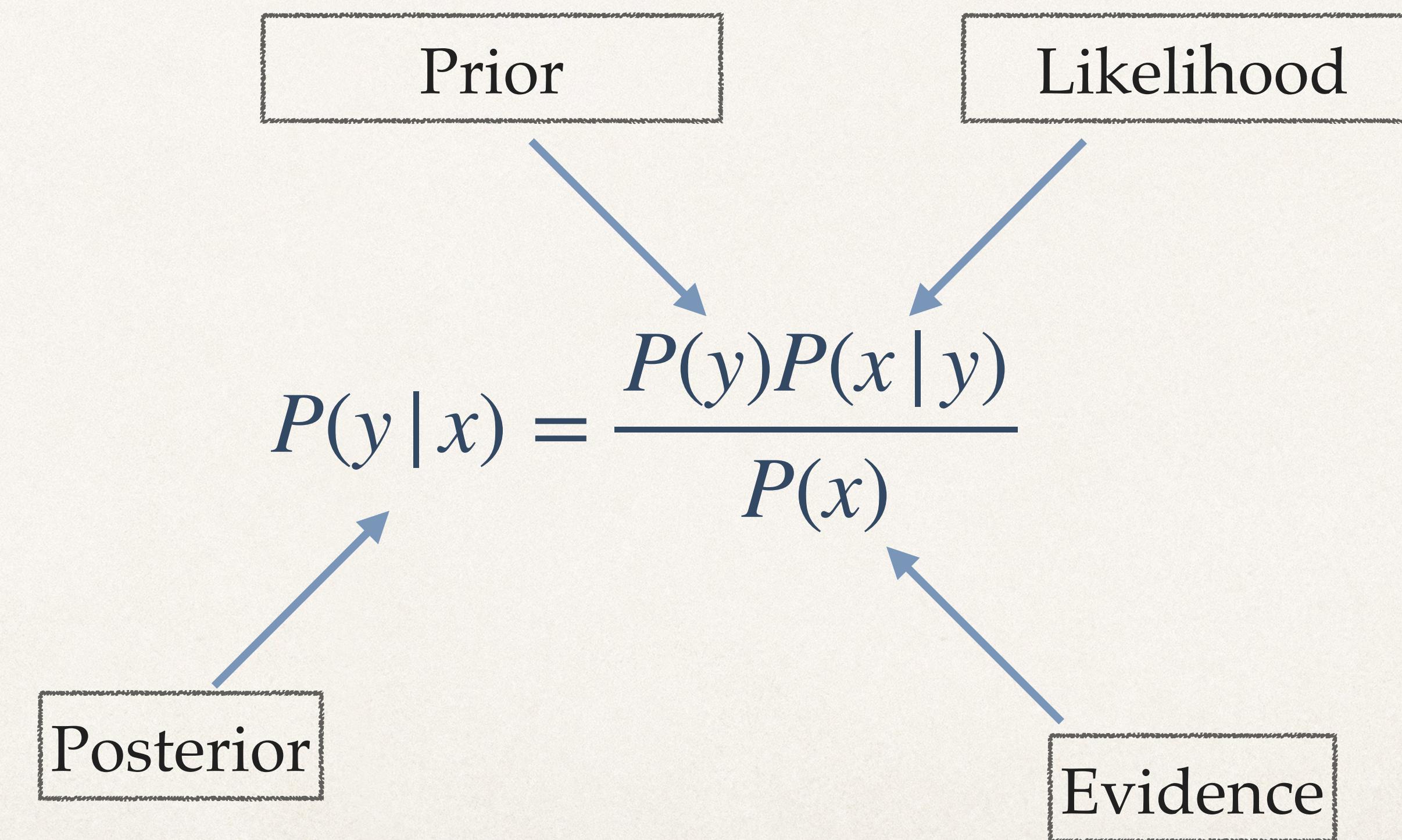
- If there are three variables, then chain rule gives:
$$P(x, y, z) = P(x)P(y | x)P(z | x, y)$$
- Just like two variables, permutations are equivalent

REFRESHER: CHAIN RULE

- In general, for multiple variables, chain rule is:

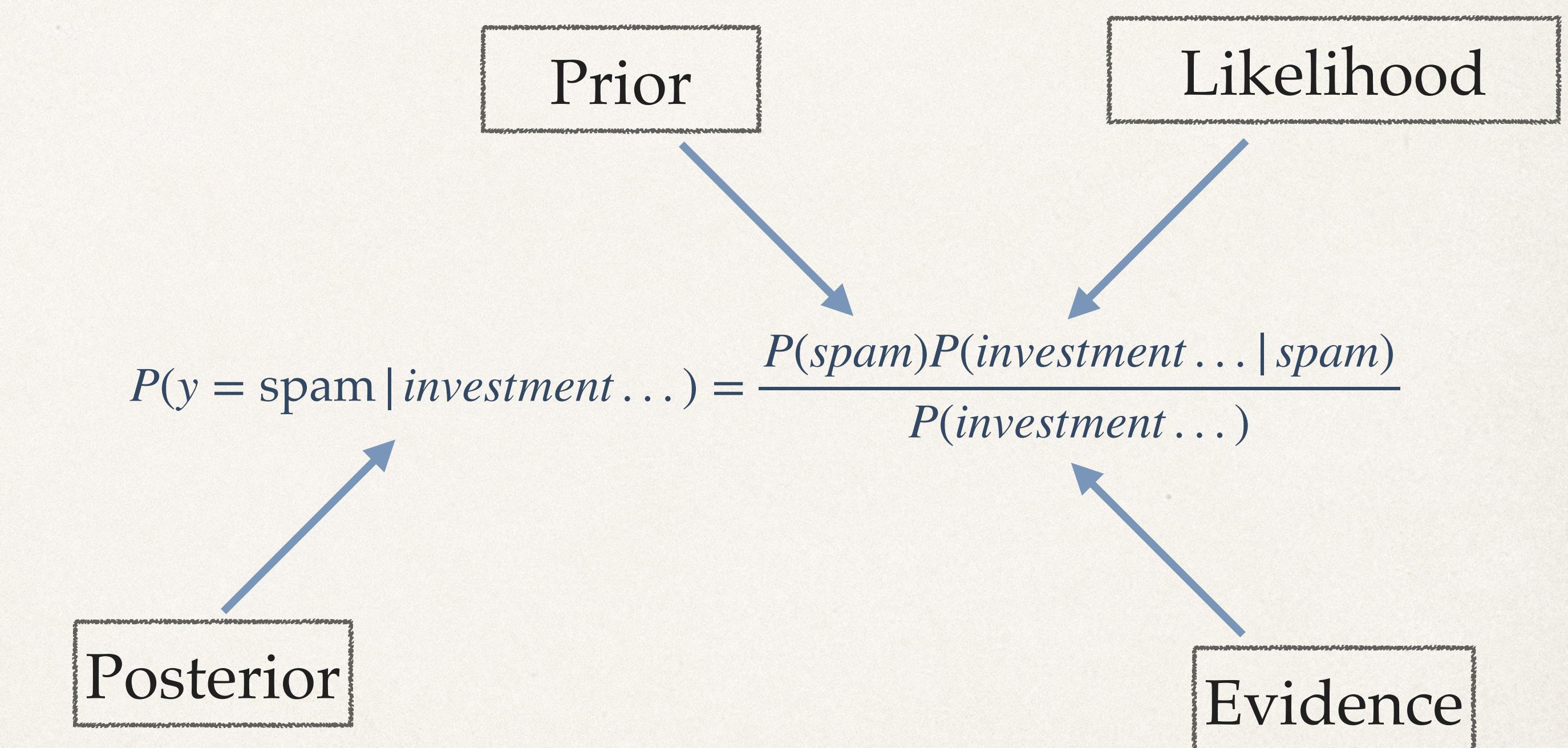
$$P(y, x_1, x_2, \dots, x_n) = P(y)P(x_1 | y)P(x_2 | x_1, y)\dots P(x_n | x_{n-1}, \dots, y)$$

REFRESHER: BAYES THEOREM



REFRESHER: BAYES THEOREM

- For spam classification:
- **Prior** how probable is to see the spam label
- **Likelihood** how likely are these words in the email for spam
- **Evidence** how probable are these words in the email
- **Posterior** how probable is spam label on this email



REFRESHER: BAYES THEOREM

- We can express the marginal probability in the denominator as:

$$P(y|x) = \frac{P(y)P(x|y)}{\sum_{y'} P(x,y')}$$

$$P(y|x) = \frac{P(y)P(x|y)}{\sum_{y'} P(y')P(x|y')}$$

TOWARDS CLASSIFICATION

- One way to learn the mapping $\hat{h}(x)$ is to learn the posterior distribution $P(y|x)$
- x is text, so not just one random variable but a bunch of variables, so we have to estimate $P(y|x_1, x_2, x_3, \dots, x_n)$

TOWARDS CLASSIFICATION

$$P(y | x_1, x_2, \dots, x_n) = \frac{P(y, x_1, x_2, \dots, x_n)}{P(x_1, x_2, \dots, x_n)}$$

TOWARDS CLASSIFICATION

- Joint probability in the numerator can be factorized using chain rule

$$P(y | x_1, x_2, \dots, x_n) = \frac{P(y)P(x_1, x_2, \dots, x_n | y)}{P(x_1, x_2, \dots, x_n)}$$

Continue to apply chain rule to the numerator (denominator is ignored for now)

$$P(y | x_1, x_2, \dots, x_n) \propto P(y)P(x_1, x_2, \dots, x_n | y)$$

$$\propto P(y)P(x_1 | y)P(x_2, \dots, x_n | y, x_1)$$

...

$$\propto P(y)P(x_1 | y)P(x_2 | x_1, y)P(x_3 | x_2, x_1, y) \dots P(x_n | x_{n-1}, \dots, x_1, y)$$

This is intractable!

Continue to apply chain rule to the numerator (denominator is ignored for now)

$$P(y | x_1, x_2, \dots, x_n) \propto P(y)P(x_1, x_2, \dots, x_n | y)$$

$$\propto P(y)P(x_1 | y)P(x_2, \dots, x_n | y, x_1)$$

...

$$\propto P(y)P(x_1 | y)P(x_2 | x_1, y)P(x_3 | x_2, x_1, y) \dots P(x_n | x_{n-1}, \dots, x_1, y)$$

This is intractable!

NAIVE BAYES

- Naive conditional independence assumption

$$P(x_i | x_{i-1}, x_{i-2}, \dots, y) = P(x_i | y)$$

- Given the category, the words (features) are independent of each other
- Under naive Bayes, $P(\text{"prince"} | \text{spam}) = P(\text{"prince"} | \text{spam, "kenyan"})$

Now we can rewrite the posterior probability as:

$$\begin{aligned} P(y | x_1, x_2, \dots, x_n) &\propto P(y)P(x_1 | y)P(x_2 | y)P(x_3 | y)\dots P(x_n | y) \\ &\propto P(y) \prod_{i=1}^n P(x_i | y) \end{aligned}$$

This is tractable!

We can estimate these probabilities by simply counting the instances in training data

$$P(y = \text{spam}) = \frac{\#\text{samples labeled spam}}{\#\text{ samples}}$$

$$P(\text{"kenyan"} | \text{spam}) = \frac{\#\text{samples labeled spam and contain "kenyan"}}{\#\text{samples labeled spam}}$$

PICKING THE LABEL

- Once you estimate the probabilities from training data, you can pick the label that maximizes the posterior

$$P(y = \text{spam} \mid \text{text}) > P(y = \text{ham} \mid \text{text})$$

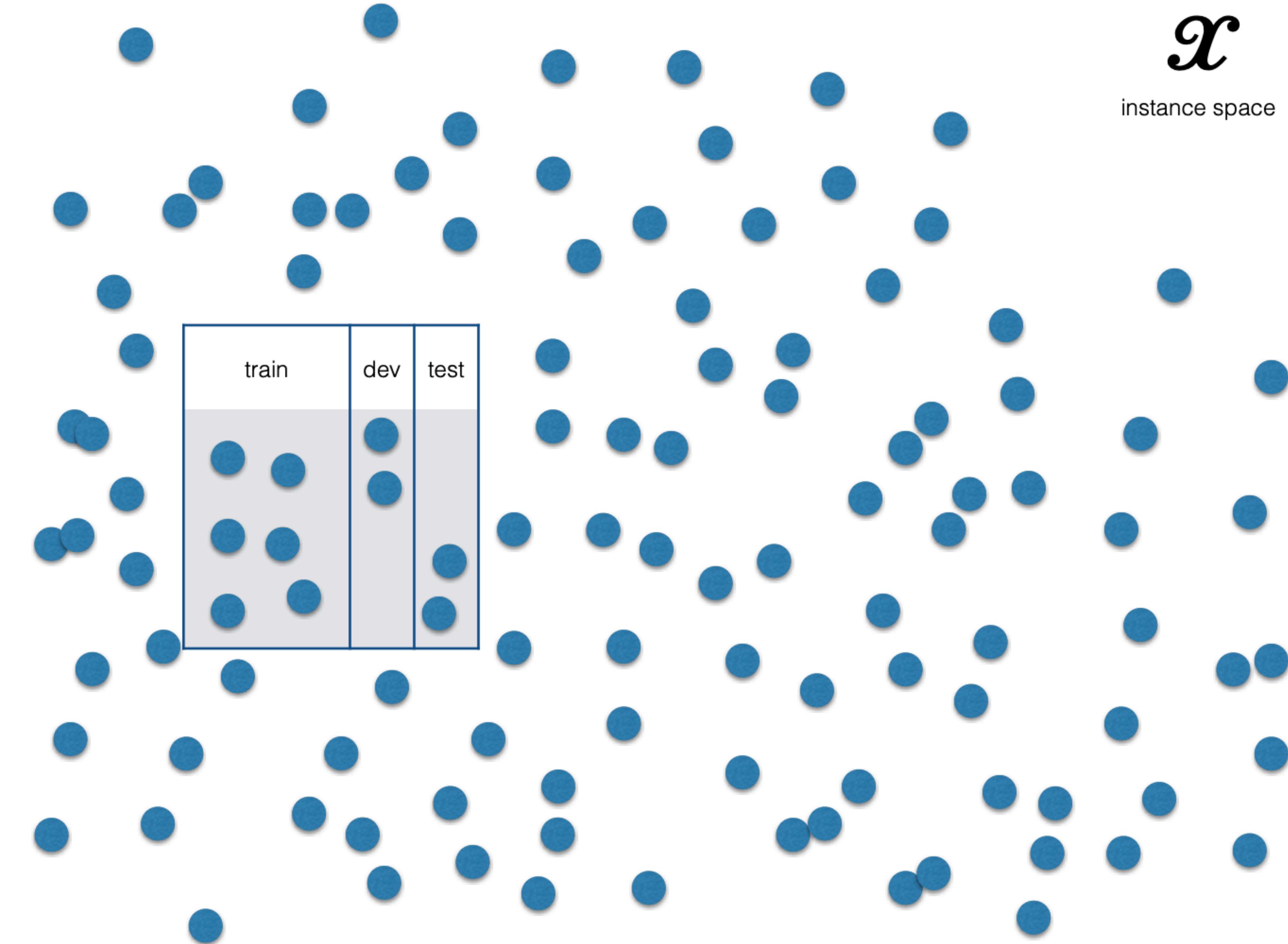
Spam

otherwise

Ham

\mathcal{X}

instance space



EXPERIMENT DESIGN

- Training set is to estimate parameters of the model
- Development set is to perform model selection
- Test set for evaluation

EXPERIMENT DESIGN

- Typically, we use 80% data for training, 10% for model selection and 10% for evaluation
- One should be careful never to use development or test data to do estimation

IN CLASS

- Naive Bayes Demo