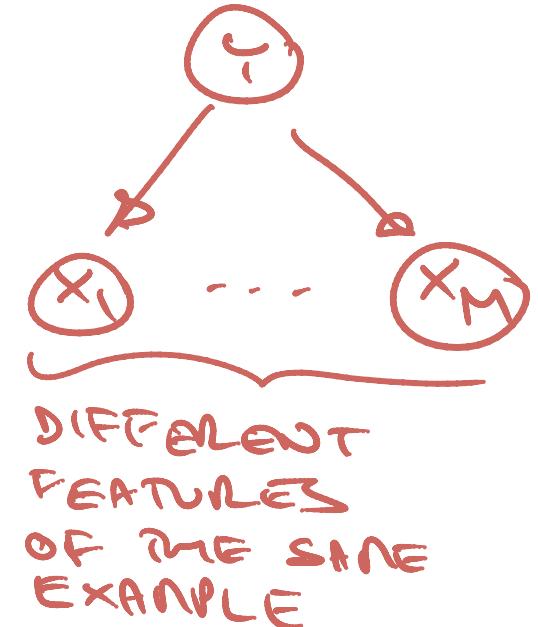


# Review

- Train-test split
- Cross-validation
- Regularization and model complexity
  - ▶  $L_1, L_2$

# Review

- Classification w/ Naïve Bayes
- Features assumed independent given class
- Prediction:  $y = \underset{y}{\operatorname{arg\,max}} \left[ p(y) \prod_j p(x_j|y) \right]$
- Variations: Bag of Words, Gaussian NB



$$P(x_j | Y=1) = \theta_{1j}^{x_j} (1-\theta_{1j})^{1-x_j}$$

$$P(x_j | Y=0) = \theta_{0j}^{x_j} (1-\theta_{0j})^{1-x_j}$$

# A closer look at NB

- $Y = 1$  if:  $P(Y=1) \prod P(x_j | Y=1) \geq P(Y=0) \prod P(x_j | Y=0)$

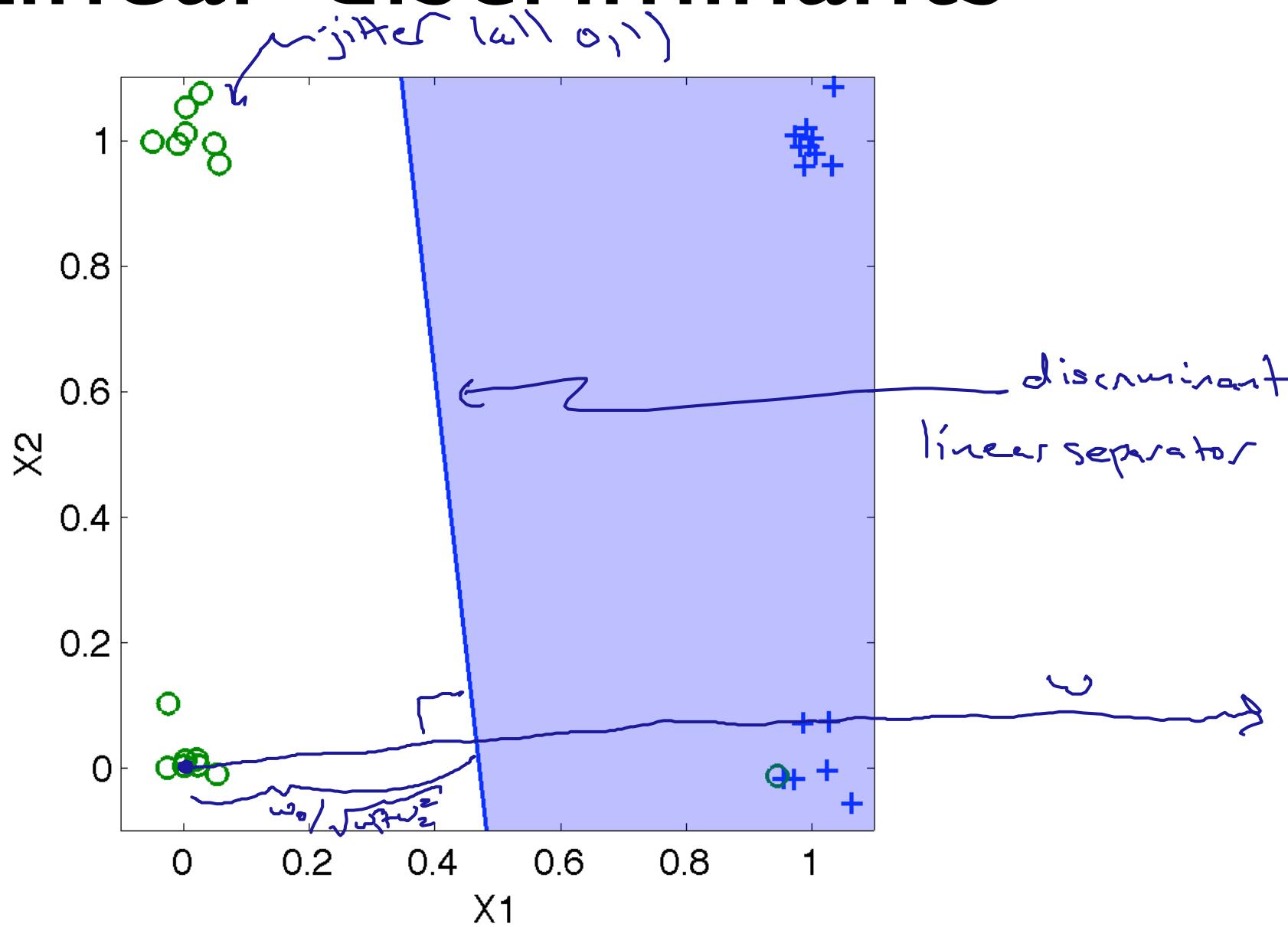
$$\log(P(Y=1)) - \log(P(Y=0)) + \sum_j [\log P(x_j | Y=1) - \log P(x_j | Y=0)] \geq 0$$

$$w_0 + \sum_j w_j x_j \geq 0$$

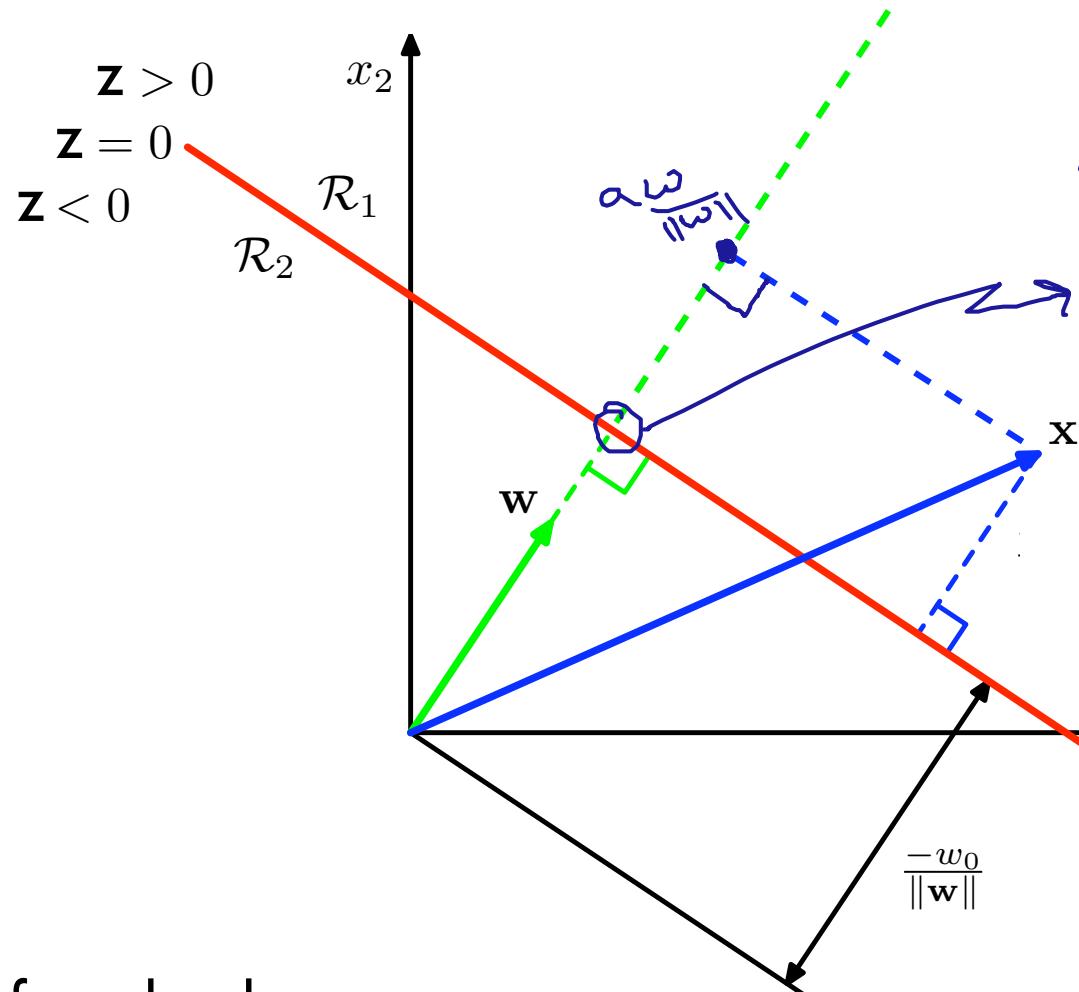
$$w_0 = \log P(Y=1) - \log P(Y=0) + \sum_j \log \frac{P(x_j | Y=1)}{P(x_j | Y=0)}$$

$$w_j = \log(\theta_{1j}) - \log(1-\theta_{1j}) - b_j(\theta_{1j}) + \log(1-\theta_{0j}) - \sum_j \log \frac{P(x_j | Y=1)}{P(x_j | Y=0)}$$

# Linear discriminants



# Geometry of a discriminant



$$\begin{aligned} \alpha \frac{w}{\|w\|} \cdot w &= x \cdot w \\ \alpha \frac{\|w\|^2}{\|w\|} &= x \cdot w \\ \alpha = x \cdot w &/ \|w\| \end{aligned}$$

$$\begin{aligned} w_0 + w^T x &= 0 \\ w_0 + w_1 x_1 + w_2 x_2 &= 0 \\ w_0 + w^T b \frac{w}{\|w\|} &= 0 \\ w_0 + b \|w\| &= 0 \\ -\frac{w_0}{\|w\|} - b & \end{aligned}$$

figure adapted from book

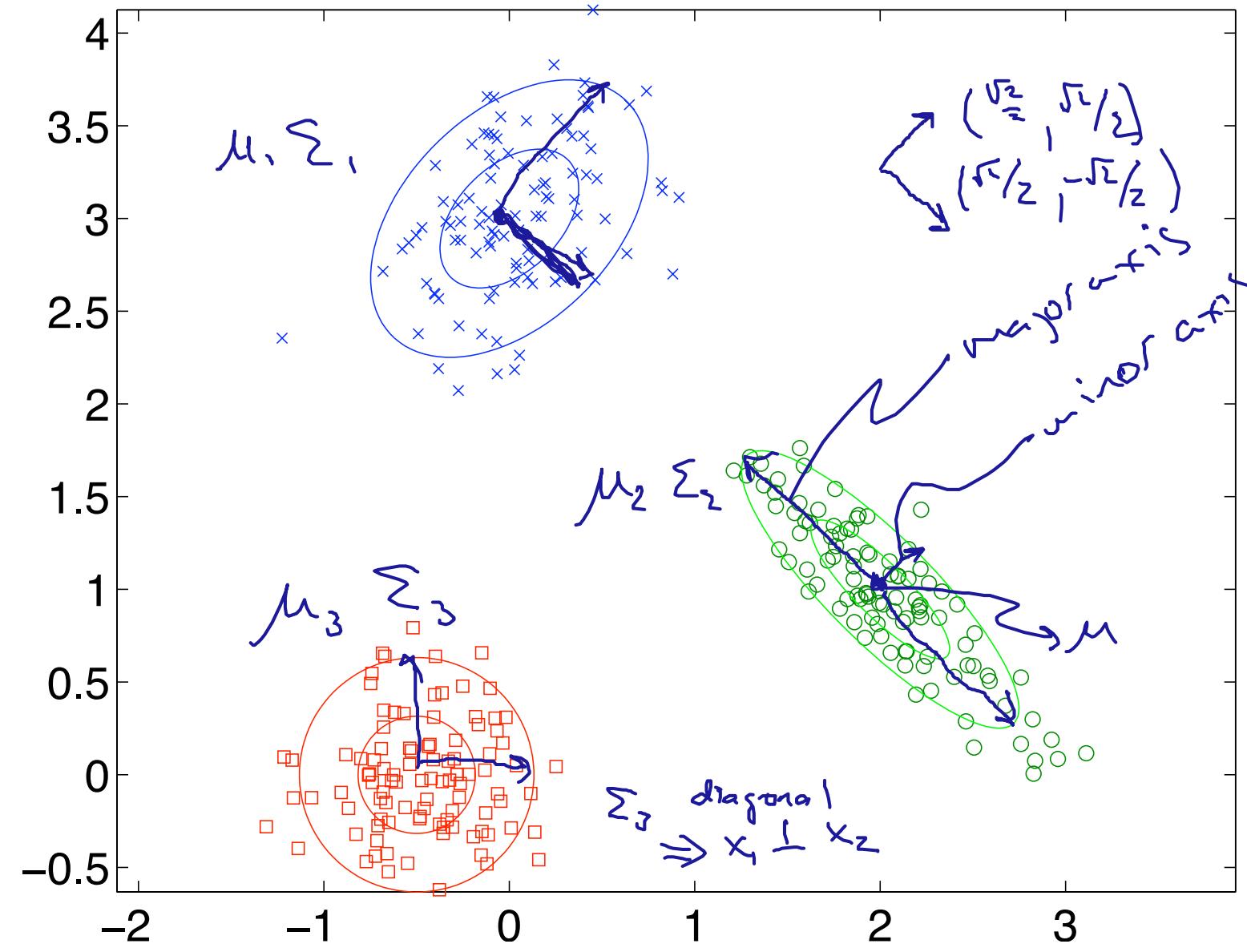
# Continuous vars

- For categorical  $X$ , NB gave us a linear discriminant
- What about continuous  $X$ ?
  - ▶ e.g., Gaussian NB
- Will turn out the same, but we'll work it out for a generalization

$$\tilde{x} \sim N_{\text{mean}}(x | \mu, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left\{-0.5(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

# Multivariate Gaussians

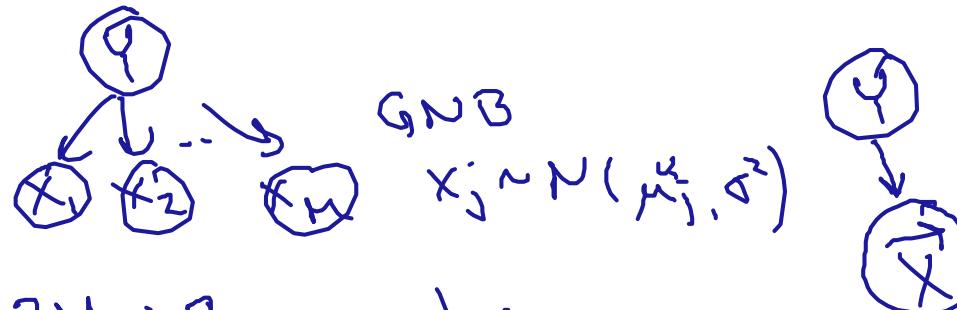
# Multivariate Gaussians



$$\begin{aligned}\mu_i &= (0, 3) \downarrow \\ \Sigma_i &= U D_i U^T \\ U &= \begin{pmatrix} \sqrt{2}/2 & \sqrt{2}/2 \\ \sqrt{2}/2 & -\sqrt{2}/2 \end{pmatrix} \\ D_i &= \begin{pmatrix} 2 & 0 \\ 0 & 0.08 \end{pmatrix} \\ \sqrt{D_i} &= \begin{pmatrix} 1.45 & 0 \\ 0 & 0.28 \end{pmatrix} \\ \Sigma_2 &= U D_2 U^T \\ D_2 &= \begin{pmatrix} 0.02 & 0 \\ 0 & 0.24 \end{pmatrix} \\ \sqrt{D_2} &= \begin{pmatrix} -1.4 & 0 \\ 0 & 0.51 \end{pmatrix} \\ \Sigma_3 &= \begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix} = D_3 \\ \Sigma_3 &= U D_3 U^T\end{aligned}$$

# A generalization of Gaussian Naïve Bayes

$$\Sigma_{ij} = \text{cov}(x_i, x_j)$$



$2M + 2$  parameters

( $y$ , binary) (shared var)

or  $4M + 1$  if no shared variances among classes, vars  
 $x_j \sim N(\mu_j, \sigma_j^2)$

$$\vec{x} \sim N(\mu_x, \Sigma_x)$$

$\mu_x$  or  $\mu_i$   
 $(M \times 1)$

$$2M + \binom{M}{2}^{+M}$$

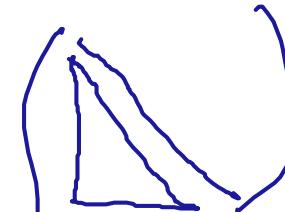
params

$$2M + 2\binom{M}{2}^{+2M}$$

if

no shared var  
 among classes

make  $\Sigma$  diagonal



$$\Sigma_0, \Sigma_1$$

# Generalizing GNB

- $P(X | Y) = N(X | \mu_Y, \Sigma)$
- ▶ if  $\Sigma = \sigma^2 I \Rightarrow \text{GNB}$
- Pick  $Y=1$  if

$$\log P(Y=1) - \log P(Y=0) + \log P(X | \mu_1, \Sigma) - \log P(X | \mu_0, \Sigma) \geq 0$$

$$\begin{aligned} & \log P(X | \mu, \Sigma) \\ & - \log [2\pi|\Sigma| - \frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)] \end{aligned}$$

$$\omega_0 + \omega^T x \geq 0$$

$$\begin{aligned} & -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1) + \frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0) \\ & \frac{1}{2} \left[ -x^T \cancel{\Sigma^{-1}} x + 2x^T \cancel{\Sigma^{-1}} \mu_1 - \mu_1^T \cancel{\Sigma^{-1}} \mu_1 \right. \\ & \left. + x^T \cancel{\Sigma^{-1}} x - 2x^T \cancel{\Sigma^{-1}} \mu_0 + \mu_0^T \cancel{\Sigma^{-1}} \mu_0 \right] \end{aligned}$$

# Fisher linear discriminant

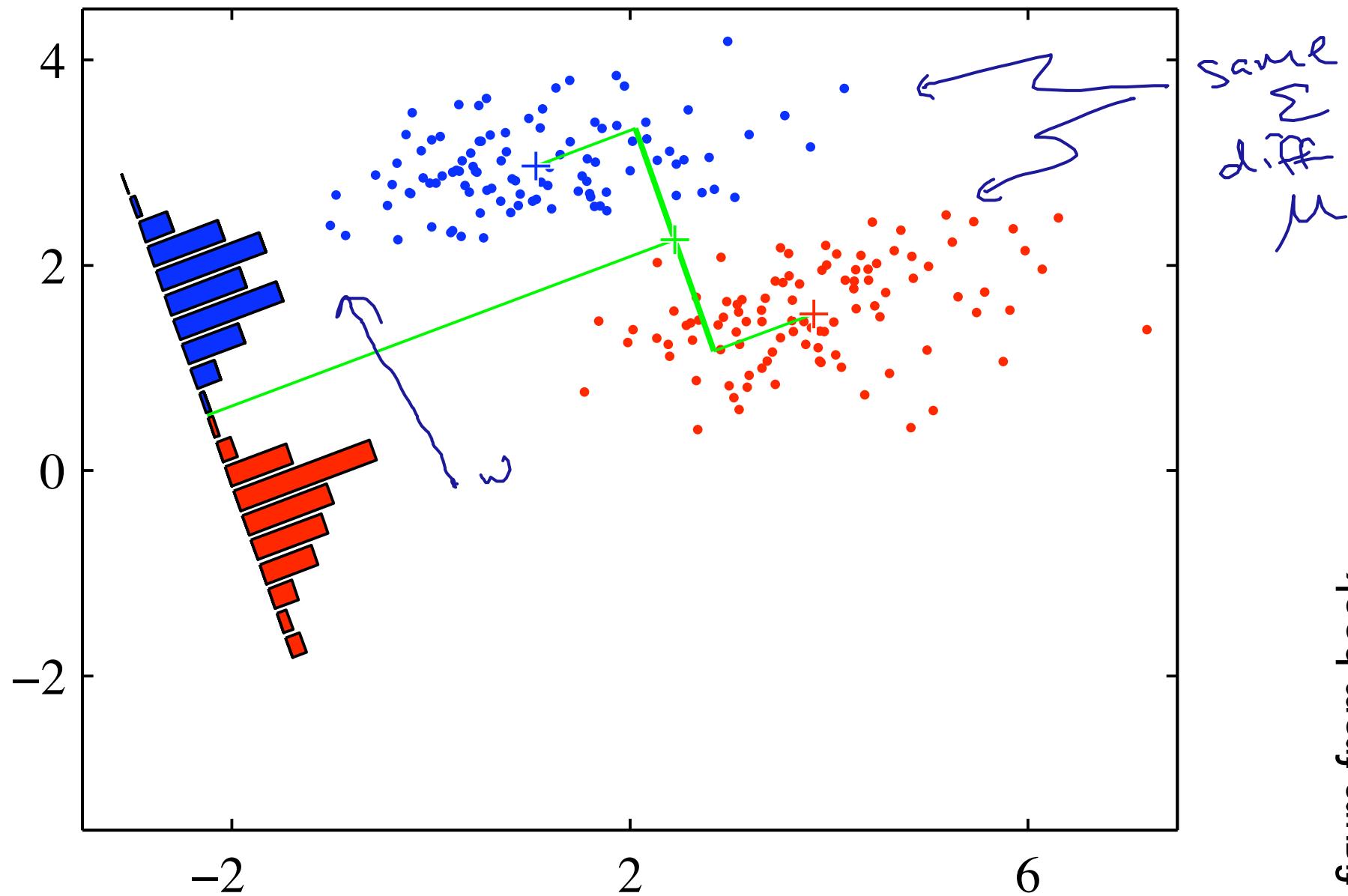


figure from book