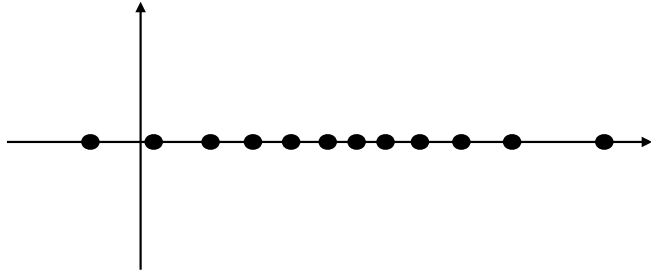


# **Bayes Classifier and Naive Bayes**

# Announcements

HW 2 is out — start early

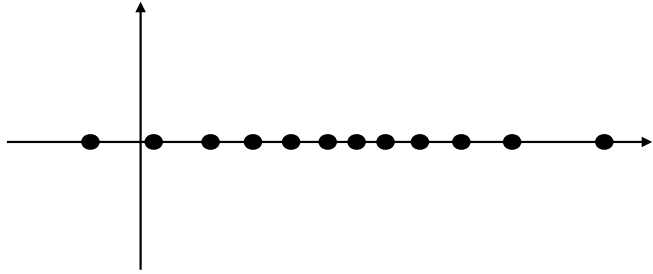
# Recap on MLE



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

Assume data is from  $\mathcal{N}(\mu^*, \sigma^2)$ , want to estimate  $\mu^*, \sigma$  from the data  $\mathcal{D}$  MLE

# Recap on MLE



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

Assume data is from  $\mathcal{N}(\mu^*, \sigma^2)$ , want to estimate  $\mu^*, \sigma$  from the data  $\mathcal{D}$  MLE

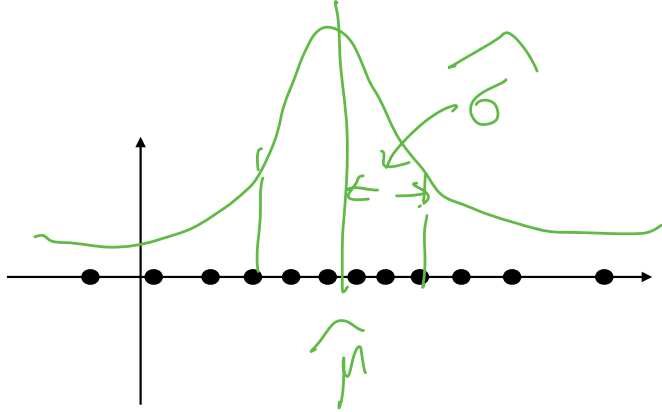
$$P(\mathcal{D} | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2/\sigma^2\right)$$

MLE:

$$\underset{\mu, \sigma}{\operatorname{argmax}} P(\mathcal{D} | \mu, \sigma)$$

$P(x | \mu, \sigma)$

# Recap on MLE



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

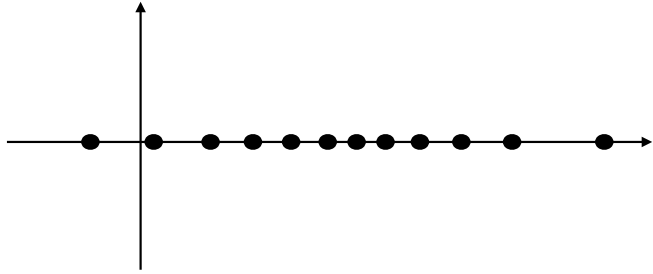
Assume data is from  $\mathcal{N}(\mu^*, \sigma^2)$ , want to estimate  $\mu^*, \sigma$  from the data  $\mathcal{D}$  MLE

$$P(\mathcal{D} | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2/\sigma^2\right)$$

The solution that maximizes the log-likelihood:

$$\hat{\mu} = \sum_{i=1}^n x_i/n, \quad \hat{\sigma}^2 = \sum_{i=1}^n (x_i - \hat{\mu})^2/n$$

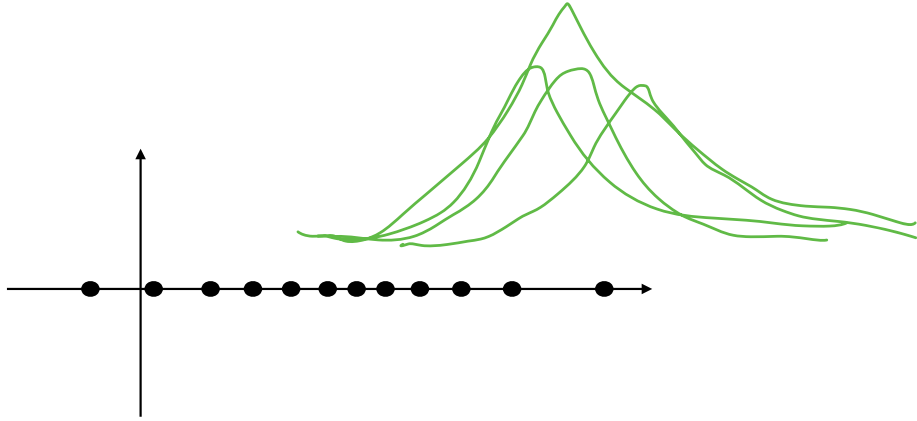
# Recap on MAP



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

Now if we want to use **MAP**:

# Recap on MAP

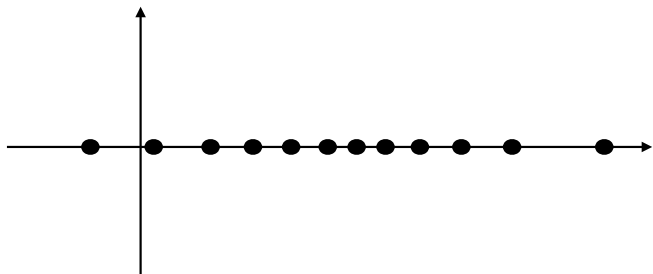


$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

Now if we want to use **MAP**:

1. Pick a prior:  $P(\mu, \sigma)$

# Recap on MAP



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

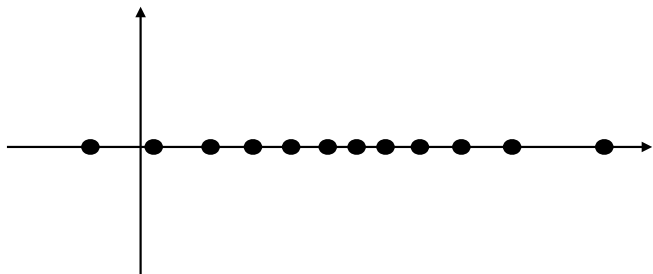
Now if we want to use **MAP**:

1. Pick a prior:  $P(\mu, \sigma)$

2. Write down data-likelihood:  $P(\mathcal{D} | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2/\sigma^2\right)$



# Recap on MAP



$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$

Now if we want to use **MAP**:

1. Pick a prior:  $P(\mu, \sigma)$

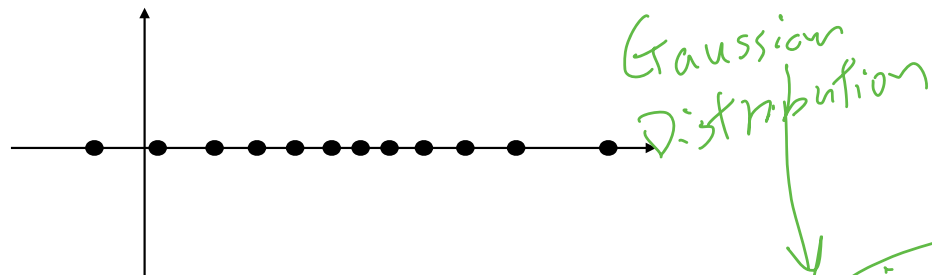
2. Write down data-likelihood:  $P(\mathcal{D} | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2/\sigma^2\right)$

3. Form posterior  $P(\mu, \sigma | \mathcal{D}) \propto P(\mu, \sigma)P(\mathcal{D} | \mu, \sigma)$

$$\underset{\mu, \sigma}{\operatorname{argmax}} P(\mu, \sigma | \mathcal{D})$$

# Recap on MAP

$$\mathcal{D} = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}$$



Now if we want to use **MAP**:

$$P(\mu)P(\sigma) \leftarrow \text{Inverse Gamma}$$

1. Pick a prior:  $P(\mu, \sigma)$

2. Write down data-likelihood:  $P(\mathcal{D} | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2/\sigma^2\right)$

3. Form posterior  $P(\mu, \sigma | \mathcal{D}) \propto P(\mu, \sigma)P(\mathcal{D} | \mu, \sigma)$

$$P(\sigma | \mathcal{D}) \leftarrow \text{Inverse Gamma}$$

# Today

Objective: learn our second classification algorithm—Naive Bayes (derived via MLE)

# Outline

1. General formulation of Naive Bayes

2. Example

3. Connection to linear classifier

# Generative modeling

Setting: binary classification w/ dataset  $\{x_i, y_i\}_{i=1}^n$ ,  $(x_i, y_i) \sim P$ , where  $x \in \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

Goal: estimate  $P(y|x)$

$$x_i \in \mathbb{R}^d$$

$$\hat{P}(y|x) \text{ from } \mathcal{D}$$

$$\text{Ex. } \hat{P}(y=1 | x=x_0)$$

$$\arg \max_{y \in \{-1, 1\}} P(y|x)$$

$$P(y|x) = \frac{P(x, y)}{P(x)} = \frac{\sum_{i=1}^n \mathbb{1}(x_i = x, y_i = y)}{\sum_{i=1}^n \mathbb{1}(x_i = x)}$$

# Generative modeling

Setting: binary classification w/ dataset  $\{x_i, y_i\}_{i=1}^n$ ,  $(x_i, y_i) \sim P$ , where  $x \in \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

Goal: estimate  $P(y | x)$

We take a generative modeling approach here:

# Generative modeling

Setting: binary classification w/ dataset  $\{x_i, y_i\}_{i=1}^n$ ,  $(x_i, y_i) \sim P$ , where  $x \in \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

Goal: estimate  $P(y | x)$

We take a generative modeling approach here:

Bayes Rule

$$P(y | x) \propto P(x | y)P(y)$$

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

# Generative modeling

Setting: binary classification w/ dataset  $\{x_i, y_i\}_{i=1}^n$ ,  $(x_i, y_i) \sim P$ , where  $x \in \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

Goal: estimate  $P(y | x)$

We take a **generative modeling** approach here:

$$P(y | x) \propto P(x | y)P(y)$$

Estimate  $P(x | y)$  &  $P(y)$  from data  
(hence generative modeling)



# Generative modeling

Setting: binary classification w/ dataset  $\{x_i, y_i\}_{i=1}^n$ ,  $(x_i, y_i) \sim P$ , where  $x \in \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

Goal: estimate  $P(y | x)$

We take a **generative modeling** approach here:

$$P(y | x) \propto P(x | y)P(y)$$

Estimate  $P(x | y)$  &  $P(y)$  from data  
(hence generative modeling)

( Discriminative modeling: directly estimate  $P(y | x)$  )

# Naive Bayes

Estimate  $P(y)$  from data:

$$P(y|x) \propto P(x|y)P(y)$$

$$y = \{+1, -1\}$$

# Naive Bayes

Estimate  $P(y)$  from data:

$$P(y|x) \propto P(x|y)P(y)$$

Estimate  $P(y)$  is easy:

$$P(y = 1) \approx \frac{\sum_{i=1}^n \mathbf{1}(y_i = 1)}{n}$$

$$P(y = -1) = \frac{\sum_{i=1}^n \mathbf{1}(y_i = -1)}{n}$$

# Naive Bayes

Estimate  $P(x | y)$  from data:

$$P(y | x) \propto P(x | y)P(y)$$

Estimate  $P(x | y)$  is not easy:

# Naive Bayes

Estimate  $P(x | y)$  from data:

$x \in \mathbb{R}^d$

Estimate  $P(x | y)$  is not easy:


$$P(y | x) \propto P(x | y)P(y)$$

$x$  can be high-dimensional, e.g.,  $d$  is large!

There may not be repetitions in  $\{x_i\}_{i=1}^n$ !

# The key assumption in Naive Bayes

The Naive Bayes assumption:

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$


# The key assumption in Naive Bayes

The Naive Bayes assumption:

$$P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$$

Conditioned on label  $y$ , feature values  
are **independent!**

# About the independence assumption

The Naive Bayes assumption:

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

Conditioned on label  $y$ , feature values are **independent!**

Q: does conditional independence imply global independence?

$$P(x) = \prod_{i=1}^d P(x[i])$$



# About the independence assumption

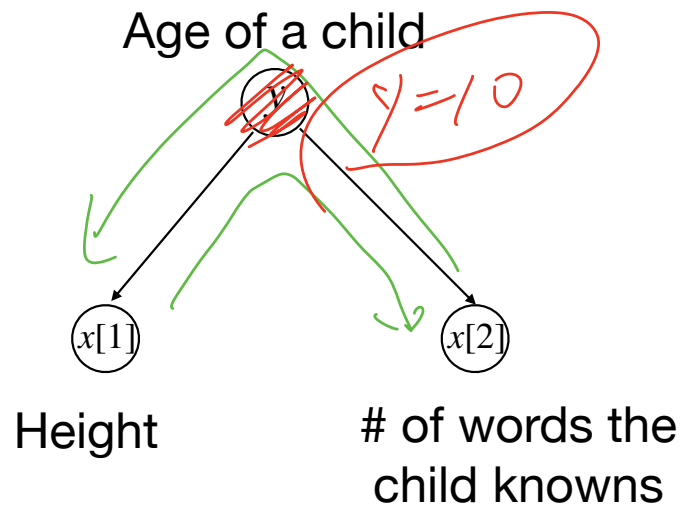
The Naive Bayes assumption:

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

Conditioned on label  $y$ , feature values are **independent!**

$$P(x^{(1)}, x^{(2)}) \neq P(x^{(1)}) P(x^{(2)})$$

Q: does conditional independence imply global independence?



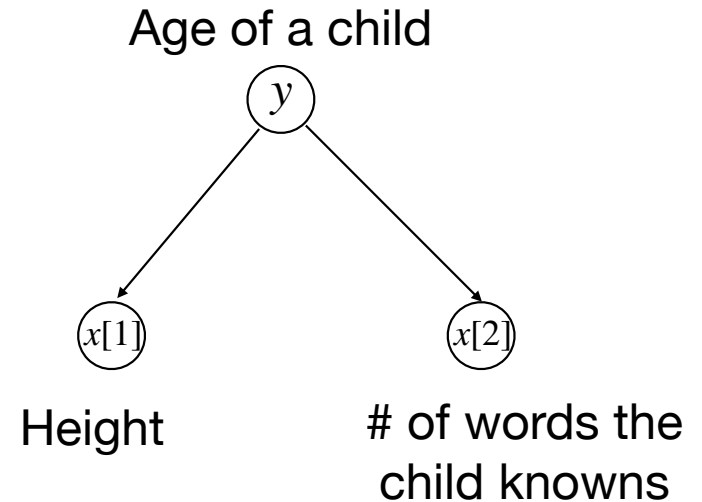
# About the independence assumption

The Naive Bayes assumption:

$$P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$$

Conditioned on label  $y$ , feature values are **independent!**

Q: why it is also a naive assumption?



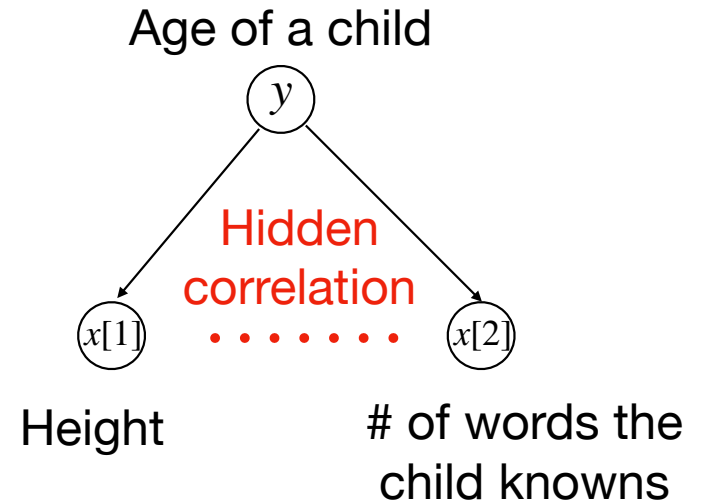
# About the independence assumption

The Naive Bayes assumption:

$$P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$$

Conditioned on label  $y$ , feature values are **independent!**

Q: why it is also a naive assumption?



# Naive Bayes

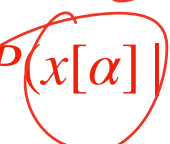
Estimate  $P(x | y)$  from data:

# Naive Bayes

Estimate  $P(x | y)$  from data:

W/ the NB assumption  $P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$

1-D scalar



# Naive Bayes

Estimate  $P(x | y)$  from data:

W/ the NB assumption  $P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$

Now we can estimate  $P(x[\alpha] | y)$  for each  $\alpha$

# Naive Bayes

Estimate  $P(x | y)$  from data:

W/ the NB assumption  $P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$

Now we can estimate  $P(x[\alpha] | y)$  for each  $\alpha$



1-dim problem!

# Naive Bayes

Once estimated  $P(y)$  and  $P(x | y)$ , we can make prediction:

In test time, given  $x$ :

$$P(y | x) \propto P(x | y)P(y)$$



# Naive Bayes

Once estimated  $P(y)$  and  $P(x | y)$ , we can make prediction:

$$P(y | x) \propto P(x | y)P(y)$$

In test time, given  $x$ :

$$\hat{y} = \arg \max_y P(y | x)$$

# Naive Bayes

Once estimated  $P(y)$  and  $P(x | y)$ , we can make prediction:

$$P(y | x) \propto P(x | y)P(y)$$

In test time, given  $x$ :

$$\begin{aligned}\hat{y} &= \arg \max_y P(y | x) \\ &= \arg \max_y P(x | y)P(y)\end{aligned}$$

# Naive Bayes

Once estimated  $P(y)$  and  $P(x | y)$ , we can make prediction:

$$P(y | x) \propto P(x | y)P(y)$$

In test time, given  $x$ :

$$\hat{y} = \arg \max_y P(y | x)$$

$$= \arg \max_y P(x | y)P(y) = \arg \max_y \left( \prod_{\alpha=1}^d P(x[\alpha] | y) \right) P(y)$$

$$= \arg \max_y \left( \sum_{\alpha=1}^d \ln P(x[\alpha] | y) + \ln P(y) \right)$$

*NB Assumption*

# Outline

1. General formulation of Naive Bayes



2. Example

3. Connection to linear classifier

# Case study

Continuous features

$x[\alpha] \in \mathbb{R}$ , for all  $\alpha \in \{1, 2, \dots, d\}$

$$x \in \mathbb{R}^d$$

# Case study

## Continuous features

$$x[\alpha] \in \mathbb{R}, \text{ for all } \alpha \in \{1, 2, \dots, d\}$$

We model each  $P(x[\alpha] | y)$  using a 1-dim Gaussian distribution:

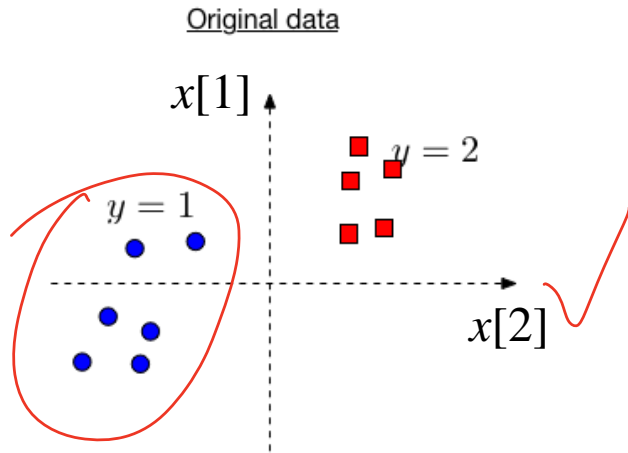
$$P(x[\alpha] | y) = \mathcal{N}(\underbrace{\mu_{\alpha,y}}_{\Delta}, \underbrace{\sigma_{\alpha,y}^2}_{\lambda})$$

# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$ :

$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha,y}^2)$$

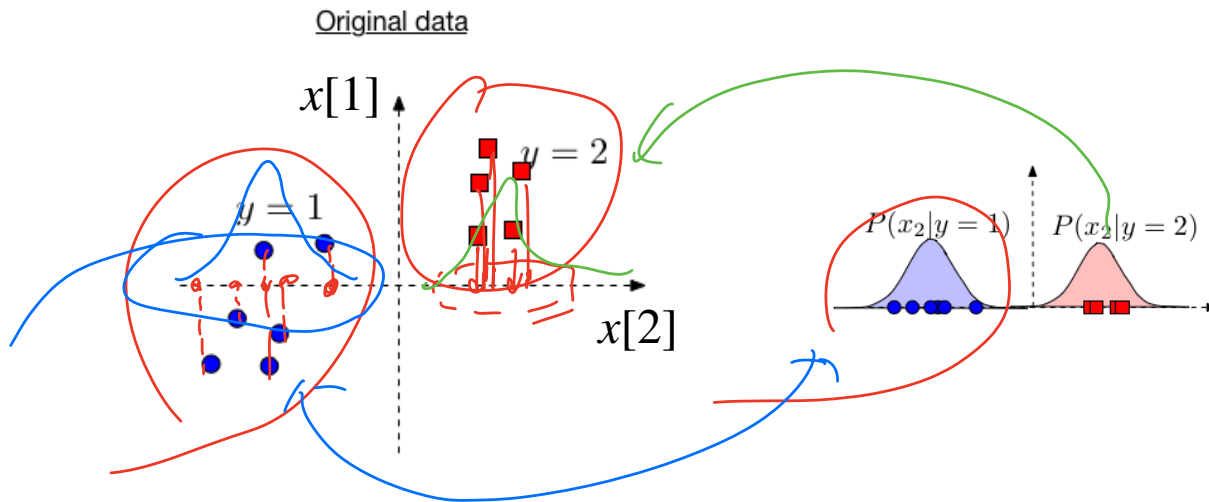
$P(x|y)$



# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$ :

$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha,y}^2)$$

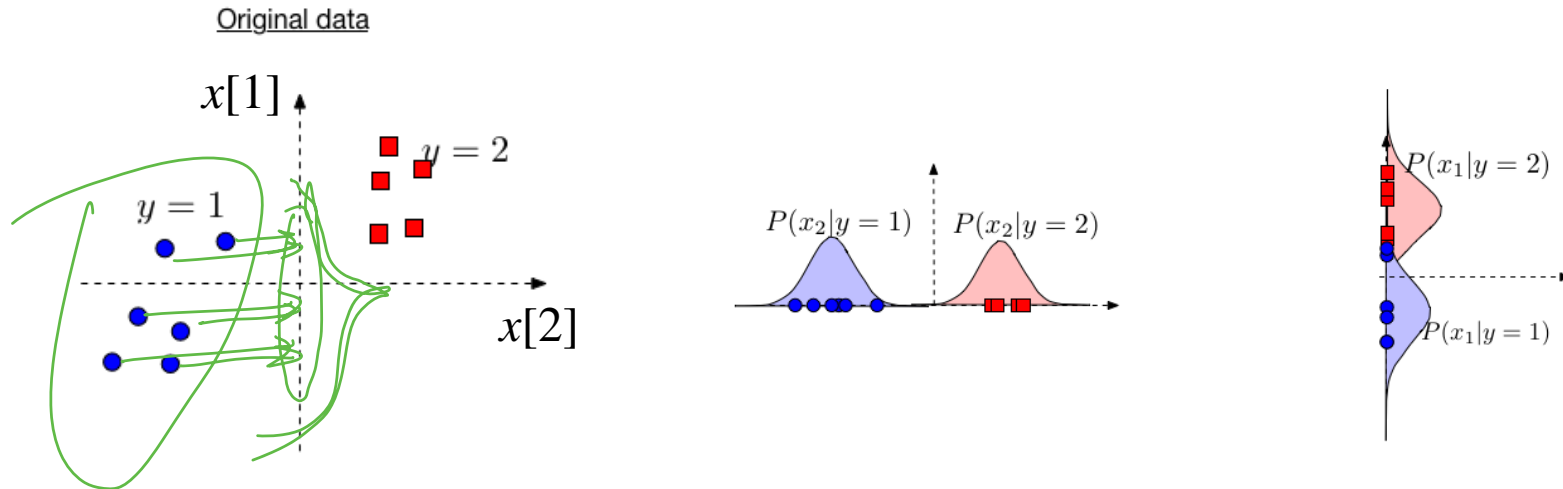




# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$ :

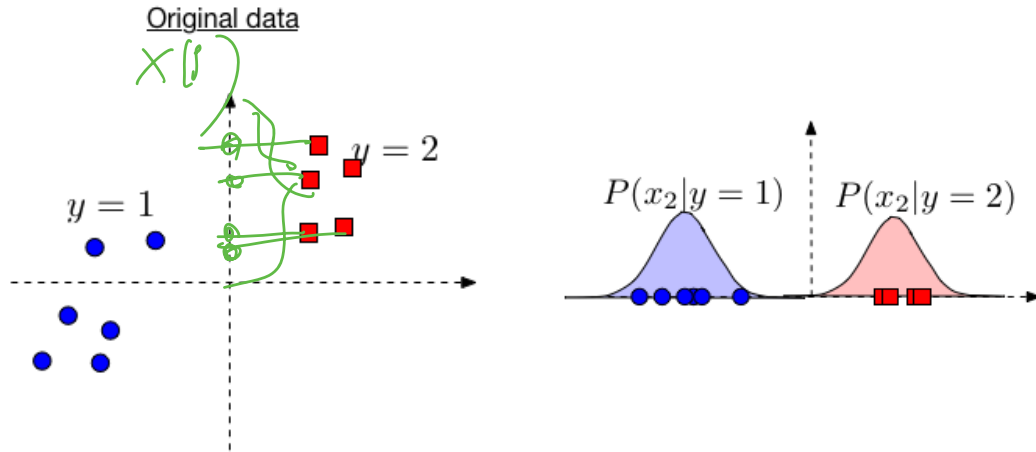
$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha,y}^2)$$



# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$  via **MLE**:

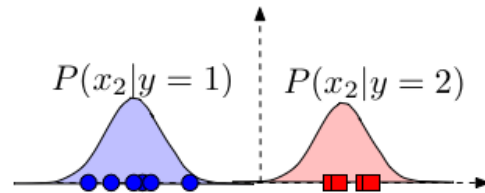
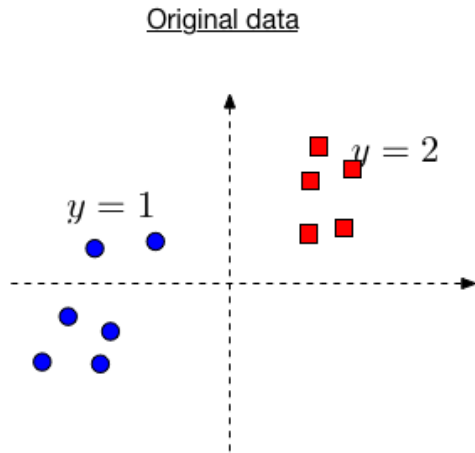
$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\sigma,\alpha}^2)$$



# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$  via **MLE**:

$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\sigma,\alpha}^2)$$

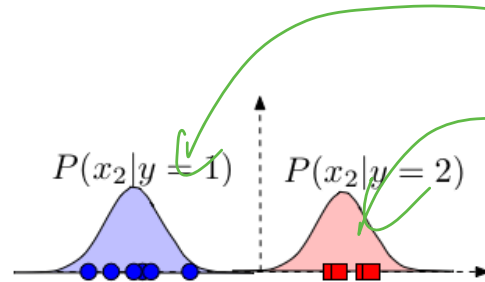
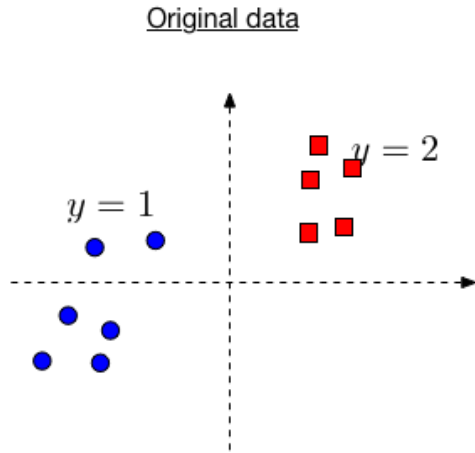


$$\mu_{\alpha,y} = \frac{\sum_{i=1}^n x_i[\alpha] \mathbf{1}(y_i = y)}{\sum_{i=1}^n \mathbf{1}(y_i = y)}$$

# Case study

Estimate the mean/std parameter  $\mu_{\alpha,y}, \sigma_{\alpha,y}$  via **MLE**:

$$P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha,y}^2)$$

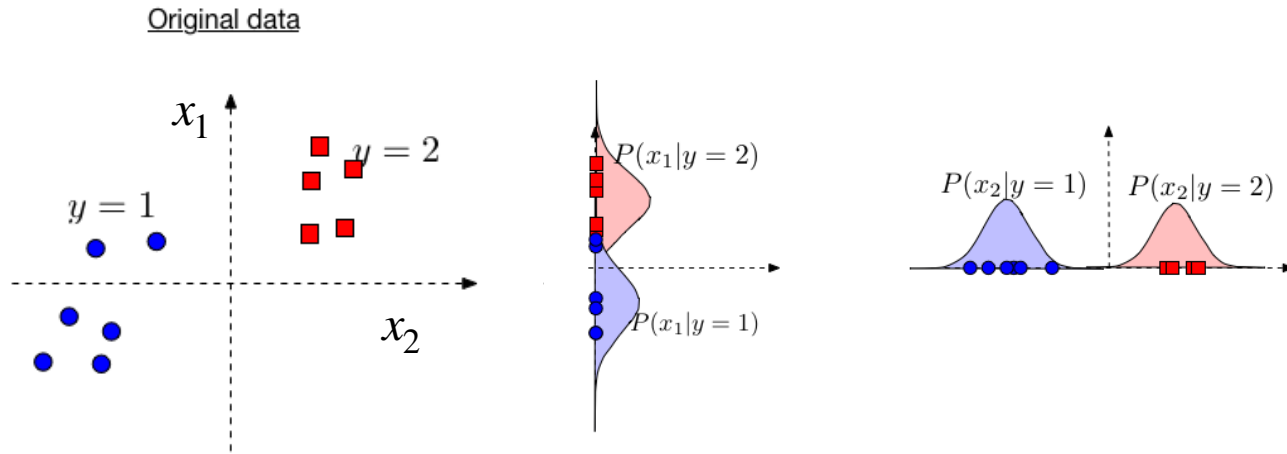


$$\mu_{\alpha,y} = \frac{\sum_{i=1}^n x_i[\alpha] \mathbf{1}(y_i = y)}{\sum_{i=1}^n \mathbf{1}(y_i = y)}$$
$$\sigma_{\alpha,y}^2 = \frac{\sum_{i=1}^n (x_i[\alpha] - \mu_{\alpha,y})^2 \mathbf{1}(y_i = y)}{\sum_{i=1}^n \mathbf{1}(y_i = y)}$$

# Case study

Formulate the joint conditional distribution

$$P(x | y) = \prod_{\alpha=1}^d P(x[\alpha] | y)$$



$$a \sim N(\mu_a, \sigma_a)$$

$$b \sim N(\mu_b, \sigma_b), \text{ Assume } a \perp b$$

## Case study

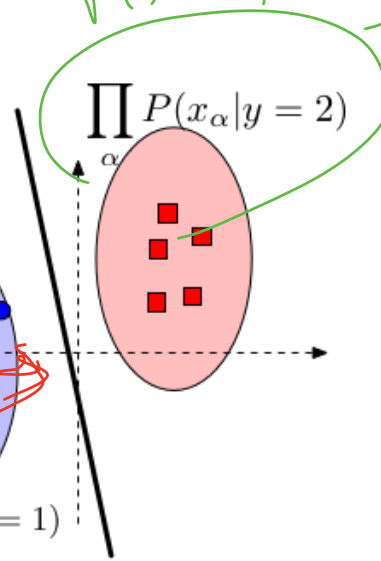
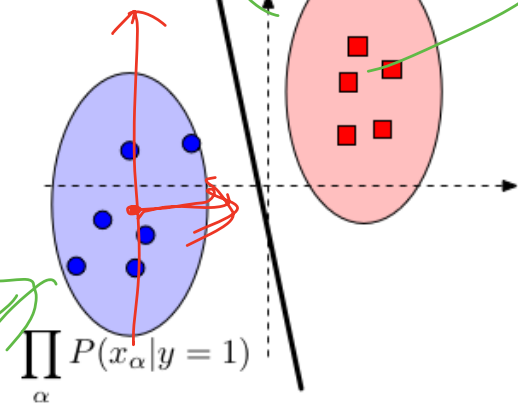
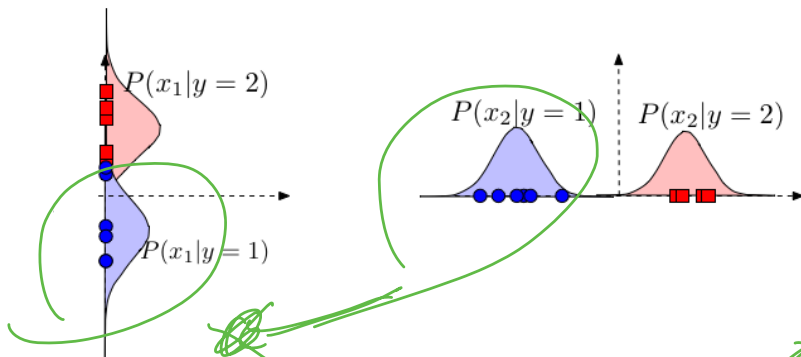
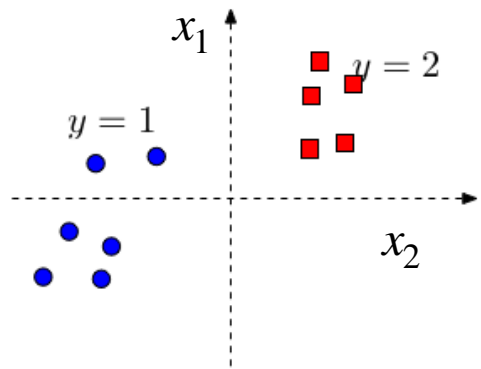
Formulate the joint conditional distribution

$$\begin{bmatrix} a \\ b \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & 0 \\ 0 & \sigma_b^2 \end{bmatrix} \right)$$

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

$$P(x^{(1)}|y=2) \cdot P(x^{(2)}|y=2)$$

Original data



# Outline

1. General formulation of Naive Bayes



2. Example



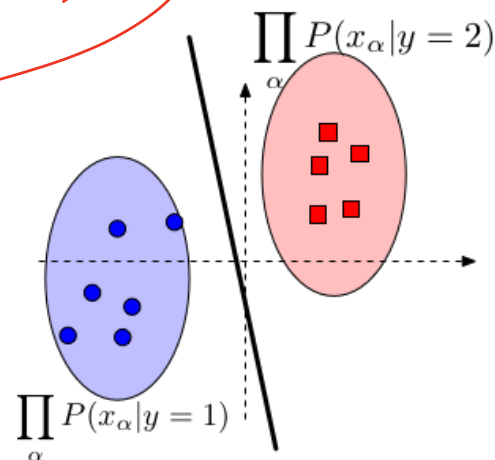
3. Connection to linear classifier

# Gaussian Naive Bayes induces a linear classifier

When  $P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha}^2)$

i.e., give  $\alpha$ , STD  $\sigma_{\alpha}$  is the same across all labels,

$$\underset{y}{\operatorname{argmax}} P(y|x)$$



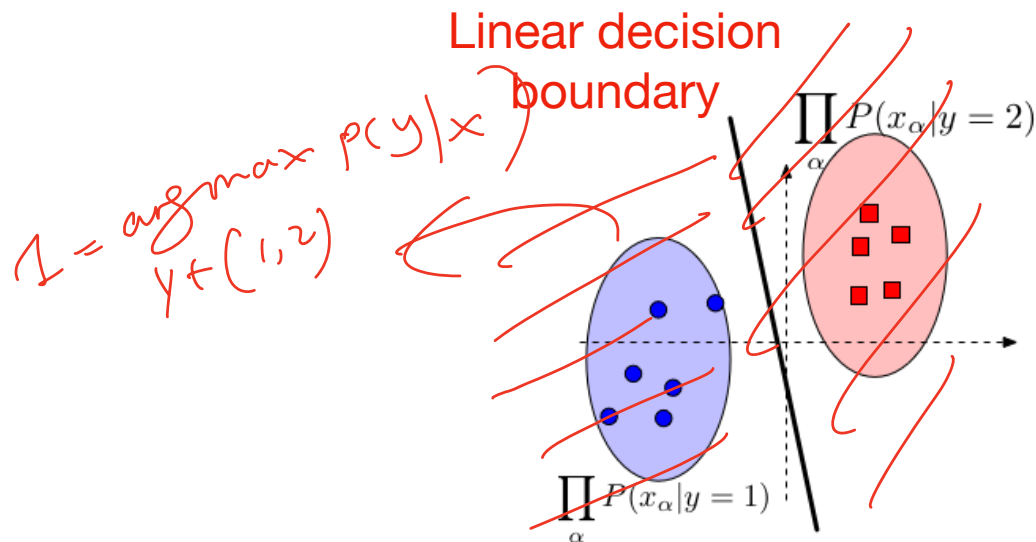


# Gaussian Naive Bayes induces a linear classifier

When  $P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha}^2)$

i.e., give  $\alpha$ , STD  $\sigma_{\alpha}$  is the same across all labels,

$\mathcal{I} = \arg \max_{y \in \{1,2\}} P(y|x)$



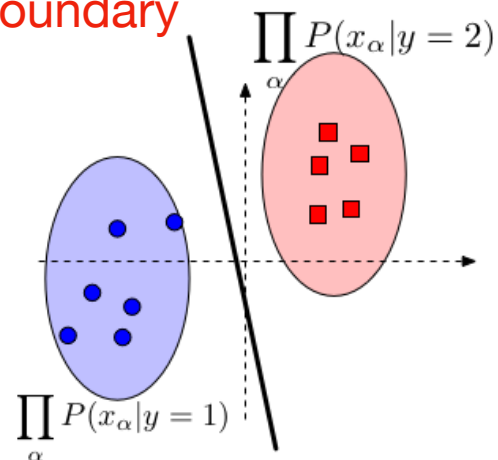
# Gaussian Naive Bayes induces a linear classifier

$$\text{When } P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha}^2)$$

i.e., give  $\alpha$ , STD  $\sigma_{\alpha}$  is the same across all labels,

$\exists w, b$  (i.e., a hyperplane), such that:

Linear decision  
boundary



# Gaussian Naive Bayes induces a linear classifier

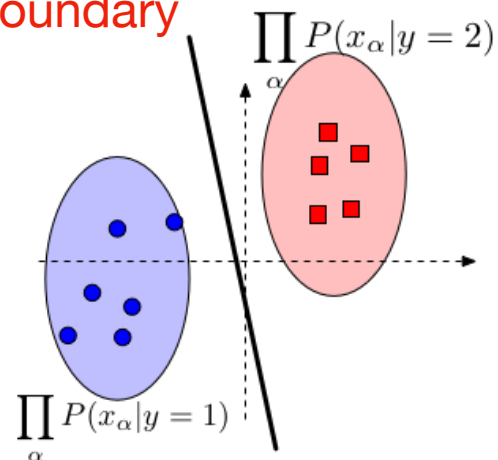
$$\text{When } P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha}^2)$$

i.e., give  $\alpha$ , STD  $\sigma_{\alpha}$  is the same across all labels,

$\exists w, b$  (i.e., a hyperplane), such that:

$$\arg \max_y P(y | x) = 1 \iff w^T x + b > 0$$

Linear decision  
boundary



# Gaussian Naive Bayes induces a linear classifier

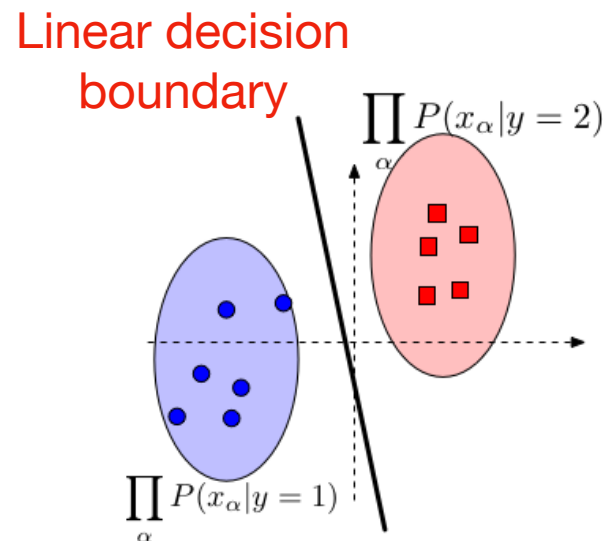
$$\text{When } P(x[\alpha] | y) = \mathcal{N}(\mu_{\alpha,y}, \sigma_{\alpha}^2)$$

i.e., give  $\alpha$ , STD  $\sigma_{\alpha}$  is the same across all labels,

$\exists w, b$  (i.e., a hyperplane), such that:

$$\arg \max_y P(y | x) = 1 \iff w^{\top} x + b > 0$$

(Try this out in HW3)




# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$

*Bayes Rule*



# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$



Easy to estimate via  
MLE

# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$

The Naive Bayes  
assumption

Easy to estimate via  
MLE

$$P(x|y) = \prod_{\alpha=1}^d P(\underline{x[\alpha]}|y)$$

*- d i n s*

# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$

The Naive Bayes  
assumption

Easy to estimate via  
MLE

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

Estimate each  $P(x[\alpha]|y)$  via MLE (or MAP)



# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$

The Naive Bayes  
assumption

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

Easy to estimate via  
MLE

NB classifier:  $\arg \max_y P(y|x)$

Estimate each  $P(x[\alpha]|y)$  via MLE (or MAP)

# Summary for today

We start from Bayes rules:

$$P(y|x) \propto P(x|y)P(y)$$

The Naive Bayes  
assumption

$$P(x|y) = \prod_{\alpha=1}^d P(x[\alpha]|y)$$

Estimate each  $P(x[\alpha]|y)$  via MLE (or MAP)

Easy to estimate via  
MLE

NB classifier:  $\arg \max_y P(y|x)$

Take-home Q: Perceptron  
VS NB classifier