

# Naive Bayes Classifier

- Generative Counter-part to Logistic Regression.
- This also outputs the  $P(y|x)$ . (But this uses Baye's Thm to do that).
- Simple, yet powerful → used in document classification  
→ spam filtering.
- This makes use of the strong (Naive) Conditional Independence assumption between the features of a given a label.

Suppose →

$$D = \left( \vec{x}, \vec{y} \right) = \left\{ \left( \vec{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^n$$

This enables us to express joint probabilities of individual features given a label.

$$\vec{x} = (x_1, x_2, \dots, x_m)$$

For binary classification

$$\begin{aligned}
 P(\vec{x} | y) &= P(x_1, x_2, \dots, x_m | y) \\
 &= P(x_1 | y) \cdot P(x_2 | y) \cdots P(x_m | y) \quad --- (*) 
 \end{aligned}$$

→ NB classifier predicts probability,  $P(y|x)$ , of a class label  $y$  given a feature vector  $x$

using Bayes' Theorem.

$$P(y|x) = \frac{P(x,y)}{P(x)} = \frac{\underbrace{P(x|y) \cdot P(y)}_{\substack{\text{class conditional} \\ \text{density}}} \cdot \underbrace{P(x)}_{\substack{\text{class prior} \\ \downarrow \text{Evidence}}}}{\underbrace{P(x)}_{\substack{\downarrow \\ \text{Posterior} \\ \text{Probability}}}}$$

If the different possible labels are

$\{y_1, y_2, \dots, y_k\}$ , then →

$$P(x) = \sum_{r=1}^k P(x, y_r) = \sum_{r=1}^k P(x|y_r) \cdot P(y_r) \quad \text{--- (†)}$$

Now, using the naive Baye's Assumption (⊗) and (†)  
we have

$$\begin{aligned} P(y=y_c|x) &= \frac{P(x|y_c) \cdot P(y_c)}{P(x)} \\ &= \frac{P(y_c) \cdot P(x_1|y_c) \cdot P(x_2|y_c) \cdots P(x_m|y_c)}{\sum_{r=1}^k P(x|y_r) \cdot P(y_r)} \end{aligned}$$

$$\frac{P(y_c) \cdot P(x_1|y_c) \cdot P(x_2|y_c) \cdots P(x_m|y_c)}{\sum_{r=1}^k P(y_r) \cdot P(x_1|y_r) \cdot P(x_2|y_r) \cdots P(x_m|y_r)}$$

$$P(y_c|x) = \frac{P(y_c) \cdot \prod_{j=1}^m P(x_j|y_c)}{\sum_{r=1}^k P(y_r) \cdot \prod_{j=1}^m P(x_j|y_r)}$$

The parameters of NB classifiers →

$'k'$  Prior probabilities :  $P(y_1), P(y_2) \dots P(y_k)$

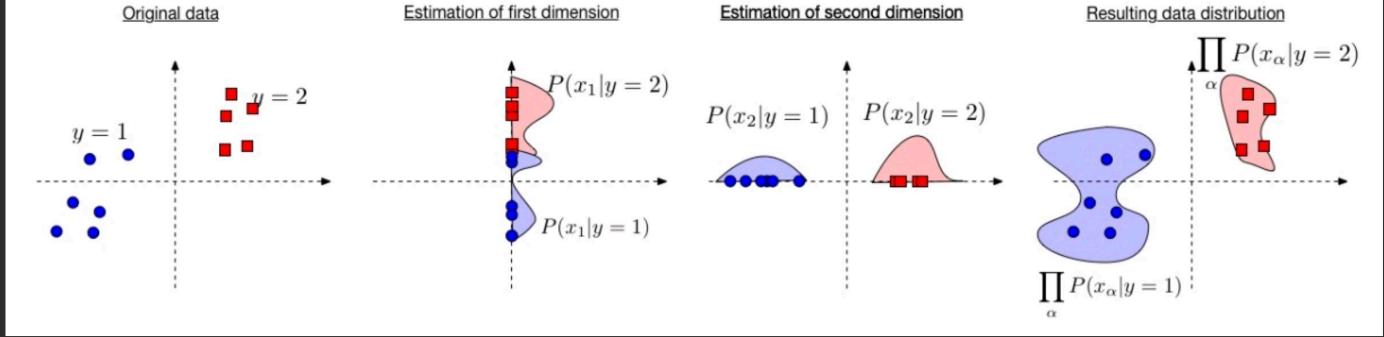
$(kxm)$  Class conditional densities  $\stackrel{\text{CD given } y=y_1}{\leftarrow} P(x_1|y_1), P(x_2|y_1) \dots P(x_m|y_1)$

$\stackrel{\text{CD given } y=y_2}{\leftarrow} P(x_1|y_2), P(x_2|y_2) \dots P(x_m|y_2)$

$\stackrel{\text{conditional density given } y=y_k}{\leftarrow} P(x_1|y_k), P(x_2|y_k) \dots P(x_m|y_k)$

\* The number of parameters for each conditional density varies and depends on its mathematical form.

# NB schematic



## Modelling Conditional Densities

$$P(x_j | y_c)$$

This depends on nature of feature  $x_j$

### Binary Feature

$$\rightarrow x_j \in \{0, 1\}$$

$$P(x_j | y_c) \sim \text{Bernoulli}(\mu_{jc})$$

Generalisation

### Categorical Feature

$$\rightarrow x_j \in \{v_1, v_2, \dots, v_e\}$$

$$P(x_j | y_c) \sim \text{Cat}(e, \mu_{j1c}, \mu_{j2c}, \dots, \mu_{je})$$

### Multinomial Feature

$\rightarrow$  Ex → Word count  $c_i$  as features with additional constraints that

$$\sum_{i=1}^m c_i = l, \text{ the length of the sequence they represent.}$$

$$\vec{x} = (c_1, c_2, \dots, c_m) \quad c_i \in \mathbb{N}$$

$$c_1 + c_2 + \dots + c_m = l.$$

$$P(\vec{x} | y_c) \sim \text{Multinomial}(l, \mu_{1c}, \mu_{2c}, \dots, \mu_{mc})$$

### Continuous Feature

→ features are real nos.

$$x_j \in \mathbb{R}$$

$$P(x_j | y_c) \sim \mathcal{N}(\mu_{jc}, \sigma_{jc})$$

→ We need to estimate parameters of relevant distributions, one for each feature, for each class label.

→ We denote the set of parameters as a single vector  $\vec{w}$  (of unknown dimension as this varies from problem to problem).

$\vec{w}$ : Set of all parameters  $\rightarrow$  Prior density  
+  
Class conditional density.

- Bernoulli Distribution.

$$P(x_j | y_c) \sim \text{Bernoulli}(u_{j|c})$$

$u_j \in \{0, 1\}$       Parameters:  $u_{j|c}$

$$P(u_j = 1 | y_c) = u_{j|c} \quad P(u_j = 0 | y_c) = 1 - u_{j|c}$$

$$P(u_j | y_c ; u_{j|c}) = u_j^{x_j} \cdot (1 - u_{j|c})^{(1-x_j)}$$

If there are ' $k$ ' classes:  
 • For each binary feature, we need ' $k$ ' parameters  $(u_{j|c})_{c=1}^k$   
 • For ' $s$ ' binary features, we need  $(k \times s)$  parameters

- Categorical Distribution.

Generalised Version of Bernoulli

Suppose  $x_j \in \{v_1, v_2, \dots, v_e\}$  (There are  $e$  possible categories)

Parameters:  $\vec{u}_{j|c} = (u_{j|c}, u_{j|2c}, \dots, u_{j|ec})$

be the parameter vector for  $P(u_j | y_c)$

$$\text{Then } \rightarrow P(u_j = v_1 | y_c ; e, \vec{u}_{j|c}) = u_{j|1c}$$

$$P(u_j = v_2 | y_c ; e, \vec{u}_{j|c}) = u_{j|2c}$$

$$P(y_j = v_e | y_c, e, \vec{u}_{jc}) = u_{jec}$$

Then →

$$P(y_j | y_c, e, \vec{u}_{jc}) = u_{j1c}^{1_{(y_j=v_1)}} \cdot u_{j2c}^{1_{(y_j=v_2)}} \cdots u_{jec}^{1_{(y_j=v_e)}}$$

This forces the constraint →  $\sum_{e=1}^E u_{jec} = 1$   
 (or  $\text{sum}(\vec{u}_{jc}) = 1$ )

- If there are  $|k|$  classes →
  - for each class  $c$  we need  $\text{len}(\vec{u}_{jc})$
  - For all  $|k|$  classes →  $k \times \text{len}(\vec{u}_{jc})$
  - for all categorical features (if all  $m$ )  
 are categorical)
 
$$k \cdot \sum_{j=1}^m \text{len}(\vec{u}_{jc})$$

## Multinomial Distribution.

Ex → When  $\vec{x}$  is a count vector, i.e. each component  $x_j$  is a "count of appearance" in the object it represents and  $\sum_{j=1}^m (\vec{x}) = \sum_{j=1}^m x_j = l$ ,

which is the length of the object.

Here, multinomial distribution is useful to model

$$P(\vec{x} | y_c) \quad \vec{x} = (x_1, x_2, \dots, x_m)$$

Parameters:  $l, \mu_{1c}, \mu_{2c}, \dots, \mu_{mc}$   
(for  $P(\vec{x} | y_c)$ )

constraint:  $\sum_{i=1}^m x_i = l \quad x_i \in \mathbb{N}$ .

$$P(\vec{x} | y_c, l, \vec{\mu}_c) = \frac{l!}{x_1! x_2! \dots x_m!} \prod_{j=1}^m \mu_j^{x_j}$$

where

$$\vec{\mu}_c = (\mu_{1c}, \mu_{2c}, \dots, \mu_{mc})$$

- If there are  $k$  classes:

- parameters per class:  $m = \text{len}(\vec{\mu}_c)$   
(for class  $y_c$ )
- Total:  $k \times m$

## Gaussian Distribution

→ when  $y_j$  is a real no. (or continuous feature)

Parameters :  $\mu_{j_c}, \sigma_{j_c}^2$   
 (for  $P(y_j | y_c)$ )

$$-\frac{1}{2} \left( \frac{(y_j - \mu_{j_c})^2}{\sigma_{j_c}^2} \right)$$

$$P(y_j | y_c, \mu_{j_c}, \sigma_{j_c}^2) = \frac{1}{\sqrt{2\pi} \sigma_{j_c}} e$$

For each class, no. of parameters :  $2 \times m$

for  $K$  classes :  $2 \times m \times K$

$\overrightarrow{\mu_c}, \overrightarrow{\sigma_c}$   
 ↗ each of length  $m$

This is the one-dimensional form.

We can also look at multi-variate  
 form to gain direct access to

$$P(\vec{y} | y_c)$$

- Multi-variate Gaussian:

Parameters:  $(\vec{\mu}_c)_{m \times 1}$ ,  $(\Sigma_c)_{m \times m}$

$\downarrow$                              $\downarrow$   
 Mean vector                      Covariance Matrix

In Naive Bayesian setting, where features are independent of each other given a label, we have the non-diagonal entries of

$\Sigma_c$  to be zero.

So,  $(\Sigma_c)_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ (\sigma_c)_{ii}^2 & \text{if } i = j \end{cases}$

$$P(\vec{x} | y_c, \vec{\mu}_c, \Sigma_c) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} e^{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

→ Assume  $C_k$  classes:

- Parameters per class:  $2 \times m$

$$\vec{\mu}_c, \Sigma_c \rightarrow \begin{bmatrix} (\bar{\mu}_c)_1 & 0 \\ (\bar{\mu}_c)_2 & \dots \\ \vdots & \vdots \\ 0 & (\bar{\sigma}_c)_{mm} \end{bmatrix}$$

- Total:  $2 \times m \times k$ .

## Loss Function and Optimization.

→ Maximum Likelihood  
 → Maximum log likelihood ) Estimators

so, basically given a sample examples  
 and labels →

$$\left( \vec{x}^{(i)}, y^{(i)} \right)^n$$

$$L(\vec{w}) = \prod_{i=1}^n P(\vec{x}^{(i)}, y^{(i)}; \vec{w}) \longrightarrow \text{Likelihood}$$

$$l(\vec{w}) = \sum_{i=1}^n \log \left( P(\vec{x}^{(i)}, y^{(i)}; \vec{w}) \right) \quad \text{log likelihood}$$

Our goal is to find the set of parameters  
 aka  $\vec{w}$  that maximizes  $L(w)$  or

$$l(\vec{w})$$

## → Loss function

$$l(\vec{w}) = \sum_{i=1}^n \log \left( P(x^{(i)}, y^{(i)}; \vec{w}) \right)$$

(

We want to maximize  
the log likelihood.

)

Instead we can minimize  
negative log-likelihood →

$$J(\vec{w}) = -l(\vec{w})$$

$$= -\sum_{i=1}^n \log \left( P(x^{(i)}, y^{(i)}; \vec{w}) \right)$$



Instead of doing that, we focus on maximizing  
the log likelihood.

|

$$\begin{aligned}
 L(\vec{w}) &= \sum_{i=1}^n \log \left( P\left(\vec{x}^{(i)}, \vec{y}^{(i)} ; \vec{w}\right) \right) \\
 &= \sum_{i=1}^n \log \left( P\left(\vec{x}^{(i)} \mid \vec{y}^{(i)} ; \vec{w}\right) \cdot P\left(\vec{y}^{(i)} ; \vec{w}\right) \right)
 \end{aligned}$$

Assuming →  $\vec{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_m^{(i)})$   
 and using N.B assumption.

$$\begin{aligned}
 &= \sum_{i=1}^n \log \left( P\left(\vec{y}^{(i)} ; \vec{w}\right) \cdot \prod_{j=1}^m P\left(x_j^{(i)} \mid \vec{y}^{(i)} ; \vec{w}\right) \right) \\
 &= \sum_{i=1}^n \log \left( P\left(\vec{y}^{(i)} ; \vec{w}\right) \right) \quad \text{Prior probability} \\
 &\quad + \sum_{i=1}^n \sum_{j=1}^m \log \left( P\left(x_j^{(i)} \mid \vec{y}^{(i)} ; \vec{w}\right) \right) \quad \text{This comes from empirical probability from the samples}
 \end{aligned}$$

This depends on the probability distribution, depending on the type of features we are dealing with (as described above) ↗ class-conditional density.

## → Optimization and Parameter estimation

Parameter estimation can be done in these 3 steps → (Maximum Likelihood Estimate)

- ( $\dagger$ )
- ① Calculate partial derivative of  $L(\vec{\omega})$  w.r.t. each parameter ( $\omega_i$ ) .
  - ② Set partial derivative to '0', which is the condition at maxima.
  - ③ Solve the resulting eq<sup>u</sup> to obtain the parameter value ( $\hat{\omega}_i$ ) .

### ① Estimating Prior Probabilities

$$P(y^{(i)}; \vec{\omega})$$

let

$$\omega_{y^{(i)}} := P(y^{(i)} = y_c) = \frac{|\{i \in [n] \mid y^{(i)} = y_c\}|}{n}$$

This is just  
the empirical  
probability .

$$= \frac{\sum_{i=1}^n \mathbb{1}(y^{(i)} = y_c)}{n}$$

$$= \frac{\# \text{examples with label } y_c}{n}$$

## ② Estimating Class-Conditional Densities.

### Bernoulli Distribution.

(Pg no. - 6)

$$L(\vec{w}) = \sum_{i=1}^n \log \left( P(y_i^{(i)}; \vec{w}) \right) + \sum_{i=1}^n \sum_{j=1}^m \log \left( P(x_j^{(i)} | y_i^{(i)}; \vec{w}) \right)$$

No.  $w_j$ ,  
 $x_j \in \{0, 1\}$

$$P(x_j | y_i; w_j) = \frac{x_j}{w_j} \frac{(1-x_j)}{(1-w_j)}$$

$$\begin{aligned} L(w) &= \sum_{i=1}^n \log \left( w_{y_i^{(i)}} \right) + \sum_{i=1}^n \sum_{j=1}^m \log \left( \frac{x_j^{(i)}}{w_j y_i^{(i)}} \frac{(1-x_j^{(i)})}{(1-w_j y_i^{(i)})} \right) \\ &= \sum_{i=1}^n \log \left( w_{y_i^{(i)}} \right) + \\ &\quad \sum_{i=1}^n \sum_{j=1}^m \left( x_j^{(i)} \cdot \log \left( \frac{1}{w_j y_i^{(i)}} \right) + (1-x_j^{(i)}) \cdot \log \left( \frac{1}{1-w_j y_i^{(i)}} \right) \right) \end{aligned}$$

(Parameters for layer  $y_r$ :  $w_{1y_r}, w_{2y_r}, w_{3y_r}, \dots, w_{ny_r}$ )

Now, we employ  $(\dagger) \rightarrow$  for  $j=t$  and  $y^{(i)}=y_r$

STEP 1:

$$\frac{\partial (\ell(\vec{w}))}{\partial w_{ty_r}} = \frac{\partial}{\partial w_{ty_r}} \left( \ell(\vec{w}) \right) = \frac{\partial}{\partial w_{ty_r}} \left( \sum_{i=1}^n \log(w_{j y^{(i)}}) \right)^0 +$$

*There are constants.*

$$+ \frac{\partial}{\partial w_{ty_r}} \sum_{i=1}^n \sum_{j=1}^m \left( x_j^{(i)} \cdot \log(w_{j y^{(i)}}) + (1 - x_j^{(i)}) \log(1 - w_{j y^{(i)}}) \right)$$

we only keep the relevant terms with  $(y^{(i)}=y_r)$

$$= \sum_{i=1}^n \underbrace{1}_{(y^{(i)}=y_r)} \left( \sum_{j=1}^m \frac{\partial}{\partial w_{ty_r}} \left( x_j^{(i)} \cdot \log(w_{j y_r}) + (1 - x_j^{(i)}) \log(1 - w_{j y_r}) \right) \right)$$

$$= \sum_{i=1}^n \underbrace{1}_{(y^{(i)}=y_r)} \left( \frac{x_t^{(i)}}{w_{ty_r}} - \frac{1 - x_t^{(i)}}{1 - w_{ty_r}} \right)$$

*Only terms with  $j=t$  survive.*

STEP 2: setting  $\frac{\partial}{\partial w_{ty_r}} (\ell(\vec{w})) = 0$

$$\Rightarrow \sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} \left( \frac{x_t^r}{w_{ty_r}} - \frac{1-x_t^r}{1-w_{ty_r}} \right) = 0$$

$$\Rightarrow \frac{\sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} \left( x_t^r (1-w_{ty_r}) - (1-x_t^r) w_{ty_r} \right)}{w_{ty_r} \cdot (1-w_{ty_r})} = 0$$

$$= \sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} \left( x_t^r (1-w_{ty_r}) - (1-x_t^r) w_{ty_r} \right) = 0$$

$$= \sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} \left( x_t^r - w_{ty_r} \right) = 0$$

$$= \sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} x_t^r = \sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} w_{ty_r}$$

$$w_{ty_r} = \frac{\sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)} x_t^r}{\sum_{i=1}^n \mathbb{1}_{(y^{(i)} = y_i)}}$$

For binary type features  
 i.e. if  $x_t^{(i)} \in \{0, 1\}$   
 we can simplify this even  
 more →

$$w_{t y_r} = \underbrace{\sum_{i=1}^n \underbrace{1}_{\text{if } y_t^{(i)} = y_r} x_t^{(i)}}_{\sum_{i=1}^n 1_{(y_t^{(i)} = y_r)}} \rightarrow \begin{array}{l} \# \text{ of examples} \\ \text{with } x_t^{(i)} = 1 \\ \text{and } y_t^{(i)} = y_r \end{array}$$

↓

$\# \text{ of examples with } y_t^{(i)} = y_r$

PROBLEM WITH ZERO COUNT :

→ Now it becomes problematic if  $\sum_{i=1}^n 1_{(y_t^{(i)} = y_r)} x_t^{(i)} = 0$   
 (i.e. if there is no example with  $x_t^{(i)} = 1$  and  $y_t^{(i)} = y_r$  at the same time)

This makes  $w_{t y_r} = 0 \Rightarrow P(x_t | y_r) = 0$   
 $\Rightarrow P(y=y_r | \vec{x}) = 0$

FIXING PROBLEM WITH ZERO COUNT: when c=1

We use something called "Laplace Smoothing".

We can control this by adding a  $(+c)$  to numerator and  $(+2c)$  to the denominator.

$$w_{tyr} = \frac{\left( \sum_{i=1}^n 1_{(y_i^r = y_r)} x_t^r \right) + c}{n_r + 2c}$$

$$n_r = \left( \sum_{i=1}^n 1_{(y_i^r = y_r)} \right)$$

$c$ : This is a hyperparameter that helps control overfitting.

Too high value of  $(c)$  leads to underfitting.

• Categorical Distribution - Pg no. - 6.

Suppose  $u_j \in \{v_1, v_2, \dots, v_e\}$  There are  $e$  possible values feature  $u_j$  can take.

Parameters:  $\vec{\mu}_{jc} = (\mu_{j1c}, \mu_{j2c}, \dots, \mu_{jec})$   
be the parameter vector for  $P(u_j | y_c)$

$$\text{Then } \rightarrow P(u_j = v_1 | y_c; e, \vec{\mu}_{jc}) = \mu_{j1c}$$

$$P(u_j = v_2 | y_c; e, \vec{\mu}_{jc}) = \mu_{j2c}$$

$$P(u_j = v_e | y_c; e, \vec{\mu}_{jc}) = \mu_{jec}$$

Then  $\rightarrow$

$$P(u_j | y_c; e, \vec{\mu}_{jc}) = \mu_{j1c}^{1_{(u_j=v_1)}} \cdot \mu_{j2c}^{1_{(u_j=v_2)}} \cdots \mu_{jec}^{1_{(u_j=v_e)}}$$

This forces the constraint  $\sum_{c=1}^e \mu_{j1c} = 1$ .

$$\left( \text{or } \sum \left( \vec{\mu}_{jc} \right) = 1 \right)$$

Here, for a fixed  $y^{(i)} = y$  :

(i.e. for a fixed class  $\rightarrow$ )

$$\vec{w}_r = \left( \vec{\mu}_{1r}, \vec{\mu}_{2r}, \dots, \vec{\mu}_{mr} \right)$$

where  $\rightarrow$

$$\vec{M}_{j,r} = \begin{pmatrix} M_{j1r} & M_{j2r} & \dots & M_{jmr} \end{pmatrix}$$

Total no. of parameters  $\rightarrow n \times m \times e$

(if each parameter has  $e$  choices)

$$M_{jv,y_r} = P(u_j=v \mid y_r) \\ = \frac{\sum_{i=1}^n I(y_i^{(i)} = y_r) \cdot I(u_j^{(i)} = v)}{n_r + ce}$$

When  $c=1 \rightarrow$  Laplace Smoothing.

$$n_r = \left( \sum_{i=1}^n I(y_i^{(i)} = y_r) \right)$$

• Multinomial Distribution. Pg - 8

$$P(\vec{x} | y_c) \quad \vec{n} = (n_1, n_2, \dots, n_m)$$

Parameters :  $\ell, M_{1c}, M_{2c}, \dots, M_{mc}$   
 (for  $P(\vec{x} | y_c)$ )

constraint :  $\sum_{i=1}^m n_i = \ell \quad n_i \in \mathbb{N}$ .

$$P(\vec{x} | y_c, \ell, \vec{M}_c) = \frac{n!}{x_1! x_2! \dots x_m!} \prod_{j=1}^m M_j^{x_j}$$

where

$$\vec{M}_c = (M_{1c}, M_{2c}, \dots, M_{mc})$$

(for bmed  $y_c^{(i)} = y_c$ )

So,

$$M_{j\gamma} = \left( \sum_{i=1}^m \frac{1}{(y_c^{(i)} - y_\gamma)} x_j^{(i)} \right) + c$$

$$\left( \sum_{i=1}^m \frac{1}{(y_c^{(i)} - y_\gamma)} \left( \sum_{j=1}^m x_j^{(i)} \right) \right) + mc$$

↳  $c_{z1}$  Laplace Smoothing.

• Gaussian / Normal Distribution

Pg - 9

Parameters :  $\mu_{j^c}, \sigma_{j^c}^2$        $u_j \in \mathbb{R}$   
 (for  $P(u_j | y_c)$ )

$$-\frac{1}{2} \left( \frac{u_j - \mu_{j^c}}{\sigma_{j^c}} \right)^2$$

$$P(u_j | y_c, \mu_{j^c}, \sigma_{j^c}^2) = \frac{1}{\sqrt{2\pi} \sigma_{j^c}} e$$

$$\mu_{j^r} = \frac{\sum_{i=1}^n I(y^{(i)} = y_r) u_j^{(i)}}{n_r}$$

$$\sigma_{j^r}^2 = \frac{\sum_{i=1}^n I(y^{(i)} = y_r) (u_j^{(i)} - \mu_{j^r})^2}{n_r}$$

$$n_r = \sum_{i=1}^n I(y^{(i)} = y_r)$$

## Generating output

estimated weights

So, now the question remains →

Now that we know the parameters.

for any  $\vec{u}$  (an input) and any class  $c$  among the  $C$  classes.

→ Posterior.

$$y = \arg \max_{y_c} P(y_c | \vec{u}; \hat{w})$$

likelihood  $P(\vec{u} | y_c; \hat{w})$  prior  $P(y_c; \hat{w})$   
 $\frac{P(\vec{u} | y_c; \hat{w})}{P(\vec{u}; \hat{w})}$   
 evidence

We can ignore the denominator  
 (because it is constant, as input is fixed)

$$= \arg \max_{y_c} P(\vec{u} | y_c; \hat{w}) \cdot P(y_c; \hat{w})$$

Simplifying this by Naive Bayes Assumption →

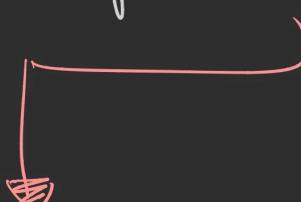
$$y = \arg \max_{y_c} P(y_c; \hat{w}) \prod_{j=1}^m P(x_j | y_c; \hat{w})$$

We can argue the product weight become very small instead on the log

$$y = \underset{y_c}{\operatorname{argmax}} \log \left( P(y_c; \vec{w}) \right) + \sum_{j=1}^m \log \left( P(x_j | y_c; \vec{w}) \right)$$

So, for any example  $\vec{x}$ , the NB classifier outputs the class with maximum value amongst  $\rightarrow (f(y_1), \dots, f(y_k))$

where  $f(y_c) = \log \left( P(y_c; \vec{w}) \right) + \sum_{j=1}^m \log \left( P(x_j | y_c; \vec{w}) \right)$



Note this is NOT probability.

We are looking at the maximum in the 'log-space'.

$$P(y_c | \vec{x}; \vec{w}) = \frac{P(y_c; \vec{w}) P(\vec{x} | y_c; \vec{w})}{P(\vec{x}; \vec{w})}$$