

# Probabilistic Generative Models

# Classification so far

First solve the inference problem of determining the posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$ , and then subsequently use decision theory to assign each new  $\mathbf{x}$  to one of the classes. Approaches that model the posterior probabilities directly are called *discriminative models*.

Find a function  $f(\mathbf{x})$ , called a discriminant function, which maps each input  $\mathbf{x}$  directly onto a class label. For instance, in the case of two-class problems,  $f(\cdot)$  might be binary valued and such that  $f = 0$  represents class  $\mathcal{C}_1$  and  $f = 1$  represents class  $\mathcal{C}_2$ . In this case, probabilities play no role.

# Generative classifiers

First solve the inference problem of determining the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$  for each class  $\mathcal{C}_k$  individually. Also separately infer the prior class probabilities  $p(\mathcal{C}_k)$ . Then use Bayes' theorem in the form

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$$

Equivalently, we can model the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  directly and then normalize to obtain the posterior probabilities

Approaches that explicitly or implicitly model the distribution of inputs as well as outputs are known as ***generative models***

# Probabilistic Generative Models

Approaches that explicitly or implicitly model the distribution of inputs as well as outputs are known as **generative models**, because by sampling from them it is possible to generate synthetic data points in the input space.

Approaches that model the posterior probabilities directly are called **discriminative models**

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} & a &= \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} & a_k &= \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k). \end{aligned}$$

# Probabilistic Generative Models

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \qquad a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \end{aligned}$$

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

# Probabilistic Generative Models

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \end{aligned}$$

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

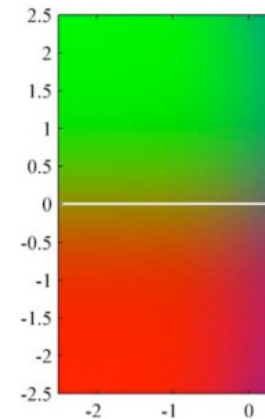
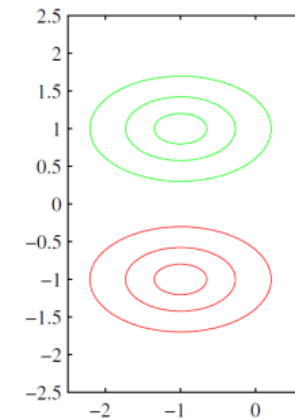
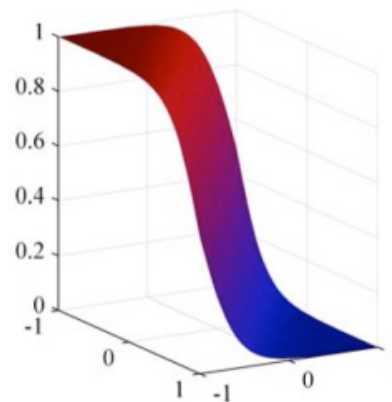
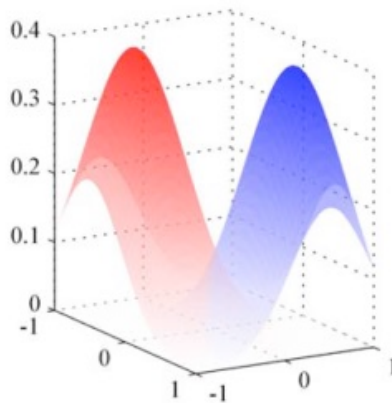
$$\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}.$$

# Probabilistic Generative Models

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}.$$

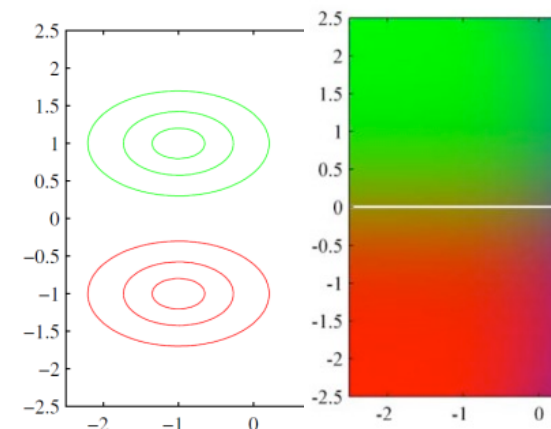
$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad \begin{aligned} \mathbf{w} &= \Sigma^{-1}(\mu_1 - \mu_2) \\ w_0 &= -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \end{aligned}$$



# Probabilistic Generative Models

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}.$$

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \sigma(\mathbf{w}^T \mathbf{x} + w_0) \\ \mathbf{w} &= \Sigma^{-1}(\mu_1 - \mu_2) \\ w_0 &= -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \end{aligned}$$



$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

$$a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k). \quad a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\begin{aligned} \mathbf{w}_k &= \Sigma^{-1} \mu_k \\ w_{k0} &= -\frac{1}{2}\mu_k^T \Sigma^{-1} \mu_k + \ln p(\mathcal{C}_k). \end{aligned}$$



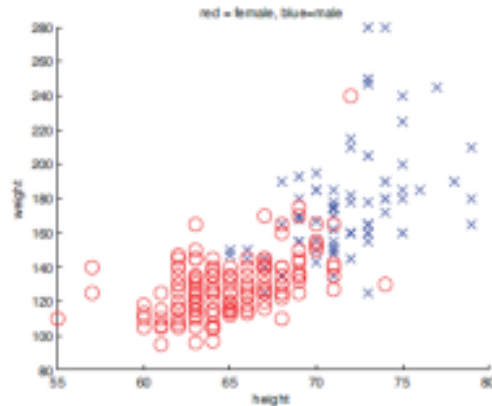
# Probabilistic Generative Models

model the **class-conditional densities**  $p(\mathbf{x}|\mathcal{C}_k)$ , as well as the **class priors**  $p(\mathcal{C}_k)$

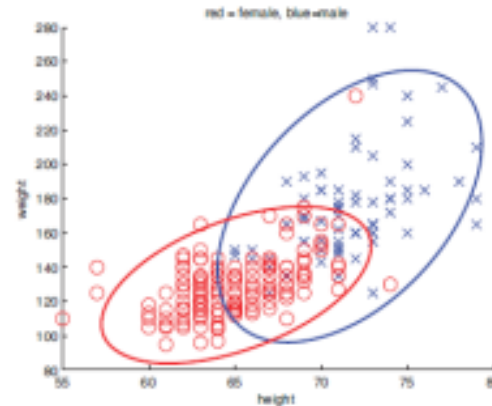
$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

**class-conditional densities**  $p(\mathbf{x}|\mathcal{C})$

$$p(\mathbf{x}|y = c, \theta) = \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)$$



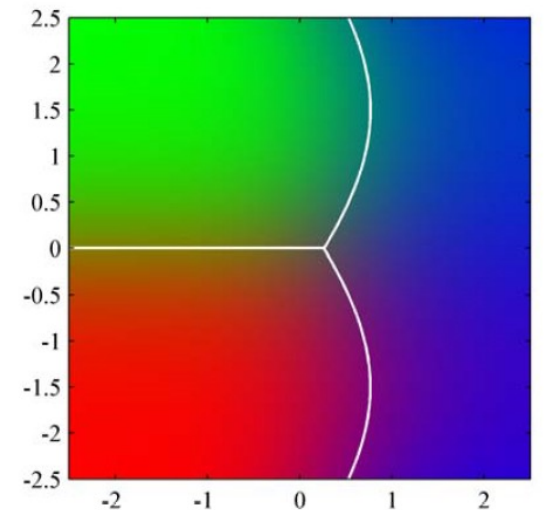
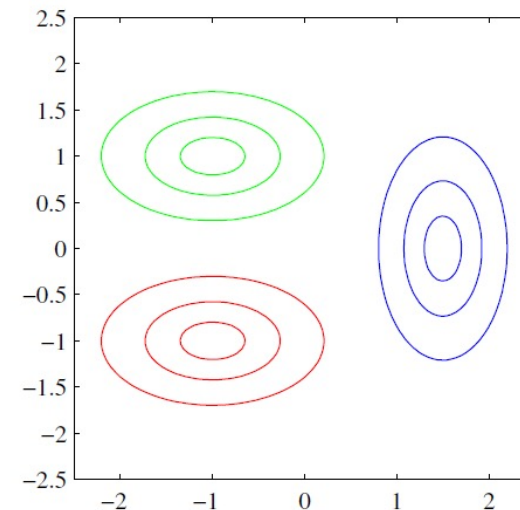
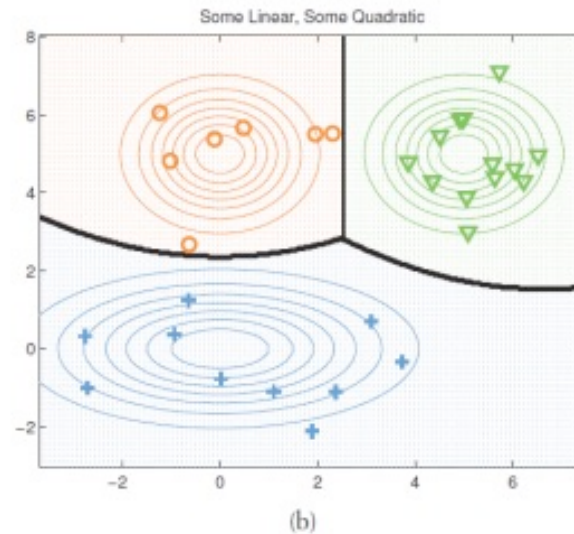
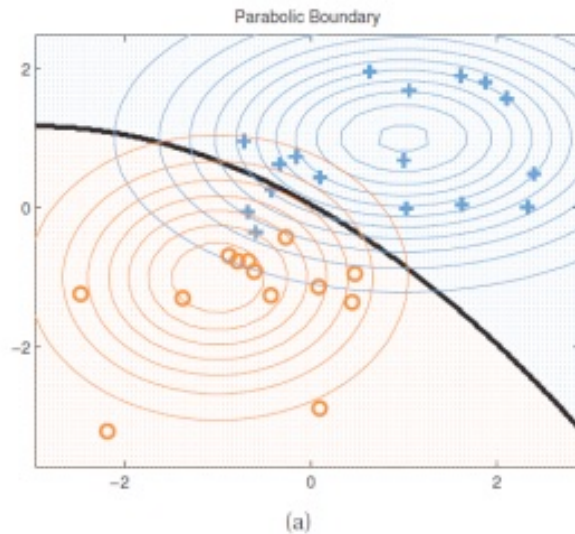
(a)



(b)

# Quadratic discriminant analysis

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \frac{\pi_c |2\pi \boldsymbol{\Sigma}_c|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right]}{\sum_{c'} \pi_{c'} |2\pi \boldsymbol{\Sigma}_{c'}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_{c'})^T \boldsymbol{\Sigma}_{c'}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{c'}) \right]}$$



# Linear Discriminant Analysis

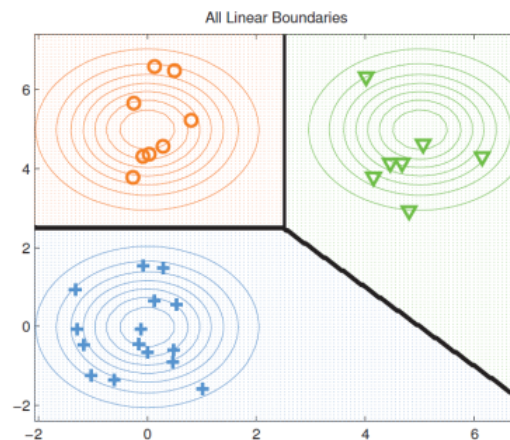
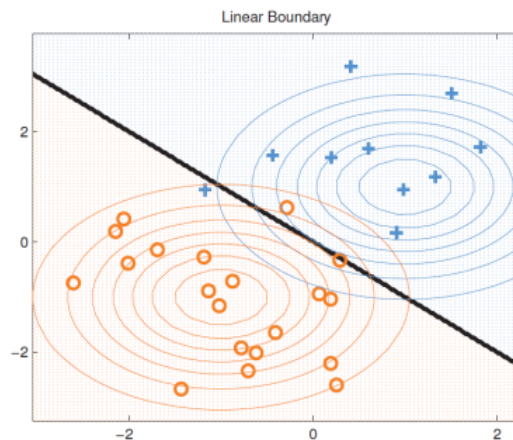
- covariance matrices are **tied** or **shared** across classes,  $\Sigma_c = \Sigma$ .

$$\begin{aligned} p(y = c | \mathbf{x}, \boldsymbol{\theta}) &\propto \pi_c \exp \left[ \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \right] \\ &= \exp \left[ \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right] \exp \left[ -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right] \end{aligned}$$

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}}}$$

$$\gamma_c = -\frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c$$

$$\boldsymbol{\beta}_c = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c$$



# Generative Models: Parameter estimation

- Models joint probability of observing input and output

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n|\mathcal{C}_1) = \pi\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}).$$

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n|\mathcal{C}_2) = (1 - \pi)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

$$p(\mathbf{t}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

# Generative Models : Parameter estimation

- Models joint probability of observing input and output

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \left[ \sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i = c) \log \pi_c \right] + \sum_{c=1}^C \left[ \sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right]$$

Use **Maximum Likelihood Estimation (MLE)** to estimate parameters

# Generative Models : Parameter estimation

- Models joint probability of observing input and output

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \left[ \sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i = c) \log \pi_c \right] + \sum_{c=1}^C \left[ \sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right]$$

$$\hat{\pi}_c = \frac{N_c}{N}, \quad \hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{x}_i, \quad \hat{\boldsymbol{\Sigma}}_c = \frac{1}{N_c} \sum_{i:y_i=c} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)^T$$

- MLE can badly overfit in high dimensions.
- Use a diagonal covariance matrix for each class, which assumes the features are conditionally independent; this is equivalent to using a **naive Bayes classifier**
- Other approaches : Use **MAP** and **Bayesian approaches**

# Naïve Bayes Classifier

- Features are discrete

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

- class-conditional density  $p(\mathbf{x}|\mathbf{y})$ 
  - Number of parameters ?

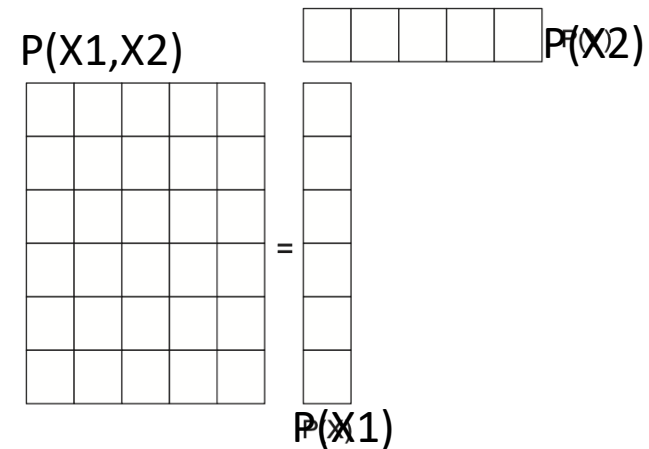
# Naïve Bayes Classifier

- Features are discrete

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

- “naive” since we do not expect the features to be independent, but results in classifiers that work well
  - model is quite simple and hence it is relatively immune to overfitting, as a lower number of parameters need to be estimated due to independence assumption.
- Class-conditional density  $p(\mathbf{x}|\mathbf{y})$

$$p(\mathbf{x}|\mathbf{y} = c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|\mathbf{y} = c, \boldsymbol{\theta}_{jc})$$





# Naïve Bayes Classifier

- Features are conditionally independent given the class label.

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y = c, \boldsymbol{\theta}_{jc})$$

- class-conditional density  $p(\mathbf{x}|y)$ 
  - In the case of real-valued features, we can use the Gaussian distribution:  $p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_{jc}, \sigma_{jc}^2)$ , where  $\mu_{jc}$  is the mean of feature  $j$  in objects of class  $c$ , and  $\sigma_{jc}^2$  is its variance.
  - In the case of binary features,  $x_j \in \{0, 1\}$ , we can use the Bernoulli distribution:  $p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_j|\mu_{jc})$ , where  $\mu_{jc}$  is the probability that feature  $j$  occurs in class  $c$ .
  - In the case of categorical features,  $x_j \in \{1, \dots, K\}$ , we can model use the multinoulli distribution:  $p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Cat}(x_j|\boldsymbol{\mu}_{jc})$ , where  $\boldsymbol{\mu}_{jc}$  is a histogram over the  $K$  possible values for  $x_j$  in class  $c$ .

# Training NBC : ML Estimation

$$p(\mathbf{x}_i, y_i | \boldsymbol{\theta}) = p(y_i | \boldsymbol{\pi}) \prod_j p(x_{ij} | \boldsymbol{\theta}_j) = \prod_c \pi_c^{\mathbb{I}(y_i=c)} \prod_j \prod_c p(x_{ij} | \boldsymbol{\theta}_{jc})^{\mathbb{I}(y_i=c)}$$

$$\log p(\mathcal{D} | \boldsymbol{\theta}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{j=1}^D \sum_{c=1}^C \sum_{i: y_i=c} \log p(x_{ij} | \boldsymbol{\theta}_{jc})$$

Bernoulli class conditional likelihood

$$\hat{\pi}_c = \frac{N_c}{N} \quad \hat{\theta}_{jc} = \frac{N_{jc}}{N_c}$$

The method is easily generalized to handle features of mixed type.

# NBC : Algorithm

---

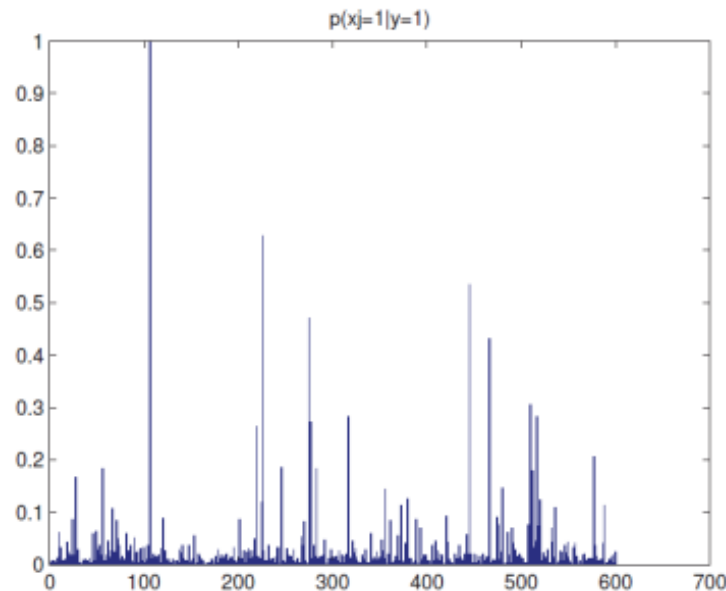
**Algorithm 3.1:** Fitting a naive Bayes classifier to binary features

---

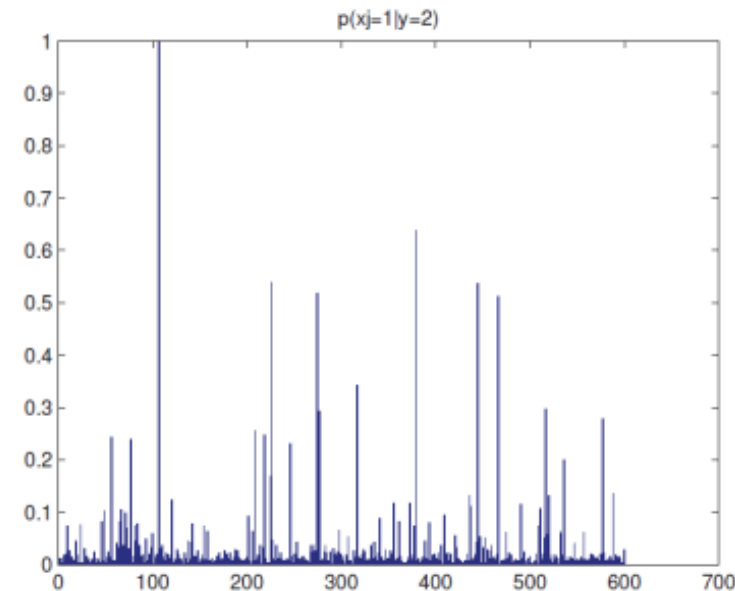
```
1  $N_c = 0, N_{jc} = 0;$ 
2 for  $i = 1 : N$  do
3    $c = y_i$  // Class label of  $i$ 'th example;
4    $N_c := N_c + 1;$ 
5   for  $j = 1 : D$  do
6     if  $x_{ij} = 1$  then
7        $N_{jc} := N_{jc} + 1$ 
8  $\hat{\pi}_c = \frac{N_c}{N}, \hat{\theta}_{jc} = \frac{N_{jc}}{N}$ 
```

---

$O(ND)$  time.



(a)



(b)

# Spam Classification example

## **HAM** examples

- d1: {good}
- d2: {very good}

$$P(y=\text{ham}) = 2/5$$

## **SPAM** examples

- d3: {bad}
- d4: {very bad}
- d5: {very bad very bad}

$$p(y=\text{spam}) = 3/5$$

- Test data: d6: {good bad very bad} - **SPAM OR HAM??**
- Vocabulary :  $V = \{\text{good, bad, very}\}$
- Use Naive Bayes classifier:
- $P(\text{ham} \mid d6) = \prod P(\text{word} \mid y=\text{ham}) P(y=\text{ham})$   
 $= P(\text{good} \mid \text{ham}) * P(\text{bad} \mid \text{ham}) * P(\text{very} \mid \text{ham})$

Word frequencies wrt class	$P(\text{good} \mid y)$	$P(\text{bad} \mid y)$	$P(\text{very} \mid y)$
Class 0: $y = \text{spam}$	0/3	3/3	2/3
Class 1: $y = \text{ham}$	2/2	0/2	1/2

- Estimate parameters using ML:
- $P(\text{ham} \mid d6) = P(\text{good} \mid \text{ham}) * P(\text{bad} \mid \text{ham}) * P(\text{very} \mid \text{ham}) * P(\text{ham})$
- $P(\text{spam} \mid d6) = P(\text{good} \mid \text{spam}) * P(\text{bad} \mid \text{spam}) * P(\text{very} \mid \text{spam}) * P(\text{spam})$
- What is the problem with this ?

Word frequencies wrt class	$P(\text{good} \mid y)$	$P(\text{bad} \mid y)$	$P(\text{very} \mid y)$
Class 0: $y = \text{spam}$	0/3	3/3	2/3
Class 1: $y = \text{ham}$	2/2	0/2	1/2

- Estimate parameters using ML:
- $P(\text{ham} \mid d6) = P(\text{good} \mid \text{ham}) * P(\text{bad} \mid \text{ham}) * P(\text{very} \mid \text{ham}) * P(\text{ham})$
- $P(\text{spam} \mid d6) = P(\text{good} \mid \text{spam}) * P(\text{bad} \mid \text{spam}) * P(\text{very} \mid \text{spam}) * P(\text{spam})$

- What is the problem with this ?



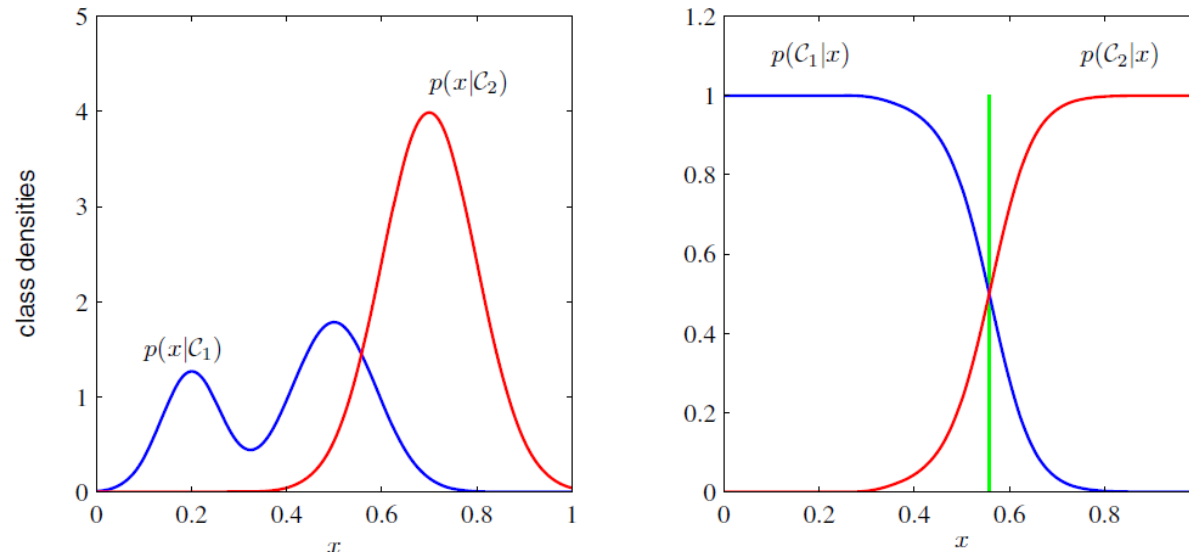
Soln : Bayesian Naïve Bayes Classifiers!

# Discriminative vs Generative

- **Generative Models** : model the joint distribution  $p(\mathbf{x}, C_k)$  directly and then normalize to obtain the posterior probabilities.

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}$$

- Most demanding because it involves finding the joint distribution over both  $\mathbf{x}$  and  $C_k$ . For many applications,  $\mathbf{x}$  will have high dimensionality, and consequently we may need a large training set. if we only wish to make classification decisions, then it can be wasteful of computational resources but can be useful for detecting **outliers or novel classes**.



# Discriminative vs Generative

- **Generative Models** : model the joint distribution  $p(\mathbf{x}, C_k)$  directly and then normalize to obtain the posterior probabilities.

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}$$

- **Easy to adapt to new data** : Use posterior probabilities as prior
- **Multi-Modal Data** : combining multiple modalities and domains

$$p(\mathbf{x}_I, \mathbf{x}_B|C_k) = p(\mathbf{x}_I|C_k)p(\mathbf{x}_B|C_k).$$

- **Easy to fit?** very easy to fit generative classifiers. we can fit a naive Bayes model and an LDA model by simple counting and averaging. logistic regression requires solving a convex optimization problem which is much slower.
- **Fit classes separately?** In a generative classifier, we estimate the parameters of each class conditional density independently, so we do not have to retrain the model when we add more classes.
- **Well-calibrated probabilities?** Some generative models, such as naive Bayes, make strong independence assumptions which are often not valid
- Generative models can easily handle unlabelled data and missing features