

# Lecture 6/6 Generative Models

HW3: GDA, Naive Bayes etc

- Fit curve to data

- Fit density to data

Want prob. density over data space

$pr(x_i | \text{param})$  = prob that  $x_i$  datapoint occurs

$$pr_{\theta}(x_i) = \frac{\# \text{ datapoints } x_i \text{ or very similar}}{\# \text{ total} = N}$$

curve = parametric with mode P for density

ex: curve = Normal = Gaussian

Next Wed 6/11

Demos for HW2B NNets

Project Proposal June 24

Project page

Smaller Project  $\approx$  Hw

(20h)

• choices

$x_i$  = patient

$pr(x_i)$  = proportions of patients like this

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

Fit curve = Gaussian to data  $X$   
 $(\mu, \sigma^2)$



Find the optimal params  $\mu, \sigma^2$   
to maximize Likelihood (data  $X$ )

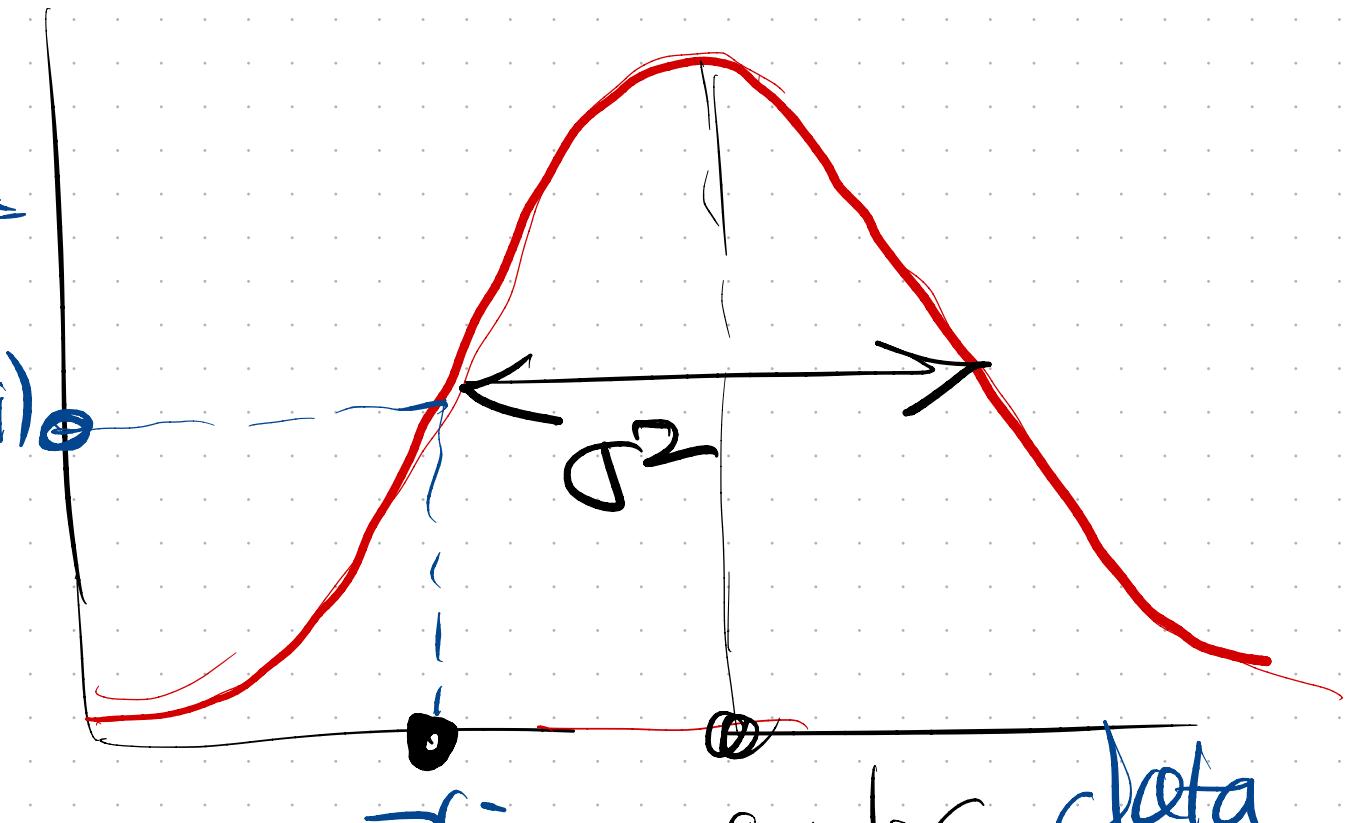
product because  
indep. datap.

$$\prod_{i=1}^N \text{pr}(x_i)$$

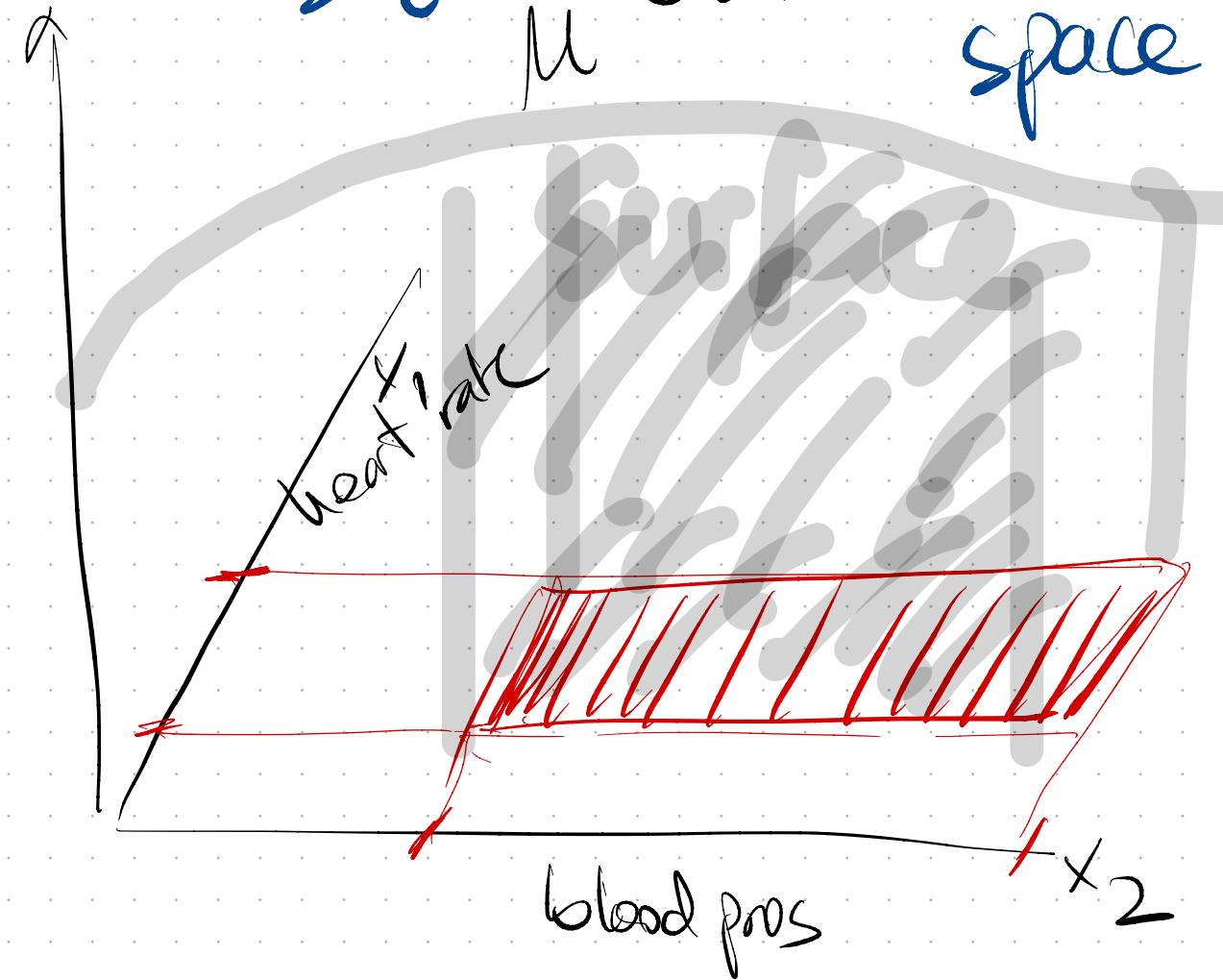
Generative formulation:

Find curve = Gaussian  $(\mu, \Sigma)$

curve = generator of data  $\Rightarrow$  data  $X$   
is most likely.



$\mu$  center data space



1-dim Gaussian:  $W(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$

$\sigma^2 = \text{VAR}$   $\rightarrow$  data center most likely  $x$   
 $\text{std} = \sqrt{\text{VAR}}$

$x_i = D\text{-dim} \Rightarrow W(x_{D\text{-dim}} | \mu_{D\text{-dim}}, \Sigma)$

$\mu$   $D\text{-dim}$  vector  
 $\Sigma$  sigma (capital)  $D \times D$  COVAR()

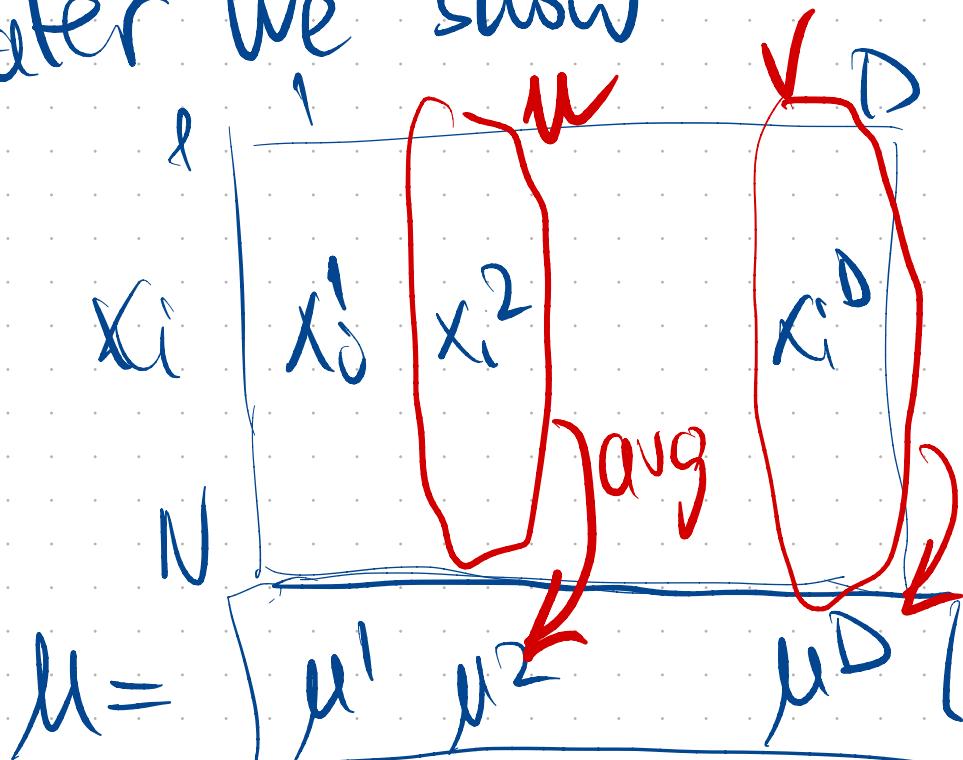
$\Sigma$   $D \times D$   
 $\Sigma = \text{COVAR. Matr}$

$\sum_{i=1}^n = \text{summation}$

$p(x_i | \mu, \Sigma) = \frac{1}{(2\pi)^{-D/2} |\Sigma|^{-1/2}} \exp \left[ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right]$

determinant ( $\Sigma$ )  $\exp$  function  
 $\Sigma$   $D \times D$   $\Sigma^{-1}$   $D \times D$  Inverse of  $\Sigma$

Later we show



if  $\mu = 0$   
 $X^T X = \sum_{i=1}^D \text{COVAR}$

[arg of each  
 $w_l = \text{feature}$

$\langle \text{col}^1 \cdot \text{col}^1 \rangle$

$\langle \text{col}^D \cdot \text{col}^1 \rangle$

$- \langle \text{col}^D \cdot \text{col}^D \rangle$

$\langle \text{col}^1 \cdot \text{col}^D \rangle$

④ Fit D-Dim Gaussian to data  $X$  (D-features)

$\Rightarrow \mu = \text{mean}(x)$  D-dim vector  $(\mu^1, \mu^2, \dots, \mu^D)$  avg for each feat

$$\Rightarrow \Sigma = \text{covar}(x) = (x - \mu)^T (x - \mu)$$

$(\mu, \Sigma)$  param that makes generator-Gaussian best for data  $X$

When used to generate  
data  $\Rightarrow$  max-chance to  
obtain  $X_{1:N}$

Data likelihood (as generated by Gaussian with param  $\mu, \Sigma$ )

all data

$i=1:N$

$$L(x) = \prod_{i=1}^N p(x_i) = \prod_{i=1}^N N(x_i | \mu, \Sigma)$$

$$\log \text{likelihood} = LL(x) = \log \left[ \prod_{i=1}^N N(x_i | \mu, \Sigma) \right]$$

$$= \sum_{i=1}^N \left[ \log(2\pi)^{-d/2} + \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu) \Sigma^{-1} (x_i - \mu)^T \right]$$

$$\text{want } \max LL(x) \rightarrow \frac{\partial LL(x)}{\partial \mu}$$

$$\frac{\partial LL(x)}{\partial \Sigma}$$

$\mu, \Sigma$   
are indep.  
(no constraint)

$$\frac{\partial \log l}{\partial \mu} = \sum_{i=1}^N \cancel{2 \frac{1}{2}} \cancel{\sum_{i=1}^{i+1}} (x_i - \mu) \stackrel{\text{want}}{=} 0 \Rightarrow \sum_{i=1}^N (x_i - \mu) = 0$$

$\Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x_i$

$$\frac{\partial \log l}{\partial \Sigma} = \sum_{i=1}^N \left[ \cancel{\frac{1}{2}} \sum_{j=1}^N - \frac{1}{2} (x_i - \mu)^T (x_i - \mu) \right] \stackrel{\text{want}}{=} 0$$

$$N \cdot \sum = \sum_{i=1}^N (x_i - \mu)^T (x_i - \mu)$$

$$1 \quad \sum_{i=1}^N (x_i - \mu)^T (x_i - \mu) = \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^T (x_i - \mu)$$

$N$        $D \times 1$        $1 \times D$

**COVAR**

**Sum of N**  $\begin{bmatrix} DXD \\ \text{mat} \end{bmatrix}$

$- D \sum_{i=1}^N (x_i^D - \mu^D)(x_i^D - \mu^D)^T$

$\sum_{i=1}^N (x_i^D - \mu^D)^2$

$$\mu = 0 \Leftrightarrow \text{X centered} = \text{X} - \mu$$

$$\text{COVAR } \Sigma = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

Classification with gen. models (Gaussian fit + data)

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

label  $\rightarrow$  datapoint

$P(X|Y)$   $\rightarrow$  prior for each  $y = \text{class } Y_1, Y_2, \dots, Y_L$   
product over  $Y_1, \dots, Y_L$

$P(Y)$   $\rightarrow$  normalization of numerators  
For  $L$  classes

density of  $X$  in certain class  $y = Y_e$

Fit curve to  $Y_e$  points  $\rightarrow$  priors

$y \in \{0, 1\}$

$x = \text{test point}$

$P(Y=1|X) = \frac{P(X|Y=1) \cdot P(Y=1)}{P(X)}$

$P(Y=0|X) = \frac{P(X|Y=0) \cdot P(Y=0)}{P(X)}$  ignore

GDA = Gauss Discriminant Analysis

- $\mathcal{N}_1(\underline{\mu}_1, \Sigma_1)$  fit on  $X$  data with  $y=+1$  subset

- separate  
 $\mathcal{N}_0(\underline{\mu}_0, \Sigma_0)$  fit on  $X$  data subset  $y=0$

- predict: compute  $P(Y=+1|z) = \mathcal{N}_1(z|\underline{\mu}_1, \Sigma_1) \cdot P(Y=+1) / P(z)$   
test point  $z$
- predict  $y = \text{argmax}$

- can do this for  $y=1, y=2, \dots, y=L$  ( $L$  classes)

- can use common  $\Sigma$  for all classes.