

# Lecture 6/11

- Recap Generative models  $\Rightarrow$  Supervised Classification
  - Last time Gen Model (Gaussian)  $\Rightarrow$  GDA
  - Gen Models ( ? other dist )  $\Rightarrow$  ? Discriminant Analysis
- Naive Bayes
  - Smoothing estimates.
- HW2B done

Generative Model: fitting curve (params) to data  $X \Rightarrow P(X)$   
distrib

$P(X)$  • Gaussian  $\mathcal{N}(\mu, \Sigma)$ : find  $\mu, \Sigma$  such that  $X \sim \mathcal{N}(\mu, \Sigma)$

$\mathcal{N}(\mu, \Sigma)$  generator for datapoints

← appears dist as ...

INDEP  $N$  times, prob observing  $(x_1, x_2, \dots, x_N)$  datapoints

given data

is as much as possible.

(best true)  $\Rightarrow \mu = \text{Mean}[X_{i=1:N}]$   $\Sigma = \text{COVAR}(X)$

$P(X)$  • model = Poisson  $\Rightarrow$  # events per fixed interval (assuming rough constant)

$\lambda$  param = expected AVG

$X$  = # events in an interval

" # buses at stop per hour "  $\rightarrow$  avg = 6

" # Trump tweets a day "

" # accidents in Boston per week "

$$P[X=k] = \frac{\lambda^k / k!}{e^{-\lambda}} \rightarrow \text{normalizer}$$

$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \underline{\text{Taylor}} e^{\lambda} \Rightarrow$  need to divide with  $e^{\lambda}$  to sum  $= 1$

- mult-dim Poisson (D dim) ?

$P(x)$  • exponential

$P(x)$  • multinomial Roll die  $k$  faces (1, 2, 3, ...,  $k$ )  $N$  times  
 prop  $p_1, p_2, p_3, \dots, p_k$   $\sum p_i = 1$

outcome  $x_1, x_2, \dots, x_k$  observed  $\sum x_k = N$   
 $X$  # 1 # 2 #  $k$

$$P(x_1, x_2, \dots, x_k) = \binom{N}{x_1, x_2, \dots, x_k} p_1^{x_1} \cdot p_2^{x_2} \cdot \dots \cdot p_k^{x_k}$$

multinom  $\frac{N!}{x_1! \cdot x_2! \cdot \dots \cdot x_k!}$

Check

$$\sum_{x_1 + x_2 + \dots + x_k = N} \binom{N}{x_1, \dots, x_k} \prod_{t=1}^k p_t^{x_t} = 1$$

Fit curve (param) to data  $X \Leftrightarrow$  Find best params.

- gaussian :  $\mu, \Sigma$

- multinomial :  $p_1, p_2, \dots, p_k$

- poisson :  $\lambda_1, \lambda_2, \dots, \lambda_k$

→ model of density  
data density

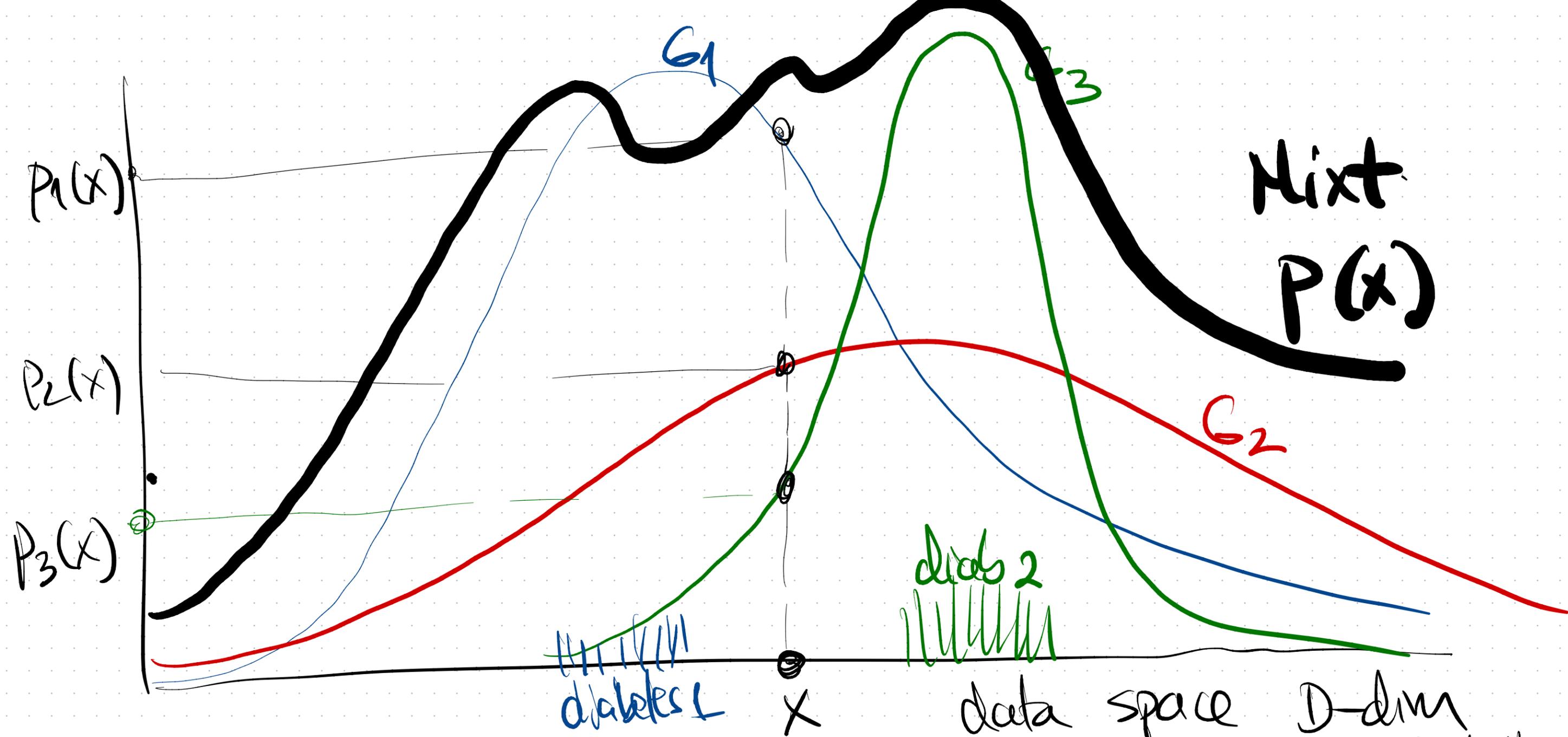
$P(x_i)$  → datapoint  
= probability to observe that point

$P(x)$  • MIXTURE (fixed weights)  
(Gaussian, Multinomial, Poisson)

HW 3: Mixture of Gaussians

$k=3$  gaussians  
each  $D$ -dim

$$P(x) = \underbrace{w_1}_{20\%} \cdot \underbrace{N_1(\mu_1, \Sigma_1)}_{D\text{-dim}} + \underbrace{w_2}_{50\%} \cdot \underbrace{N_2(\mu_2, \Sigma_2)}_{D\text{-dim}} + \underbrace{w_3}_{30\%} \cdot \underbrace{N_3(\mu_3, \Sigma_3)}_{D\text{-dim}}$$



Mixture adv: ability to fit data with multiple "hills"  
 "bumps"

# Generative Models $\Rightarrow$ supervised classification.

① For each class  $y_1, y_2, \dots, y_L$  ( $L$  labels)

fit a model on  $X$  that class

$$P_1(x) = P(x | Y=1)$$

model for class  $y=1$

$$P_2(x) = P(x | Y=2)$$

model for class  $y=2$

$$P_L(x) = P(x | Y=L)$$

model for class  $y=L$

not ind. dist.  $X =$  filtered by class

models trained per class

label  $y \in \{1, 2, \dots, L\}$

Predict ARGMAX

prior (labels)

$$P(z | y_e) \cdot P(y = y_e)$$

normalized

gauss  $\leftarrow$

not nec. same model  $\leftarrow$

exp  $\leftarrow$

②

$z =$  test point  $\Rightarrow$  predict label  $y \in \{1, 2, \dots, L\}$

$$P(Y | z) = \frac{P(z | Y) \cdot P(Y)}{P(z)}$$

calculate

for each  $y = l$

$$P(Y=l | z) =$$

product:  $P(Y=l|Z)$  = confidence label  
over all  $l=1:L$  classes weight by  $Y=l$

Naive Bayes = same prediction method  
naive assumption about D-dim

D-dim  
 $z = \text{test point} = (z^1 z^2 \dots z^D)$   
prediction  $P(Y=c|Z) = \frac{P(Z|Y) \cdot P(Y)}{\text{normalized } P(Z)}$   $\rightarrow$  prior

$\Rightarrow$  need  $P(Z|Y) = P(z^1, z^2, \dots, z^D | Y)$

before  $P(z^1 z^2 \dots z^D | Y) \stackrel{\text{FIT}}{\sim} \mathcal{N}(z^{\text{Ddim}} | \mu, \Sigma)$   
D-dim curve

Now: assume feat 1, 2, ..., D relatively (??) independent

NAIVE  
 $P(z^1 z^2 \dots z^D | Y) = P(z^1 | Y) P(z^2 | Y) \dots P(z^D | Y)$

ex:  $z = \text{person}$   $z^1 = \text{height}$   $z^2 = \text{weight}$   $z^3 = \text{shoe size}$

Independent? NO assumption wrong!

but classifier still works  $\Rightarrow$  produces output.

Fit  $P(z|y=l) \Rightarrow$  Fit separately every feature (1-dim)  
 class  $y=l$  D-dim

TRAINING: (like before, separate for each class)

class  $y=l$   $\Rightarrow$  FILTER X data =  $X(y=l)$

• for each feature  $j=1:D \rightarrow$  separately

Fit 1-dim  $P(x^j) \approx$

vector  
 $x^j_{1:M}$

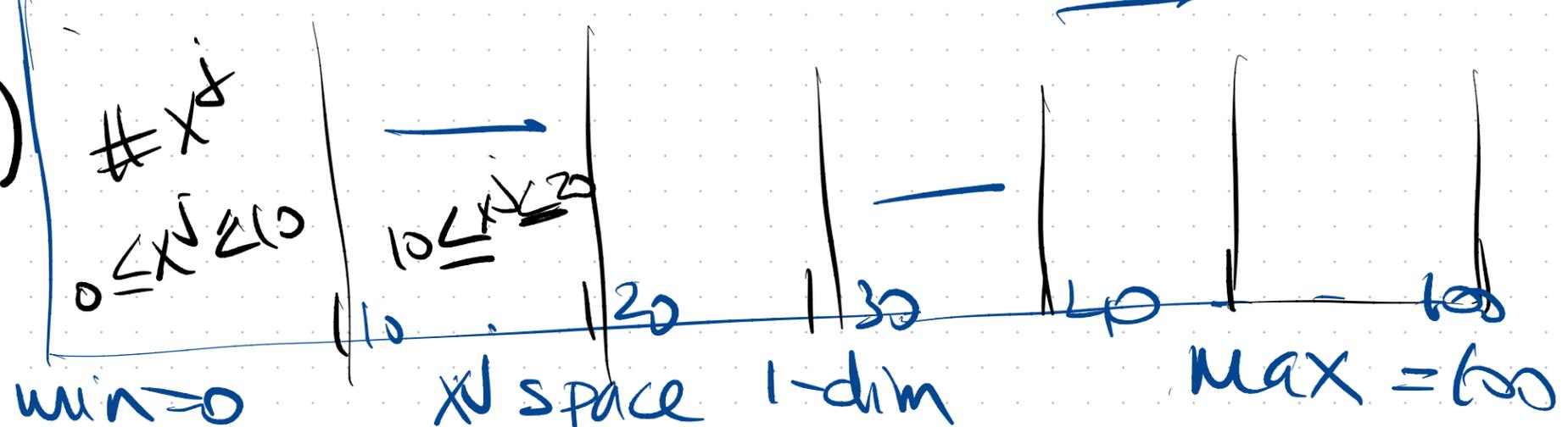
	$f_1$	$f_2$	$f_j$	$f_D$	$y$
$x_1$			$x_{1j}$		$l$
$x_2$			$x_{2j}$		$l$
$x_M$			$x_{Mj}$		$l$

fit param curve (1-dim)  
 $P(x^j) \approx \mathcal{N}(\mu, \sigma^2)$

options

no assumptions

fit histogram  
 • buckets ( $B=10$ )  
 • count per bucket  
 normalized



ex:  $x^i = \text{weight}$

buckets	$B_1: 0-40LB$	$B_2: 40LB-60LB$	$B_3: 60-90$	...	$B_{\text{last}} > 250LB$
	count <sub>1</sub> / num	count <sub>2</sub> / num	count <sub>3</sub> / num		count <sub>last</sub> / num

More granular  $\Rightarrow$  more buckets  
requires more data for accuracy

product: for  $Z = \text{test point} = (z^1, z^2, \dots, z^D)$

$$P(Y|Z) = \frac{P(Z|Y) P(Y)}{P(Z)} = \frac{P(z^1|Y) P(z^2|Y) \dots P(z^D|Y) \cdot P(Y)}{P(Z)}$$

*trained* (red arrow pointing to  $P(z^D|Y)$ )

*num product very small* (blue arrow pointing to the denominator  $P(Z)$ )

*use logs*

$$\log(P(Y|Z)) = \log(P(Y)) + \sum_{d=1}^D \log(P(z^d|Y)) - \log(P(Z))$$

# Naive Bayes for Text using Bag-of-Words representation

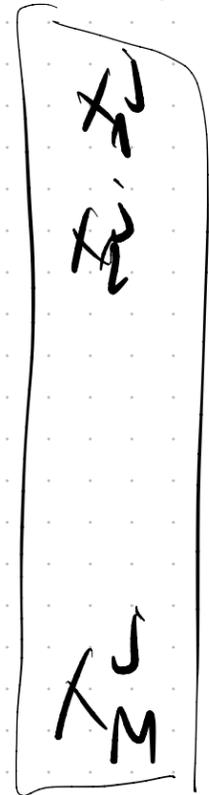
- class  $y = e$  as filter.  $X = X[y=e]$  data with label  $y=e$
- want fit  $P(x) = P(x|y=e)$

$$P(x) = P(x^1) \cdot P(x^2) \cdot \dots \cdot P(x^D)$$

- fit each separately

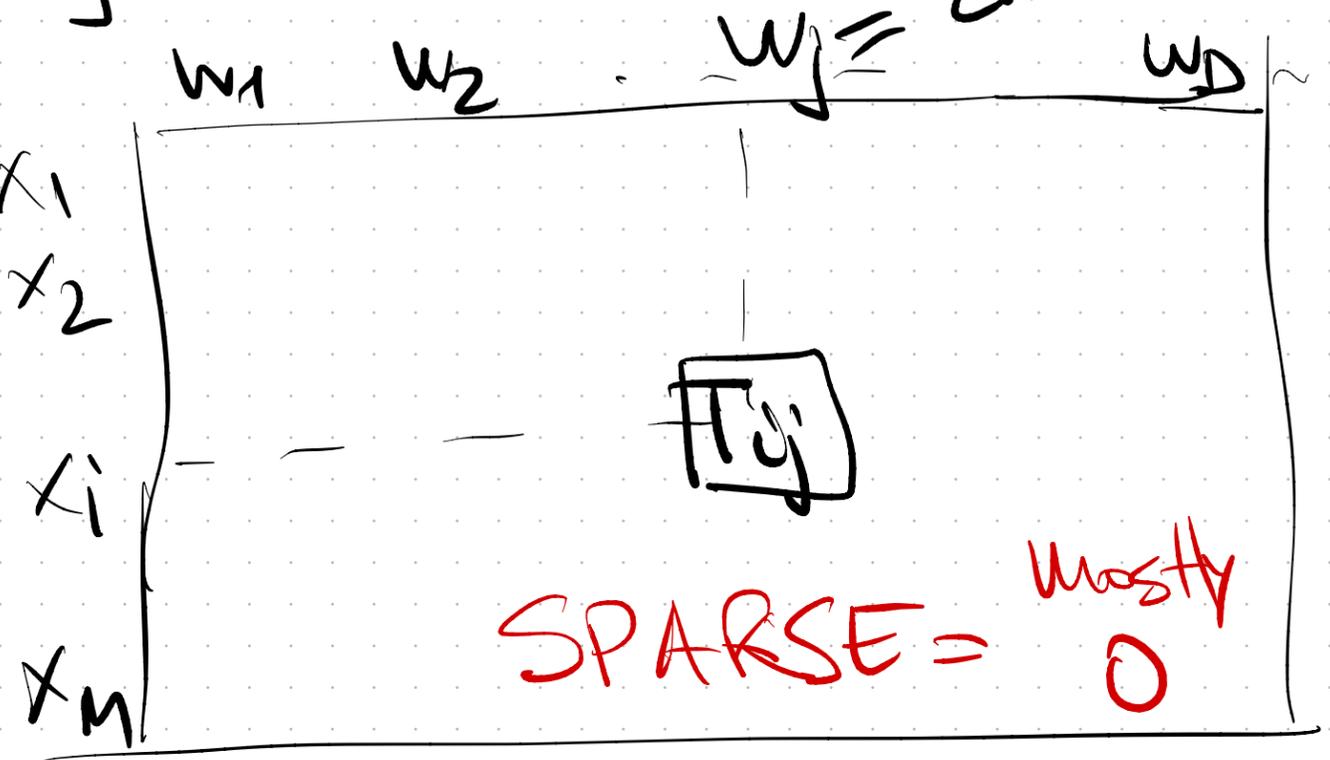
For feature word  $j = \text{"china"}$

counts  $\Rightarrow$  distribution (1-dim)



$$\text{OBVIOUS} = P(x_i^j) = \frac{\text{TF}_{ij}}{\text{DOCLength}}$$

docs  $y=e$



$T_{ij} = \text{TF}_{ij} = \frac{\text{term freq (word } j \text{ in doc } i)}{\text{\# occ of word } j}$

useless words: "for", "a", "about", "the", "in"  
 "STOPWORDS"  
 ↓ exclude

VISUALIZE BR ADOC  $x_i$

$$x_i = \frac{w_1 \quad w_2 \quad \dots \quad w_i \quad \dots \quad w_D}{x_i^1 \quad x_i^2 \quad \dots \quad x_i^i \quad \dots \quad x_i^D} \Rightarrow \left[ \frac{x_i^1}{|x_i|} \quad \frac{x_i^2}{|x_i|} \quad \dots \quad \frac{x_i^D}{|x_i|} \right]$$

term freq in doc  $x_i$

doclength( $x_i$ ) =  $|x_i|$  = # words in doc

discrete (counts)  
very sparse (many 0)

• dont like prob = 0! Fix = "smoothing"  
**LAPLACE** - add +1 numerators  
add +D denominator

$$x_i = \frac{w_1 \quad w_2 \quad \dots \quad w_i \quad \dots \quad w_D}{x_i^1 \quad x_i^2 \quad \dots \quad x_i^i \quad \dots \quad x_i^D} \Rightarrow \left[ \frac{x_i^1 + 1}{|x_i| + D} \quad \frac{x_i^2 + 1}{|x_i| + D} \quad \dots \quad \frac{x_i^D + 1}{|x_i| + D} \right]$$

**LAPLACE** • do not have prob = 0 (min prob =  $\frac{1}{D}$ )  
• prior for words in doc if they do not actually occur  $\frac{1}{|x_i| + D}$   
• deal with words  $w_{new}$  new in TEST SET (never seen in TRAIN)

• math nice, including for logs.

• general:  $\frac{+ \epsilon \text{ to numerator}}{+ D \epsilon \text{ to denominator}} \Rightarrow$

$\epsilon$  = strength of prior

$\epsilon = 1$  default in text

• still a probab (sum = 1)

---

prod  $\Rightarrow 0 \Rightarrow$  BAD

•  $P(x) = P(x^1) \cdot P(x^2) \cdot \dots$

$P(x^i)$

$\Rightarrow 0$   
 $\rightarrow$  makes all other features irrelevant.

• never seen  $x^i$  word in training

$0/0$  in overall formula for each class  $Y_c$ .