

# Naïve Bayes Classifiers

# Review

- Let event  $D$  = data we have observed.
- Let events  $H_1, \dots, H_n$  be events representing the  $n$  hypotheses we want to choose between.
- Use  $D$  to pick the "best"  $H$ .
- There are two "standard" ways to do this, depending on what information we have available.

# Maximum likelihood hypothesis

- The maximum likelihood hypothesis ( $H^{\text{ML}}$ ) is the hypothesis that **maximizes the probability of the data given that hypothesis.**

$$H^{\text{ML}} = \underset{i}{\operatorname{argmax}} P(D | H_i)$$

- How to use it: compute  $P(D | H_i)$  for each hypothesis (1 through  $n$ ) and select the one with the greatest value.

# Maximum a posteriori (MAP) hypothesis

- The MAP hypothesis is the hypothesis that **maximizes the posterior probability:**

$$\begin{aligned} H^{\text{MAP}} &= \operatorname{argmax}_i P(H_i | D) \\ &= \operatorname{argmax}_i \frac{P(D | H_i)P(H_i)}{P(D)} \\ &\propto \operatorname{argmax}_i P(D | H_i)P(H_i) \end{aligned}$$

- The  $P(D | H_i)$  terms are now *weighted* by the hypothesis prior probabilities.

# Posterior probability

- If you need the actual posterior probability for some hypothesis  $H_i$ :

$$\begin{aligned} P(H_i | D) &= \frac{P(D | H_i)P(H_i)}{P(D)} \\ &= \frac{P(D | H_i)P(H_i)}{\sum_j P(D, H_j)} \\ &= \frac{P(D | H_i)P(H_i)}{\sum_j P(D | H_j)P(H_j)} \end{aligned}$$

# Combining evidence

- If we have multiple pieces of data/evidence (say two pieces), then we need to compute or estimate

$$P(D_1, D_2 \mid H_i)$$

which is often hard.

- Instead, we **assume** all pieces of evidence are conditionally independent given a hypothesis:

$$P(D_1, D_2 \mid H_i) = P(D_1 \mid H_i)P(D_2 \mid H_i)$$

- This assumption is **most likely not true**, but we do it to make our lives easier.

# Combining evidence ( $m$ pieces)

$$\begin{aligned} P(H_i | D_1, \dots, D_m) &= \frac{P(D_1, \dots, D_m | H_i)P(H_i)}{P(D_1, \dots, D_m)} \\ &= \frac{\left[ P(D_1 | H_i) \cdots P(D_m | H_i) \right] P(H_i)}{P(D_1, \dots, D_m)} \\ &= \frac{\left[ \prod_{j=1}^m P(D_j | H_i) \right] P(H_i)}{P(D_1, \dots, D_m)} \end{aligned}$$

where

$$P(D_1, \dots, D_m) = \sum_{k=1}^n \left( \left[ \prod_{j=1}^m P(D_j | H_k) \right] P(H_k) \right)$$

# Classification

- Classification is the problem of identifying which of a set categories (called classes) a particular item belongs in.
- Lots of real-world problems are classification problems:
  - spam filtering (classes: spam/not-spam)
  - handwriting recognition & OCR (classes: one for each letter, number, or symbol)
  - text classification, image classification, music classification, etc.
- Almost any problem where you are assigning a label to items can be set up as a classification task.



# Classification

- An algorithm that does classification is called a **classifier**. Classifiers take an item as input and output the class it thinks that item belongs to. That is, the classifier *predicts* a class for each item.
- Lots of classifiers are based on probabilities and statistical inference:
  - The **classes** become the hypotheses being tested.
  - The item being classified is turned into a collection of data called **features**. Useful features are attributes of the item that are strongly correlated with certain classes.
  - The classification algorithm is usually ML or MAP, depending on what data we have available.

# Example: Spam classification

- New email arrives: is it spam or not spam?
- A useful set of features might be the presence or absence of various words in the email:
  - F1,  $\sim$ F1: the word "luxury" appears/does not appear
  - F2,  $\sim$ F2: the word "save" appears/does not appear
  - F3,  $\sim$ F3: the word "brands" appears/does not appear
- Let's say our new email contains "luxury" and "brands," but not "save."
- The features for this email are F1,  $\sim$ F2, and F3.
- Let's use MAP for classification.

# Example: Spam classification

- Features = Data =  $D = F_1, \sim F_2, F_3$ .

$$H^{\text{MAP}} = \underset{i}{\operatorname{argmax}} P(D \mid H_i) P(H_i)$$

$$H^{\text{MAP}} = \underset{i \in \{\text{spam}, \text{not-spam}\}}{\operatorname{argmax}} P(F_1, \neg F_2, F_3 \mid H_i) P(H_i)$$

# Example: Spam classification

- Features = Data =  $D = F_1, \sim F_2, F_3$ .

$$H^{\text{MAP}} = \underset{i}{\operatorname{argmax}} P(D \mid H_i) P(H_i)$$

$$H^{\text{MAP}} = \underset{i \in \{\text{spam}, \text{not-spam}\}}{\operatorname{argmax}} P(F_1, \neg F_2, F_3 \mid H_i) P(H_i)$$

- Let's assume all the features are conditionally independent given the hypothesis.

$$H^{\text{MAP}} = \underset{i \in \{\text{spam}, \text{not-spam}\}}{\operatorname{argmax}} P(F_1 \mid H_i) P(\neg F_2 \mid H_i) P(F_3 \mid H_i) P(H_i)$$

# Naïve Bayes

- This is called a Naïve Bayes classifier.
- Assumes the data is a collection of features, and each feature is conditionally independent of all other features given the hypothesis.
- Classifies using MAP hypothesis.

# Naïve Bayes

- Hypotheses:  $H_1$  through  $H_n$ .
- Features (data):  $F_1$  through  $F_m$ .

$$\begin{aligned} H^{\text{MAP}} &= \operatorname{argmax}_i P(D \mid H_i) P(H_i) \\ &= \operatorname{argmax}_i P(F_1, \dots, F_m \mid H_i) P(H_i) \\ &= \operatorname{argmax}_i \left[ P(F_1 \mid H_i) \cdots P(F_m \mid H_i) \right] P(H_i) \\ &= \operatorname{argmax}_i \left[ \prod_{j=1}^m P(F_j \mid H_i) \right] P(H_i) \end{aligned}$$

# Probabilities needed

- $P(H_i)$  for  $i = 1$  to  $n$ .
- $P(F_j | H_i)$  for  $j = 1$  to  $m$  and  $i = 1$  to  $n$ .

# Example

- Suppose I know 80% of my email is spam.
- I have three features, "luxury," "brands," and "save."
- I know:
  - $P(\text{luxury} \mid \text{spam}) = 0.4$        $P(\text{luxury} \mid \sim\text{spam}) = 0.01$
  - $P(\text{brands} \mid \text{spam}) = 0.3$        $P(\text{brands} \mid \sim\text{spam}) = 0.2$
  - $P(\text{save} \mid \text{spam}) = 0.4$        $P(\text{save} \mid \sim\text{spam}) = 0.1$
- Suppose a new, incoming email contains "luxury" and "save" but not "brands." Should it be classified as spam or  $\sim\text{spam}$ ?



# Learning probabilities from data

- To use MAP, we need to calculate or estimate  $P(H_i)$  and  $P(F_1, \sim F_2, F_3 \mid H_i)$  for each  $i$ .
- In other words, we need to know:
  - $P(\text{spam})$
  - $P(\text{not-spam})$
  - $P(F_1, \sim F_2, F_3 \mid \text{spam})$
  - $P(F_1, \sim F_2, F_3 \mid \text{not-spam})$
- In the previous example, these were given to us, but what if they weren't?

# Learning probabilities from data

- Let's assume we have access to a large number of old emails that are correctly labeled as spam/not-spam.
- How can we estimate  $P(\text{spam})$ ?

$$P(\text{spam}) = \frac{\# \text{ of emails labeled as spam}}{\text{total } \# \text{ of emails}}$$

# Learning probabilities from data

- Let's assume we have access to a large number of old emails that are correctly labeled as spam/not-spam.
- How can we estimate  $P(F_1, \sim F_2, F_3 \mid \text{spam})$ ?

$$P(F_1, \neg F_2, F_3 \mid \text{spam}) = \frac{\# \text{ of spam emails with those exact features}}{\text{total } \# \text{ of spam emails}}$$

- Why is this probably going to be a very rough estimate?

## Conditional independence to the rescue!

- It is unlikely that our set of old emails contains many messages with that exact set of features.
- Let's assume that all of our features are conditionally independent of each other, given the hypothesis (spam/not-spam).

$$P(F_1, \neg F_2, F_3 \mid \text{spam}) =$$

$$P(F_1 \mid \text{spam}) \cdot P(\neg F_2 \mid \text{spam}) \cdot P(F_3 \mid \text{spam})$$

- These probabilities are easier to get good estimates for!

# Learning probabilities from data

- So now we need to estimate  $P(F_1 \mid \text{spam})$  instead of  $P(F_1, \sim F_2, F_3 \mid \text{spam})$ .
- Equivalently, how can we estimate the probability of seeing "luxury" in an email given that the email is spam?
- $P(F_1 \mid \text{spam}) = \frac{\text{\# of spam emails with "luxury"}}{\text{total \# of spam emails}}$

# Example

Suppose I have 20 emails that have been already classified into spam (15 emails) and non-spam (5 emails). Suppose I only care about the presence or absence of the words **luxury**, **brands**, and **save**.

Suppose 6 of the spam emails contain "luxury," 3 of the spam emails contain "brands," and 7 of the spam emails contain "save."

Suppose 1 of the non-spam emails contains "luxury," 2 of the non-spam emails contain "brands," and 2 of the non-spam emails contain "save."

Suppose a new email arrives that contains the words "luxury" and "save" but not "brands." Should this be classified as spam or not spam?

# Another problem to handle...

- What if we a word *never* appears in any spam emails? What happens to its probability estimate? (and why is this bad?)

$$P(F_j \mid \text{spam}) = \frac{\# \text{ of spam emails with word } F_j}{\text{total } \# \text{ of spam emails}}$$

$$P(\text{spam} \mid F_1, \dots, F_m) = \frac{\left[ \prod_{j=1}^m P(F_j \mid \text{spam}) \right] P(\text{spam})}{P(F_1, \dots, F_m)}$$

- Probability of zero destroys the entire calculation!

# Another problem to handle...

- Fix the estimates:

$$P(F_j \mid \text{spam}) = \frac{\# \text{ of spam emails with word } F_j + 1}{\text{total } \# \text{ of spam emails} + 2}$$

- This is called ***smoothing***. Removes the possibility of a zero probability wiping out the entire calculation.
- Simulates adding two additional spam emails, one containing every word, and containing no words.
  - We would also smooth for non-spam: adding two non-spam emails, one with all words, one with no words.



# Summary of Naïve Bayes

- Assumes the data is a collection of features, and each feature is conditionally independent of all other features given the hypothesis.
- Classifies using MAP hypothesis.

# Example

Suppose I have 20 emails that have been already classified into spam (15 emails) and non-spam (5 emails). Suppose I only care about the presence or absence of the words **luxury**, **brands**, and **save**.

Suppose 6 of the spam emails contain "luxury," 3 of the spam emails contain "brands," and 7 of the spam emails contain "save."

Suppose 1 of the non-spam emails contains "luxury," 2 of the non-spam emails contain "brands," and 2 of the non-spam emails contain "save."

Suppose a new email arrives that contains the words "luxury" and "save" but not "brands." Should this be classified as spam or not spam?

# Summary of Naïve Bayes

- Hypotheses:  $H_1$  through  $H_n$ .
- Features (data):  $F_1$  through  $F_m$ .

$$H^{\text{MAP}} = \operatorname{argmax}_i P(D | H_i) P(H_i)$$

$$= \operatorname{argmax}_i P(F_1, \dots, F_m | H_i) P(H_i)$$

$$= \operatorname{argmax}_i \left[ P(F_1 | H_i) \cdots P(F_m | H_i) \right] P(H_i)$$

$$= \operatorname{argmax}_i \left[ \prod_{j=1}^m P(F_j | H_i) \right] P(H_i)$$

# Summary of Naïve Bayes

- Probabilities needed to be determined (either given to you or estimated from data):
  - $P(H_i)$  for  $i = 1$  to  $n$ .
  - $P(F_j | H_i)$  for  $j = 1$  to  $m$  and  $i = 1$  to  $n$ .

# Summary of Naïve Bayes (for email)

- Naïve Bayes classifies using MAP:

$$\begin{aligned} H^{\text{MAP}} &= \operatorname{argmax}_i P(D | H_i)P(H_i) \\ &= \operatorname{argmax}_{i \in \{\text{spam}, \text{not-spam}\}} P(F_1, \dots, F_m | H_i)P(H_i) \\ &= \operatorname{argmax}_{i \in \{\text{spam}, \text{not-spam}\}} \left[ P(F_1 | H_i) \cdots P(F_m | H_i) \right] P(H_i) \\ &= \operatorname{argmax}_{i \in \{\text{spam}, \text{not-spam}\}} \left[ \prod_{j=1}^m P(F_j | H_i) \right] P(H_i) \end{aligned}$$

- Compute this for spam and for not-spam; see which is bigger.

# Summary of Naïve Bayes (for email)

- Estimating the **prior** for each hypothesis:

$$P(H_i) = \frac{\# \text{ of emails labeled as } H_i}{\text{total } \# \text{ of emails}}$$

- Estimating the probability of a feature given a class (aka **likelihood**) with smoothing:

$$P(F_j | H_i) = \frac{\# \text{ of } H_i \text{ emails with word } F_j + 1}{\text{total } \# \text{ of } H_i \text{ emails} + 2}$$