

# Naive Bayesian Classifier

AW

# Lecture Overview

1 Assumptions of Bayesian Classifier

2 Estimating Feature Probabilities

# Bayes Predictor

- Approach so far was distribution free learning: no assumptions on the underlying distribution over the data
  - As a consequence - discriminative approach in which our goal is not to learn the underlying distribution but rather to learn an accurate predictor
- Change of strategy: generative approach, in which it is assumed that the underlying distribution over the data has a specific parametric form and our goal is to estimate the parameters of the model
  - This task is called parametric density estimation.
- If we succeed in learning the underlying distribution  $\mathcal{D}$  over  $X \times \{0, 1\}$  accurately, then we can predict by using the Bayes optimal classifier:

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \Pr_{\mathcal{D}}(y = 1|x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } \frac{\Pr_{\mathcal{D}}(y=1|x)}{\Pr_{\mathcal{D}}(y=0|x)} \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

*for proof of optimality see materials section*

- Another way to describe Bayes predictor for a data point  $\bar{x} \in X$  is 
$$h_{\text{Bayes}} = \arg \max_{y \in \{0, 1\}} \Pr_{\mathcal{D}}(Y = y | X = \bar{x})$$

# Conditional Independence of Features

- Let  $X$  have  $n$  features with respective finite domains  $\mathbb{D}_1, \dots, \mathbb{D}_n$
  - $$\Pr_{\mathcal{D}}(Y = y | X = \bar{x}) = \Pr_{\mathcal{D}}(Y = y | x_1 = a_1, \dots, x_n = a_n)$$
$$= \frac{\Pr_{\mathcal{D}}(x_1 = a_1, \dots, x_n = a_n | Y = y) \Pr_{\mathcal{D}}(Y = y)}{\Pr_{\mathcal{D}}(x_1 = a_1, \dots, x_n = a_n)} .$$
  - Bayes predictor  $h_{Bayes} = \arg \max_{y \in \{0,1\}} \Pr_{\mathcal{D}}(Y = y | X = \bar{x})$  would need to be optimized with respect to  $d = |\mathbb{D}_1| \times \dots \times |\mathbb{D}_n|$  parameters
- $\Pr_{\mathcal{D}}(x_1 = a_1, \dots, x_n = a_n | Y = y)$
- For example if all features are Boolean then  $d = 2^n$ . So the number of parameters grows exponentially with the number of features - a bit too much!
  - The Naive Bayes approach makes the assumption about distribution (aka *generative assumption*) that given a class, all features are independent of each other given class  $y$ , i.e.

$$\Pr_{\mathcal{D}}(x_1 = a_1, \dots, x_n = a_n | Y = y) = \prod_{i=1}^n \Pr(x_i = a_i | Y = y)$$

- For 2 class predictor under these assumptions we have

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \frac{\Pr_{\mathcal{D}}(y=1|x)}{\Pr_{\mathcal{D}}(y=0|x)} = \frac{\Pr(Y=1) \prod_{i=1}^n \Pr(x_i=a_i|Y=1)}{\Pr(Y=0) \prod_{i=1}^n \Pr(x_i=a_i|Y=0)} \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

# Lecture Overview

1 Assumptions of Bayesian Classifier

2 Estimating Feature Probabilities

# Likelihood of a Parameter

We need to estimate  $\Pr(x_i = a_i | Y = y)$  and  $\Pr(Y = y)$ . From our training data for feature  $x_i$  we have a sample  $S_i^0$  of pairs  $(s_{i_1}, 0), \dots, (s_{i_k}, 0)$  of class 0 and a sample  $S_i^1$  of pairs  $(s_{i_{k+1}}, 1), \dots, (s_{i_m}, 1)$  of class 1. We make more assumptions:

- Our samples are coming from unknown pdf (or pmf)  $f_0(\cdot) = \mathcal{D}|_{X_i \times Y}$  (where  $\mathcal{D} \sim X \times Y$ ) that belongs to a certain family of distributions  $\{f(\cdot|\theta), \theta \in \Theta\}$ , so  $f_0(\cdot)$  is in fact  $f_0(\cdot|\theta_0)$  for some parameter  $\theta_0$  that we need to estimate
- $S_i^0 = \{s_{i_1}, \dots, s_{i_k}\}$  (resp.  $S_i^1$ ) is drawn independently from  $\mathcal{D}|_{X,0}$

Let  $f(s_{i_1}, s_{i_2}, \dots, s_{i_k} | \theta) = f(s_{i_1} | \theta) \times f(s_{i_2} | \theta) \times \dots \times f(s_{i_k} | \theta)$ . In this joint pdf (pmf)

- Observed values  $s_{i_1}, \dots, s_{i_k}$  are considered fixed "parameters"
- $\theta$  is the function's free variable

Define **log likelihood** (or when more convenient  $\ln$ -likelihood)

$$\begin{aligned}\mathcal{L}(\theta | s_{i_1}, \dots, s_{i_k}) &= \log(f(s_{i_1}, s_{i_2}, \dots, s_{i_k} | \theta)) \\ &= \log(\prod_{j=1}^k f(s_{i_j} | \theta)) = \sum_{j=1}^k \log(f(s_{i_j} | \theta))\end{aligned}$$

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(Y = y)$ .

- We assume Bernoulli distribution with unknown probability  $\theta$  of 1. Given sample  $S = \{y_1, \dots, y_m\}$
- Then  $f(y_i|\theta) = \begin{cases} 1 - \theta & \text{if } y_i = 0 \\ \theta & \text{otherwise} \end{cases}$
- $$\begin{aligned} \mathcal{L}(\theta|S) &= \log \left( \prod_{i=1}^m \theta^{y_i} \cdot (1 - \theta)^{1-y_i} \right) \\ &= \log(\theta) \sum_{i=1}^m y_i + \log(1 - \theta) \sum_{i=1}^m (1 - y_i) \end{aligned}$$
- $$\frac{d}{d\theta}(\mathcal{L}(\theta|S)) = \frac{\sum_{i=1}^m y_i}{\theta \ln 2} - \frac{\sum_{i=1}^m (1-y_i)}{(1-\theta) \ln 2} = 0 \text{ so } \theta = \frac{1}{m} \sum_{i=1}^m y_i$$



# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(X = x \mid Y = 0)$  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has finite domain  $[d] = \{0, 1, \dots, d\}$  and the sample  $S = S_0 \cup S_1$  where  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$  and  $S_1 = \{(z_{k+1}, 1) \dots, (z_m, 1)\}$ .

- Let our subsample  $S_0$  consists of  $x_0$  occurrences of  $X = 0$ ,  $x_1$  occurrences of  $X = 1$ ,  $\dots$ ,  $x_d$  occurrences of  $X = d$
- By law of total probability  $\Pr(Y = 0) = \sum_{i=0}^d \Pr(X = i \wedge Y = 0) = \sum_{i=0}^d \Pr(X = i \mid Y = 0) \Pr(Y = 0)$  so  $\sum_{i=0}^d \Pr(X = i \mid Y = 0) = 1$ . Let  $\theta_i = \Pr(X = i \mid Y = 0)$ . So  $\sum_i \theta_i = 1$ .
- Then  $f(S \mid \theta_1, \dots, \theta_d) = f(\{z_1, \dots, z_k\} \mid \theta_1, \dots, \theta_d) = \prod_{i=0}^d (\theta_i \cdot \Pr(Y = 0))^{x_i}$
- The likelihood of our sample  $S$  then is

$$\mathcal{L}(\theta|S) = \log \left( \prod_{i=0}^d (\theta_i \cdot \Pr(Y = 0))^{x_i} \right)$$

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

**Estimate  $\Pr(X = x \mid Y = 0)$**  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has finite domain  $[d] = \{0, 1, \dots, d\}$  and the sample  $S = S_0 \cup S_1$  where  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$  and  $S_1 = \{(z_{k+1}, 1) \dots, (z_m, 1)\}$ .

- Sample  $S_0$  consists of  $x_i$  occurrences of  $X = i$ . Denote  $\theta_i = \Pr(X = i \mid Y = 0)$ . So  $\sum_i \theta_i = 1$ .
- The likelihood of our sample  $S$  then is

$$\begin{aligned}\mathcal{L}(\theta|S) &= \log \left( \prod_{i=0}^d (\theta_i \cdot \Pr(Y = 0))^{x_i} \right) \\ &= k \log(\Pr(Y = 0)) + \sum_{i=0}^d x_i \log(\theta_i)\end{aligned}$$

- Thus maximize  $\mathcal{L}(\theta|S)$  is to maximize  $\sum_{i=0}^d x_i \log(\theta_i)$  since first term doesn't depend on parameters. Equivalently this means to minimize  $-\sum_{i=0}^d x_i \log(\theta_i)$ .

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(X = x \mid Y = 0)$  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has finite domain  $[d] = \{0, 1, \dots, d\}$  and the sample  $S = S_0 \cup S_1$  where  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$  and  $S_1 = \{(z_{k+1}, 1) \dots, (z_m, 1)\}$ . Let sample  $S_0$  consists of  $x_i$  occurrences of  $X = i$  and  $\theta_i = \Pr(X = i \mid Y = 0)$ . So  $\sum_i \theta_i = 1$ .

To maximize  $\mathcal{L}(\theta|S)$  we must solve the optimization problem:

$$\begin{aligned} &\text{minimize} && -\sum_{i=0}^d x_i \log(\theta_i) \\ &\text{subject to} && \sum_{i=0}^d \theta_i - 1 = 0 \end{aligned}$$

We use Lagrangian method: introduce a new variables  $\lambda_i$ , one per constraint, called a **Lagrange multipliers** and study a **Lagrangian** function

$\mathcal{L}(x_1, \dots, x_n, \lambda_1, \dots, \lambda_k) = f(x_1, \dots, x_n) \pm \sum_{i=1}^k \lambda_i \cdot g_i(x_1, \dots, x_n)$  where  $f(x_1, \dots, x_n)$  is the goal and  $g_i(x_1, \dots, x_n)$  are constraint functions

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(X = x \mid Y = 0)$  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has finite domain  $[d] = \{0, 1, \dots, d\}$  and the sample  $S = S_0 \cup S_1$  where  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$  and  $S_1 = \{(z_{k+1}, 1) \dots, (z_m, 1)\}$ . Let sample  $S_0$  consists of  $x_i$  occurrences of  $X = i$  and  $\theta_i = \Pr(X = i \mid Y = 0)$ . So  $\sum_i \theta_i = 1$ .

To maximize  $\mathcal{L}(\theta|S)$  we must solve the optimization problem:

$$\begin{aligned} &\text{minimize} && - \sum_{i=0}^d x_i \log(\theta_i) \\ &\text{subject to} && \sum_{i=0}^d \theta_i - 1 = 0 \end{aligned}$$

In our case  $\mathcal{L}(\theta_0, \dots, \theta_d, \lambda) = - \sum_{i=0}^d x_i \log(\theta_i) + \lambda \left( \sum_{i=0}^d \theta_i - 1 \right)$   
Stationary point of  $\mathcal{L}$  is a point where the partial derivatives of  $\mathcal{L}$  are zero. If  $f(x_1^*, \dots, x_n^*)$  is a minimum (resp. maximum) of the goal for the original constrained problem, then there exists  $\lambda_1^*, \dots, \lambda_k^*$  such that  $(x_1^*, \dots, x_n^*, \lambda_1^*, \dots, \lambda_k^*)$  is a stationary point for the Lagrangian.

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(X = x \mid Y = 0)$  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has finite domain  $[d] = \{0, 1, \dots, d\}$  and the sample  $S = S_0 \cup S_1$  where  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$  and  $S_1 = \{(z_{k+1}, 1) \dots, (z_m, 1)\}$ . Let sample  $S_0$  consists of  $x_i$  occurrences of  $X = i$  and  $\theta_i = \Pr(X = i \mid Y = 0)$ . So  $\sum_i \theta_i = 1$ .

To maximize  $\mathcal{L}(\theta|S)$  we must solve the optimization problem:

$$\begin{aligned} &\text{minimize} && -\sum_{i=0}^d x_i \log(\theta_i) \\ &\text{subject to} && \end{aligned}$$

$$\sum_{i=0}^d \theta_i - 1 = 0$$

So  $\max_{\lambda} \min_{\theta_1, \dots, \theta_d} \left[ -\sum_{i=0}^d x_i \log(\theta_i) + \lambda \left( \sum_{i=0}^d \theta_i - 1 \right) \right]$ . Partial derivatives for  $\theta_i$  is  $\frac{x_i}{\theta_i} - \lambda = 0$  for all  $i$ , or  $\theta_i = \frac{x_i}{\lambda}$ . Partial derivative for  $\lambda$  is  $\sum_{i=0}^d \theta_i - 1 = 0$  which yields  $\lambda = \sum_{i=0}^d x_i = \begin{cases} k & \text{if } Y = 0 \\ m - k & \text{if } Y = 1 \end{cases}$ , so e.g. for  $Y = 0$  we get  $\theta_i = \frac{x_i}{k}$ .

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

Estimate  $\Pr(X = x \mid Y = 0)$  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has continuous domain, assume that  $\Pr(X = x \mid Y = 0)$  is normally distributed and the sample is  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$ . In other words we know pdf  $p(X = x \mid Y = 0) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  where parameter  $\bar{\theta}$  is  $(\mu, \sigma)^T$ . Then  $f(\{z_1, \dots, z_k\} \mid \bar{\theta}) = \prod_{i=1}^k \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(z_i-\mu)^2}{2\sigma^2}\right)$  and

$$\begin{aligned}\mathcal{L}(\bar{\theta}|S) &= \log\left(\prod_{i=1}^k \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(z_i-\mu)^2}{2\sigma^2}\right)\right) \\ &= -\frac{\sum_{i=1}^k (z_i-\mu)^2}{2\sigma^2} - k \log(\sigma\sqrt{2\pi})\end{aligned}$$

So obtaining  $\hat{\theta}$  involves solving  $\frac{\partial}{\partial \mu} (\mathcal{L}(\bar{\theta}|S)) = 0$  and  $\frac{\partial}{\partial \sigma} (\mathcal{L}(\bar{\theta}|S)) = 0$ ,

# Maximum Likelihood Estimation

Given a sample  $S$  maximum likelihood estimator of parameter  $\theta$  is

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S)$$

For our purposes:

**Estimate  $\Pr(X = x \mid Y = 0)$**  (or resp.  $\Pr(X = x \mid Y = 1)$ ) when feature  $X$  has continuous domain, assume that  $\Pr(X = x \mid Y = 0)$  is normally distributed and the sample is  $S_0 = \{(z_1, 0) \dots, (z_k, 0)\}$ .

$$\begin{aligned}\mathcal{L}(\bar{\theta}|S) &= \log \left( \prod_{i=1}^k \frac{1}{\sigma\sqrt{2\pi}} \exp \left( -\frac{(z_i - \mu)^2}{2\sigma^2} \right) \right) \\ &= -\frac{\sum_{i=1}^k (z_i - \mu)^2}{2\sigma^2} - k \log(\sigma\sqrt{2\pi})\end{aligned}$$

So obtaining  $\hat{\theta}$  involves solving  $\frac{\partial}{\partial \mu} (\mathcal{L}(\bar{\theta}|S)) = 0$  and  $\frac{\partial}{\partial \sigma} (\mathcal{L}(\bar{\theta}|S)) = 0$ ,

$$\begin{aligned}\frac{\partial}{\partial \mu} (\mathcal{L}(\bar{\theta}|S)) &= \frac{\sum_{i=1}^k (z_i - \mu)}{\sigma^2} = 0 & \frac{\partial}{\partial \sigma} (\mathcal{L}(\bar{\theta}|S)) &= \frac{\sum_{i=1}^k (z_i - \mu)^2}{\sigma^3} - \frac{k}{\sigma} = 0 \\ \Downarrow & & \Downarrow & \\ \hat{\mu} &= \frac{\sum_{i=1}^k z_i}{k} & \hat{\sigma} &= \sqrt{\frac{\sum_{i=1}^k (z_i - \mu)^2}{k}}\end{aligned}$$