

CPSC 340: Machine Learning and Data Mining

Generative Models

Admin

- **Assignment 0** was due last Wednesday.
 - Because of late days, we have to wait 3 days to post solutions.
 - At that time you will also gain read access to your classmates' work.
 - We voted on this during the first lecture.
 - No one approached me privately with objections.
- **Assignment 1** is out.
 - This is a representative assignment w.r.t. length/difficulty/format/style.
- **Registration:**
 - Keep checking your registration, it could change quickly.
 - As of last night, waitlist was down to 14 people.
- **Probability:**
 - If you are struggling with probability concepts towards the end of class today, check out the posted notes on probability.

Last Time: Training, Testing, and Validation

- Training step:

Input: set of ' n ' training examples x_i with labels y_i

Output: a model that maps from arbitrary x_i to a y_i

- Prediction step:

Input: set of ' t ' testing examples \hat{x}_i and a model.

Output: predictions \hat{y}_i for the testing examples.

- What we are interested in is the **test error**:

- Error made by prediction step on new data.

- Validation set or cross-validation can be used to estimate test error.

Should you trust them?

- Scenario 1:
 - “I built a model based on the data you gave me.”
 - “It classified your data with 98% accuracy.”
 - “It should get 98% accuracy on the rest of your data.”
- Probably not:
 - They are reporting training error.
 - This might have nothing to do with test error.
 - E.g., they could have fit a very deep decision tree.
- Why ‘probably’?
 - If they only tried a few very simple models, the 98% might be reliable.
 - E.g., they only considered decision stumps with simple 1-variable rules.

Should you trust them?

- Scenario 2:
 - “I built a model based on half of the data you gave me.”
 - “It classified the other half of the data with 98% accuracy.”
 - “It should get 98% accuracy on the rest of your data.”
- Probably:
 - They computed the validation error once.
 - This is an unbiased approximation of the test error.
 - Trust them if you believe they didn’t violate the golden rule.

Should you trust them?

- Scenario 3:
 - “I built 10 models based on half of the data you gave me.”
 - “One of them classified the other half of the data with 98% accuracy.”
 - “It should get 98% accuracy on the rest of your data.”
- Probably:
 - They computed the validation error a small number of times.
 - Maximizing over these errors is a biased approximation of test error.
 - But they only maximized it over 10 models, so bias is probably small.
 - They probably know about the golden rule.

Should you trust them?

- Scenario 4:
 - “I built 1 billion models based on half of the data you gave me.”
 - “One of them classified the other half of the data with 98% accuracy.”
 - “It should get 98% accuracy on the rest of your data.”
- Probably not:
 - They computed the validation error a huge number of times.
 - Maximizing over these errors is a biased approximation of test error.
 - They tried so many models, one of them is likely to work by chance.
 - This is the “multiple comparisons problem” in statistics
- Why ‘probably’?
 - If the 1 billion models were all extremely-simple, 98% might be reliable.

Should you trust them?

- Scenario 5:
 - “I built 1 billion models based on the first third of the data you gave me.”
 - “One of them classified the second third of the data with 98% accuracy.”
 - “It also classified the last third of the data with 98% accuracy.”
 - “It should get 98% accuracy on the rest of your data.”
- Probably:
 - They computed the first validation error a huge number of times.
 - But they had a second validation set that they only looked at once.
 - The second validation set gives unbiased test error approximation.
 - This is ideal, as long as they didn't violate golden rule on second set.
 - And assuming you are using IID data in the first place.

The 'Best' Machine Learning Model

- Decision trees are not always most accurate.
- What is the 'best' machine learning model?
- First we need to define generalization error:
 - Test error on new examples (excludes test examples seen during training).
- No free lunch theorem:
 - There is **no** 'best' model achieving the best generalization error for every problem.
 - If model A generalizes better to new data than model B on one dataset, there is another dataset where model B works better.
- This question is like asking which is 'best' among "rock", "paper", and "scissors".

The 'Best' Machine Learning Model

- Implications of the lack of a 'best' model:
 - We need to learn about and **try out multiple models**.
- So which ones to study in CPSC 340?
 - We'll usually motivate a method by a specific application.
 - But we'll focus on **models that are effective in many applications**.
- Caveat of no free lunch (NFL) theorem:
 - The world is very structured.
 - **Some datasets are more likely than others**.
 - Model A really could be better than model B on every real dataset in practice.
- Machine learning research:
 - Large focus on models that are **useful across many applications**.

Application: E-mail Spam Filtering

- Want a build a system that filters spam e-mails.

<input type="checkbox"/>			Jannie Keenan	ualberta	You are owed \$24,718.11
<input type="checkbox"/>			Abby	ualberta	USB Drives with your Logo
<input type="checkbox"/>			Rosemarie Page		Re: New request created with ID: ##62
<input type="checkbox"/>			Shawna Bulger		RE: New request created with ID: ##63
<input type="checkbox"/>			Gary	ualberta	Cooperation



Gary <jaiwasie@mail.com>

to schmidt ▾



Be careful with this message. Similar messages were used to steal people's personal information. [Learn more](#)

Hey,

Do you have a minute today?

Are you interested to use our email marketing and lead generation solutions?

We have worked on a number of projects and campaigns in many industries since 2007

Please reply today so we can go over options for you.

I am sure we can help to grow your business soon by using our mailing services.

Best regards,

Gary

Contact: abelfong@sina.com

- We have a big collection of e-mails, labeled by users.
- Can we formulate as supervised learning?

First a bit more supervised learning notation

- We have been using the notation 'X' and 'y' for supervised learning:

The diagram shows two mathematical structures: a matrix X and a vector y . The matrix X is represented by large square brackets containing a horizontal red oval. An arrow points from this oval to the label x_i . Below the matrix, a red oval highlights a specific element, with an arrow pointing to the label x_{ij} . The vector y is represented by large square brackets containing a single red oval. An arrow points from this oval to the label y_i .

- X is matrix of all features, y is vector of all labels.
- Need a way to refer to the features and label of **specific object 'i'**.
 - We use y_i for the label of object 'i' (element 'i' of 'y').
 - We use x_i for the features object 'i' (row 'i' of 'X').
 - We use x_{ij} for feature 'j' of object 'i'.

Feature Representation for Spam

- How do we make label ' y_i ' of an individual e-mail?
 - ($y_i = 1$) means 'spam', ($y_i = 0$) means 'not spam'.
- How do we construct features ' x_i ' for an e-mail?
 - Use **bag of words**:
 - "hello", "vicodin", "\$".
 - "vicodin" feature is 1 if "vicodin" is in the message, and 0 otherwise.
 - Could add phrases:
 - "be your own boss", "you're a winner", "CPSC 340".
 - Could add regular expressions:
 - <recipient>, <sender domain == "mail.com">

Probabilistic Classifiers

- For years, best spam filtering methods used **naïve Bayes**.
 - Naïve Bayes is a **probabilistic** classifier based on **Bayes rule**.
 - It's "naïve" because it makes a **strong conditional independence assumption**.
 - But it tends **to work well with bag of words**.
- Probabilistic classifiers model the **conditional probability**, $p(y_i | x_i)$.
 - "If a message has words x_i , what is probability that message is spam?"
- If $p(y_i = \text{'spam'} | x_i) > p(y_i = \text{'not spam'} | x_i)$, classify as spam.

- Recall our spam filtering setup:
 - y_i : whether or not the e-mail was spam.
 - x_i : the **set of words/phrases/expressions** in the e-mail.
- To model conditional probability, **naïve Bayes** uses Bayes rule:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- Easy part #1: $p(y_i = \text{'spam'})$ is the **probability that an e-mail is spam**.
 - Count of number of times ($y_i = \text{'spam'}$) divided by number of objects 'n'.

- Recall our spam filtering setup:
 - y_i : whether or not the e-mail was spam.
 - x_i : the **set of words/phrases/expressions** in the e-mail.
- To model conditional probability, **naïve Bayes** uses Bayes rule:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- Easy part #2: We **don't need** $p(x_i)$.

To test $p(y_i = \text{"spam"} | x_i)$ we just need to know if $p(y_i = \text{"spam"} | x_i) > p(y_i = \text{"not spam"} | x_i)$.

By Bayes rule this is equivalent to $\frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)} > \frac{p(x_i | y_i = \text{"not spam"}) p(y_i = \text{"not spam"})}{p(x_i)}$

Denominators are the same so we just test $p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"}) > p(x_i | y_i = \text{"not spam"}) p(y_i = \text{"not spam"})$

Generative Classifiers

- The **hard part is estimating $p(x_i \mid y_i = \text{'spam'})$** :
 - the probability of seeing the words/expressions x_i if the e-mail is spam.
- Classifiers based on Bayes rule are called **generative classifier**:
 - It needs to know the **probability of the features, given the class**.
 - How to “generate” features.
 - You need a model that knows what spam messages look like.
 - And a second that knows what non-spam messages look like.
 - This work well with **tons of features compared to number of objects**.

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} \mid x_i) = \frac{p(x_i \mid y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- What do these terms mean?

ALL E-MAILS
(including duplicates)

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} \mid x_i) = \frac{p(x_i \mid y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i)$ is probability that a random e-mail has features x_i .

ALL E-MAILS
(including duplicates)

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} \mid x_i) = \frac{p(x_i \mid y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i)$ is probability that a random e-mail has features x_i .

ALL E-MAILS
(including duplicates)

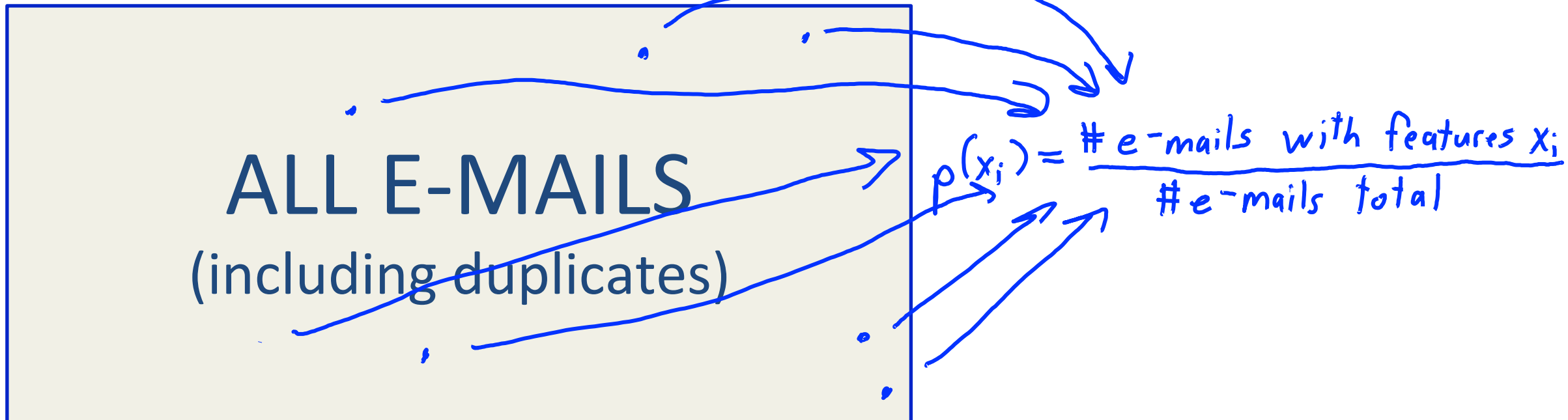
$$p(x_i) = \frac{\# \text{ e-mails with features } x_i}{\# \text{ e-mails total}}$$

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i)$ is probability that a random e-mail has features x_i .



Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} \mid x_i) = \frac{p(x_i \mid y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i)$ is probability that a random e-mail has features x_i .

ALL E-MAILS
(including duplicates)

$$p(x_i) = \frac{\# \text{e-mails with features } x_i}{\# \text{e-mails total}}$$

- Hard, but not needed to classify using:
 $p(y_i = \text{'spam'} \mid x_i) > p(y_i = \text{'not spam'} \mid x_i)$

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(y_i = \text{'spam'})$ is probability that a random e-mail is spam.



$$p(y_i = \text{"spam"}) = \frac{\# \text{ spam messages}}{\# \text{ total messages}}$$

- Hard to compute exactly.
- But is easy to approximate from data:
 - Count (#spam in data)/(#messages)

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i | y_i = \text{'spam'})$ is probability that spam has features x_i .



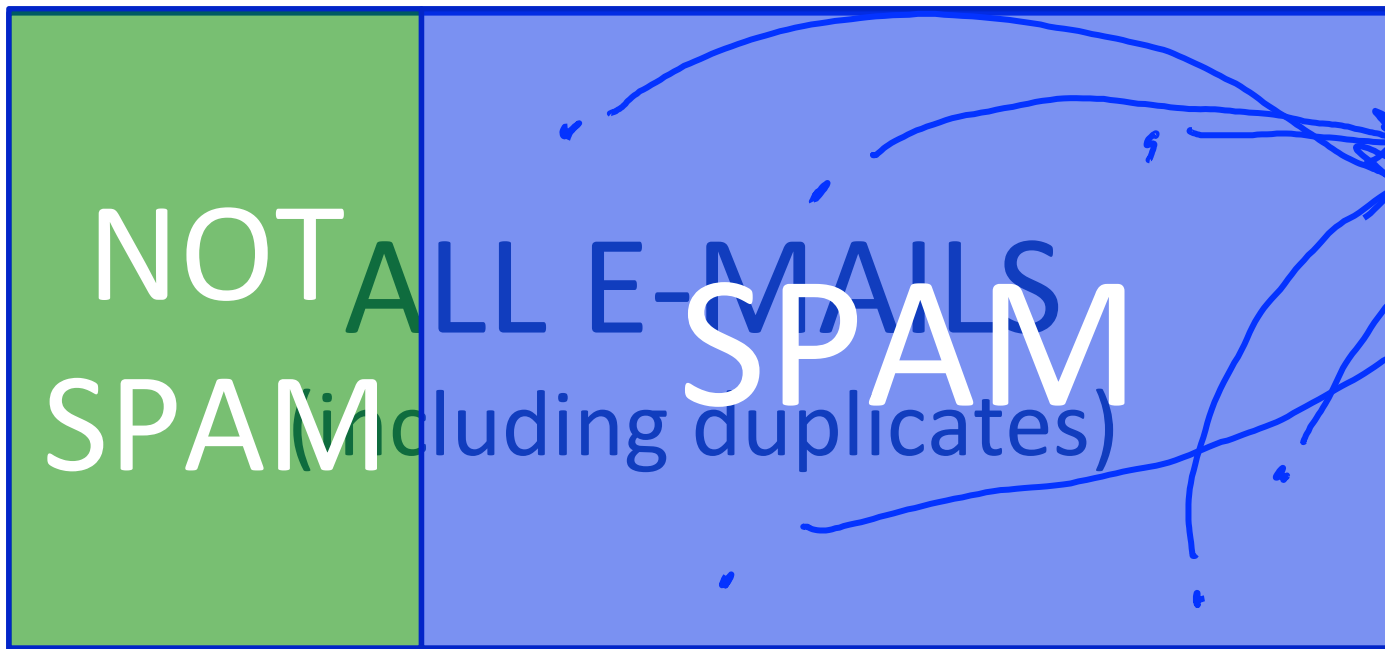
$$p(x_i | y_i = \text{"spam"}) = \frac{\# \text{ spam messages with features } x_i}{\# \text{ spam messages}}$$

Generative Models

- Spam filtering methods based on **generative models**:

$$p(y_i = \text{"spam"} | x_i) = \frac{p(x_i | y_i = \text{"spam"}) p(y_i = \text{"spam"})}{p(x_i)}$$

- $p(x_i | y_i = \text{'spam'})$ is probability that spam has features x_i .



- Very hard to estimate:**
 - Too many possible x_i .

Naïve Bayes

- How the naïve Bayes model deals with the hard terms:

$$p(\text{spam} | \text{hello}, \text{vicodin}, \text{CPSC 340}) = \frac{p(\text{hello}, \text{vicodin}, \text{CPSC 340} | \text{spam}) p(\text{spam})}{p(\text{hello}, \text{vicodin}, \text{CPSC 340})} \quad (\text{Bayes rule})$$

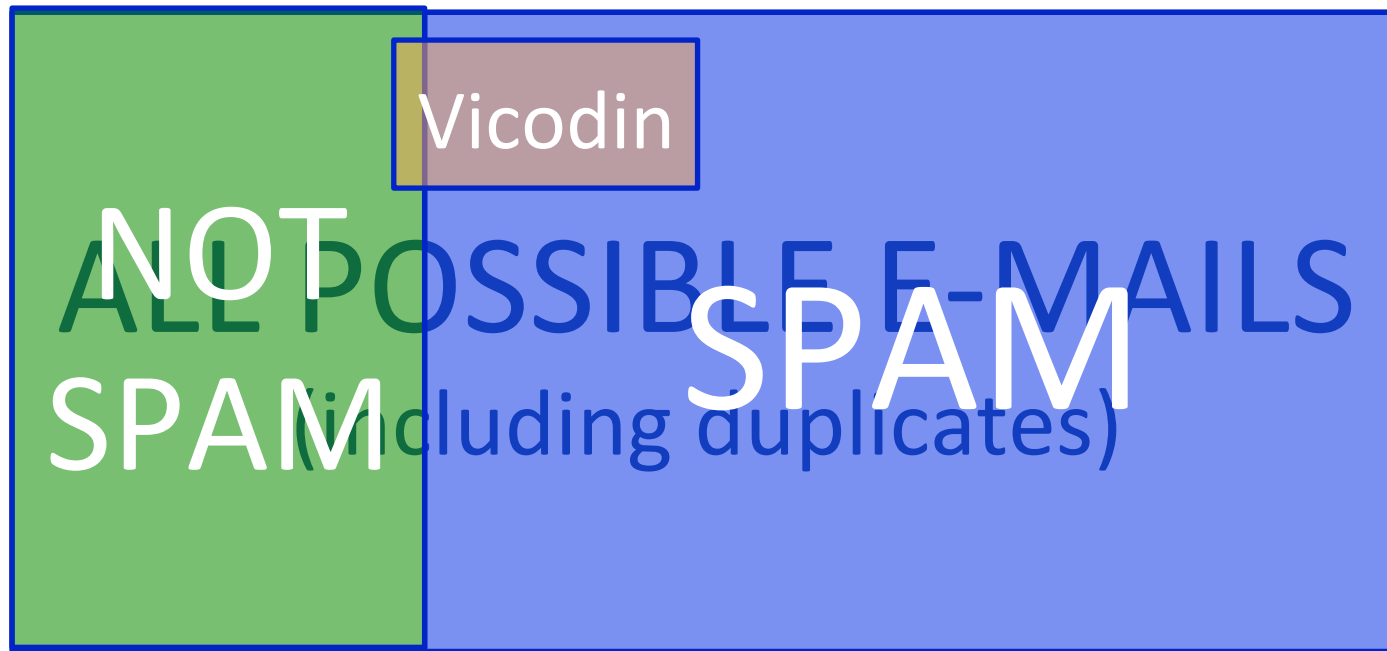
"equal up to constant not depending on spam" \propto $\underbrace{p(\text{hello}, \text{vicodin}, \text{CPSC 340} | \text{spam})}_{\text{HARD}} \underbrace{p(\text{spam})}_{\text{easy}}$

(naive Bayes assumption) $\approx \underbrace{p(\text{hello} | \text{spam})}_{\text{easy}} \underbrace{p(\text{vicodin} | \text{spam})}_{\text{easy}} \underbrace{p(\text{CPSC 340} | \text{spam})}_{\text{easy}} \underbrace{p(\text{spam})}_{\text{easy}}$

- Now **only** need easy quantities like $p(\text{'vicodin'} = 1 | y_i = \text{'spam'})$.

Naïve Bayes Models

- $p(\text{vicodin} = 1 \mid \text{spam} = 1)$ is probability of seeing 'vicodin' in spam.



$$p(\text{vicodin}=1 \mid \text{spam}=1) = \frac{\# \text{ spam messages w/ vicodin}}{\# \text{ spam messages}}$$

- Easy to estimate:
 - $\#(\text{spam w/ Vicodin}) / \# \text{spam}$
 - "Maximum likelihood estimate"

Naïve Bayes

- Naïve Bayes more formally:

$$p(y_i | x_i) = \frac{p(x_i | y_i) p(y_i)}{p(x_i)}$$

$$\propto p(x_i | y_i) p(y_i)$$

$$\approx \prod_{j=1}^d [p(x_{ij} | y_i)] p(y_i)$$

- Assumption: all x_i are conditionally independent given y_i .

Independence of Random Variables

- **Events** A and B are **independent** if $p(A, B) = p(A)p(B)$.
 - Equivalently: $p(A | B) = p(A)$.
 - “Knowing B happened tells you nothing about A”.
 - We use the notation:
$$A \perp B$$
- **Random variables** are **independent** if $p(x, y) = p(x)p(y)$ for all x and y.
 - Flipping two coins:
$$p(C_1 = \text{'heads'}, C_2 = \text{'heads'}) = p(C_1 = \text{'heads'})p(C_2 = \text{'heads'})$$
$$p(C_1 = \text{'tails'}, C_2 = \text{'heads'}) = p(C_1 = \text{'tails'})p(C_2 = \text{'heads'})$$
$$\dots$$

Conditional Independence

- A and B are **conditionally independent** given C if
$$p(A, B \mid C) = p(A \mid C)p(B \mid C).$$
 - Equivalently: $p(A \mid B, C) = p(A \mid C)$.
 - “Knowing C happened, also knowing B happened says nothing about A”.
 - Example: $p(\text{Pizza} \mid D_1, \text{Survive}) = p(\text{Pizza} \mid \text{Survive})$.
 - Knowing you survived, dice 1 gives no information about chance of pizza.
 - We use the notation: $A \perp B \mid C$
- Semantics of $p(A, B \mid C, D)$:
 - “probability of A and B happening, if we know that C and D happened”.

Naïve Bayes

- In naïve Bayes: assume **features are independent given label**.
 - “Once you know it’s spam, there is no dependency between features.”
 - Not true, but sometimes a good approximation.

Training:

1. Set n_c to the number of times $(y_i = c)$.

2. Estimate $p(y = c)$ as $\frac{n_c}{n}$.

3. Set n_{cjk} as the number of times $(y_i = c, X_{ij} = k)$

4. Estimate $p(x_i = k \mid y = c) = \frac{p(x_i = k, y = c)}{p(y = c)}$ as $\frac{\frac{n_{cjk}}{n}}{\frac{n_c}{n}} = \frac{n_{cjk}}{n_c}$.

Naïve Bayes

- In naïve Bayes: assume **features are independent given label**.
 - “Once you know it’s spam, there is no dependency between features.”
 - Not true, but sometimes a good approximation.

Prediction:

Given a new example x_i we want to find the 'c' maximizing $p(x_i | y_i)$.

Under the naive Bayes assumption we thus maximize

$$p(y=c | x_i) \propto \prod_{j=1}^d [p(x_{ij} | y=c)] p(y=c)$$

Summary

- **No free lunch theorem**: there is no “best” ML model.
- **Joint probability**: probability of A and B happening.
- **Conditional probability**: probability of A if we know B happened.
- **Generative classifiers**: build a probability of seeing the features.
- Next time:
 - A “best” machine learning model as ‘n’ goes to ∞ .



- All the remaining slides are “bonus”.
- We may go through them briefly, if time permits.

Generative Classifiers

- But does it need to know language to model $p(x_i | y_i)$???
- To fit generative models, usually make BIG assumptions:
 - Gaussian discriminant analysis (GDA):
 - Assume that $p(x_i | y_i)$ follows a multivariate normal distribution.
 - Naïve Bayes (NB):
 - Assume that each variables in x_i is independent of the others in x_i given y_i .