

UCB - CS189
Introduction to Machine Learning
Fall 2015

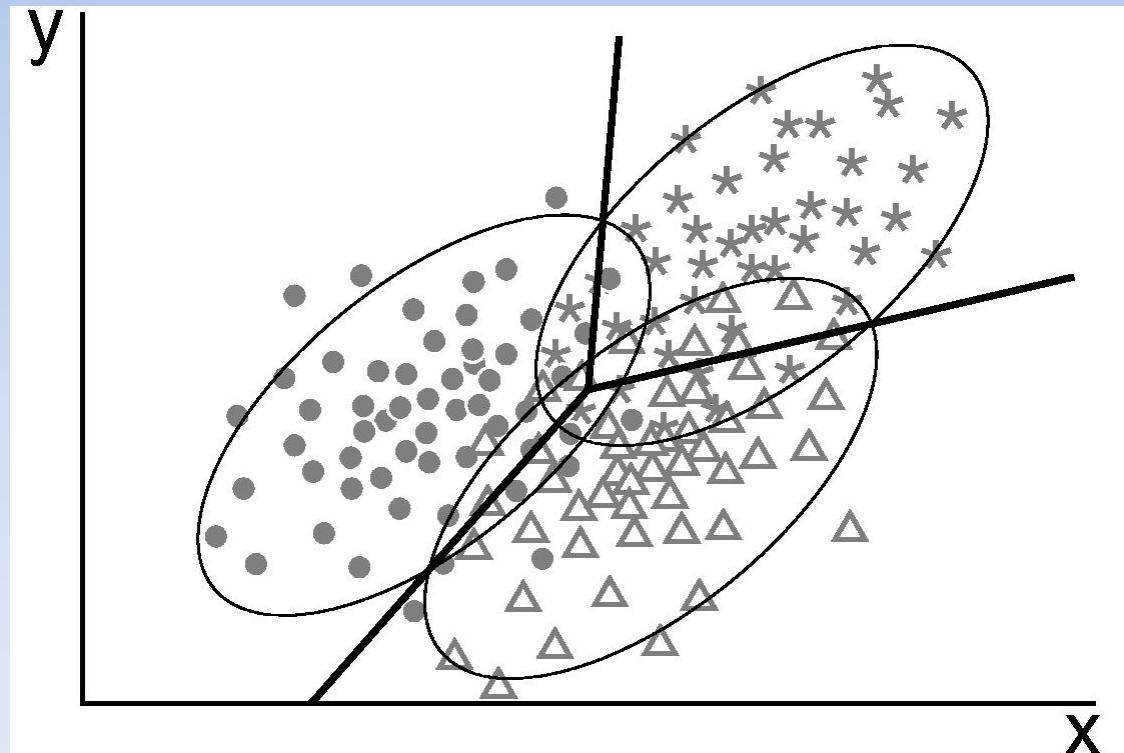
Lecture 14: Mixture models

Isabelle Guyon
ChaLearn

Come to my office hours...

Wed 1:30-3:30 Soda 329

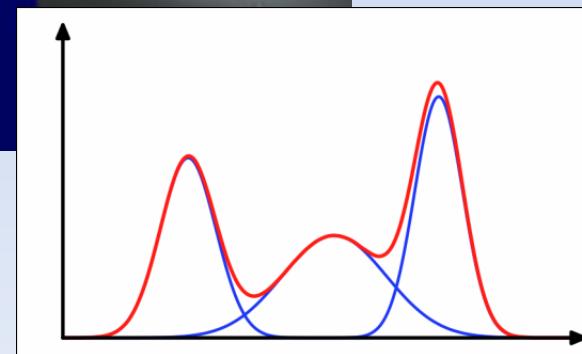
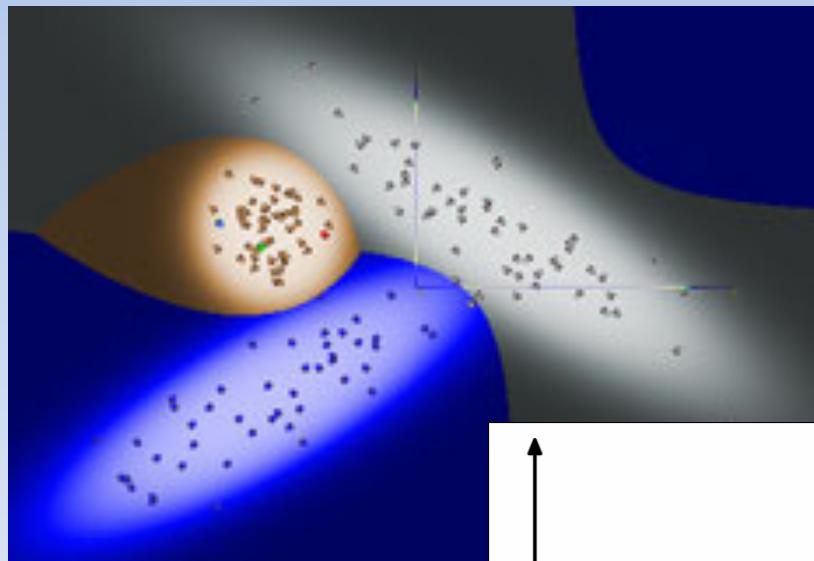
Last time: LDA



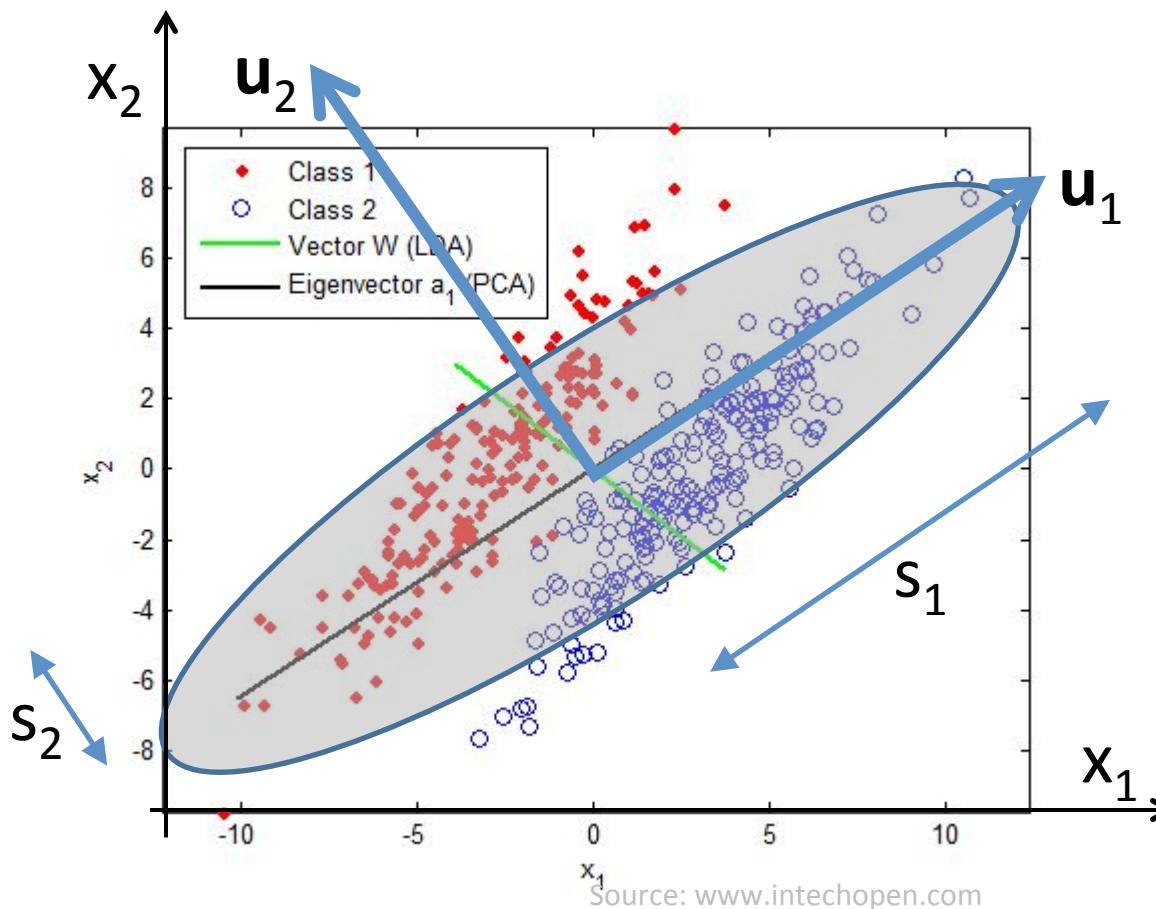
Come to my office hours...

Wed 1:30-3:30 Soda 329

Today: Mixture models



Covariance matrix $\Sigma = X^T X$



Source: www.intechopen.com

$$\Sigma = X^T X \text{ covariance matrix (centered data)} \quad \sigma_{ij} = (1/N) \sum_{k=1:N} (x_i^k - \mu_i)(x_j^k - \mu_j)$$

$$\Sigma = US^2U^T$$

$U = [u_i \text{ eigen vectors}], S = [s_i \text{ diagonal singular values}]$

$$U^T U = I$$

N \uparrow
 \downarrow d X^T N samples

$$X = [x_i^k]$$

data matrix

$$x^k$$

sample dim(1, d)

feature dim(N, 1)

$$x_i$$

feature dim(N, 1)

$$x_j$$

d features

$$\Sigma = X^T X$$
$$\text{dim}(d, d)$$

Σ^{-1} for BIG data



$X^T X$ dim(d,d)
 XX^T dim(N,N)
 which one do I
 rather invert?

Singular value decomposition:

$X = VSU^T$, with $U^T U = 1$ and $V^T V = 1$

S diagonal dim(r, r): singular values, $r = \text{rank}(X) \leq \min(d, N)$

$X^T X = US^2 U^T$ and $XX^T = VS^2 V^T$ $\text{dim}(U) = (d, r)$ $\text{dim}(V) = (N, r)$

cheap $\Xi = XU$, $\Sigma^{1/2} = USU^T$, $\Sigma^{-1} = US^{-2}U^T$, etc.

r is small + keep only largest singular values.

Kernel trick:

$d \gg N \rightarrow XX^T = VS^2 V^T$ $XX^T \rightarrow K$ (kernel trick)

Instead of $\Xi = XU$ use $\Xi = VS$, much cheaper!

$\Xi_{\text{new}} = X_{\text{new}} U?$

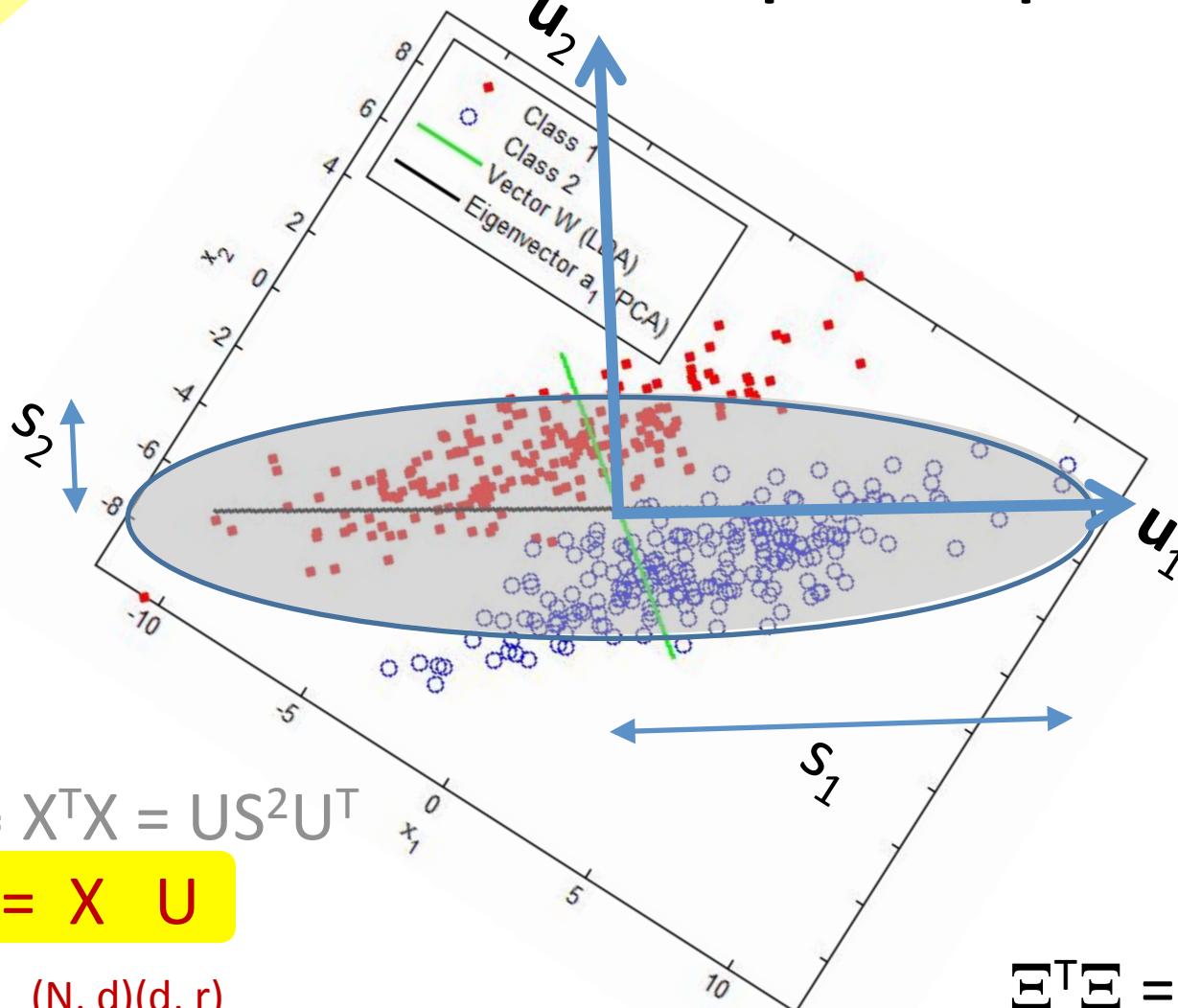
$\Xi_{\text{new}} = X_{\text{new}} X^T VS^{-1}$

$X = VSU^T$ so $U = X^T VS^{-1}$

$X_{\text{new}} X^T \rightarrow [k(x^h, x^k)]$ (kernel trick)

Review

Rotation in principal axes



$U = [u_i \text{ eigen vectors}]$

PCA consists in keeping only the axes with largest eigen values ⁷

Review

Benefits of rotation in principal axes



$$\Sigma = X^T X = U S^2 U^T$$

$$\Sigma = X U$$

dim(N, r)

$$\Sigma^T \Sigma = S^2$$

$U = [u_i]$ eigen vectors]

Benefits of rotation in principal axes



$$\Sigma = X^T X = U S^2 U^T$$

$$\Xi = X U$$

dim(N, r)

$$\Xi^T \Xi = S^2$$

$U = [u_i]$ eigen vectors]

- In the new Ξ -space the covariance matrix $\Xi^T \Xi$ is **diagonal**.
- **Inverting a diagonal matrix (or taking its square root) is trivial.**

Benefits of rotation in principal axes



$$\Sigma = X^T X = U S^2 U^T$$

$$\Xi = X U$$

dim(N, r)

$$\Xi^T \Xi = S^2$$

$U = [u_i]$ eigen vectors]

- In the new Ξ -space the covariance matrix $\Xi^T \Xi$ is **diagonal**.
- **Inverting a diagonal matrix (or taking its square root)** is trivial.

Application to ridge regression:

$$f(x) = x w^T$$

$$w^T = (X^T X + \lambda I)^{-1} X^T y \quad \lambda > 0$$

(d,1) (d,N)(N,d) (d,d) (d,N)(N,1)

In the rotated space:

$$f(x) = \xi \omega^T \quad \xi = x U$$

$$\omega^T = (S^2 + \lambda I)^{-1} \Xi^T y$$

$$w^T = U (S^2 + \lambda I)^{-1} U^T X^T y$$

$$S^2 + \lambda I = \begin{pmatrix} S_1^2 + \lambda & & & \\ & S_2^2 + \lambda & & \\ & & \ddots & \\ & & & S_d^2 + \lambda \end{pmatrix}$$

How to transform new samples

$$\Sigma = X^T X = U S^2 U^T$$

$$\begin{aligned}\Xi &= X \quad U & \text{dim}(N, r) \\ (N, r) \quad (N, d) &(d, r) \\ U &= [\mathbf{u}_i \text{ eigen vectors}]\end{aligned}$$

$$\Xi^T \Xi = S^2$$

Now we get a new sample \mathbf{x}_{new} not in the training set.

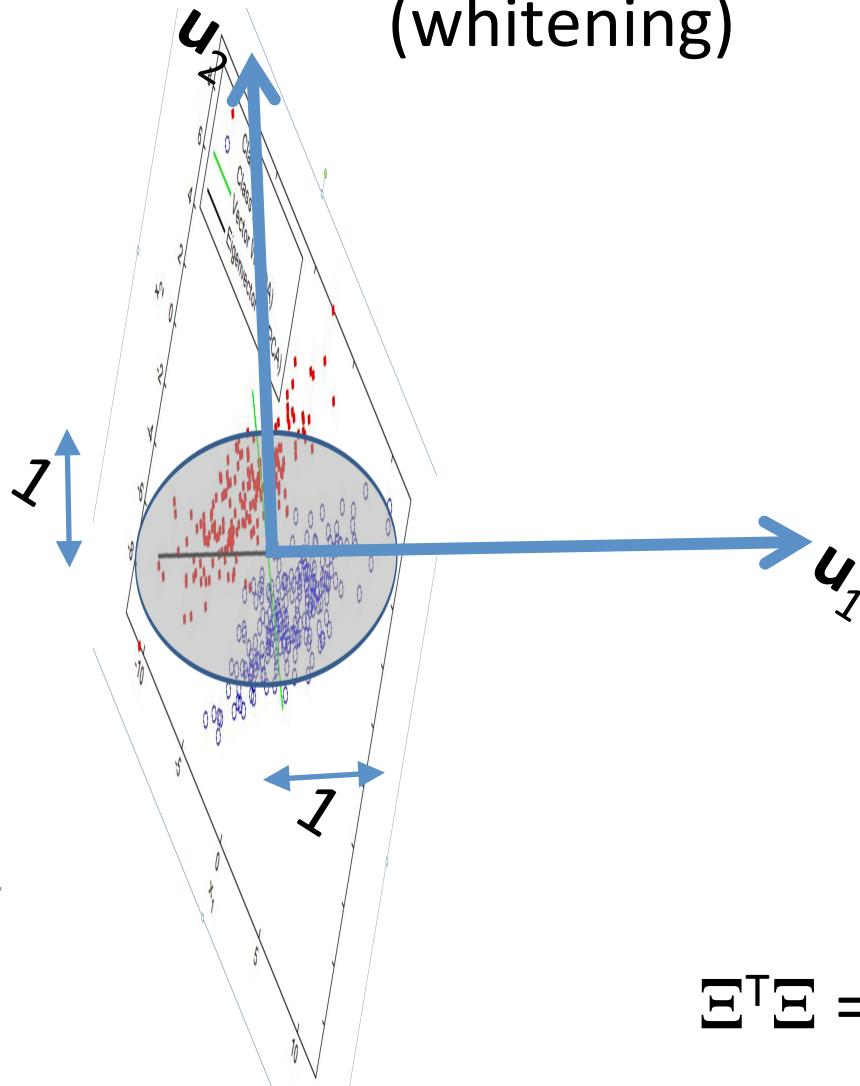
$$\phi_{\text{new}} = \mathbf{x}_{\text{new}} U$$

Likewise on a test matrix X' of $\text{dim}(N', d)$

$$\begin{aligned}\Xi' &= X' \quad U & \text{dim}(N', r) \\ (N', r) \quad (N', d) &(d, r)\end{aligned}$$

Do NOT compute U on all the data $[X ; X']$ and do $[X ; X']U$.

Rotation and scaling (whitening)



$$\Sigma = X^T X = U S^2 U^T$$

$$\Xi = X \Sigma^{-1/2}$$

$(N, r) \quad (N, d)(d, r)$

$$\Sigma^{1/2} = U S U^T$$

$$\Xi^T \Xi = I$$

Benefits of whitening



$$\Sigma = X^T X = U S^2 U^T$$

$$\Xi = X \Sigma^{-1/2}$$

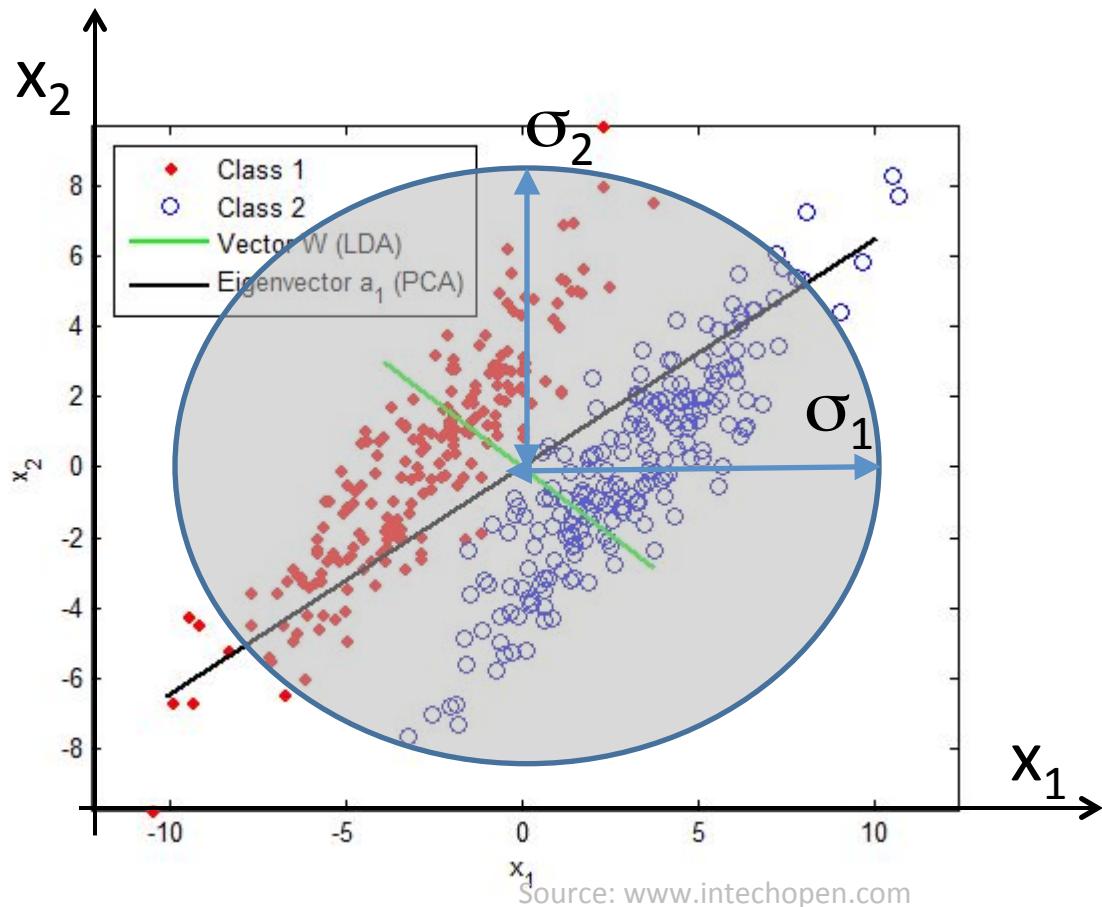
dim(N, r)

$$\Xi^T \Xi = I$$

$$\Sigma^{1/2} = U S U^T$$

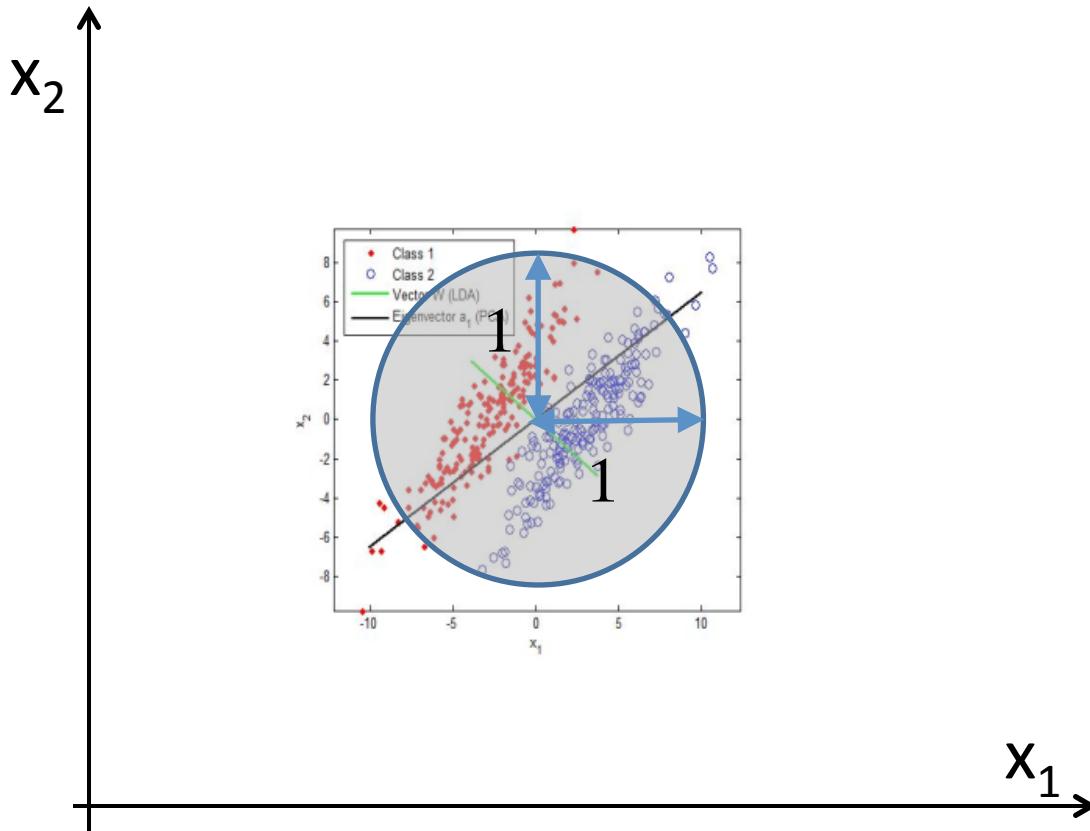
Same as rotating in principal axes +
The weights of any linear method in Ξ -space
are on the same scale.

Spherering



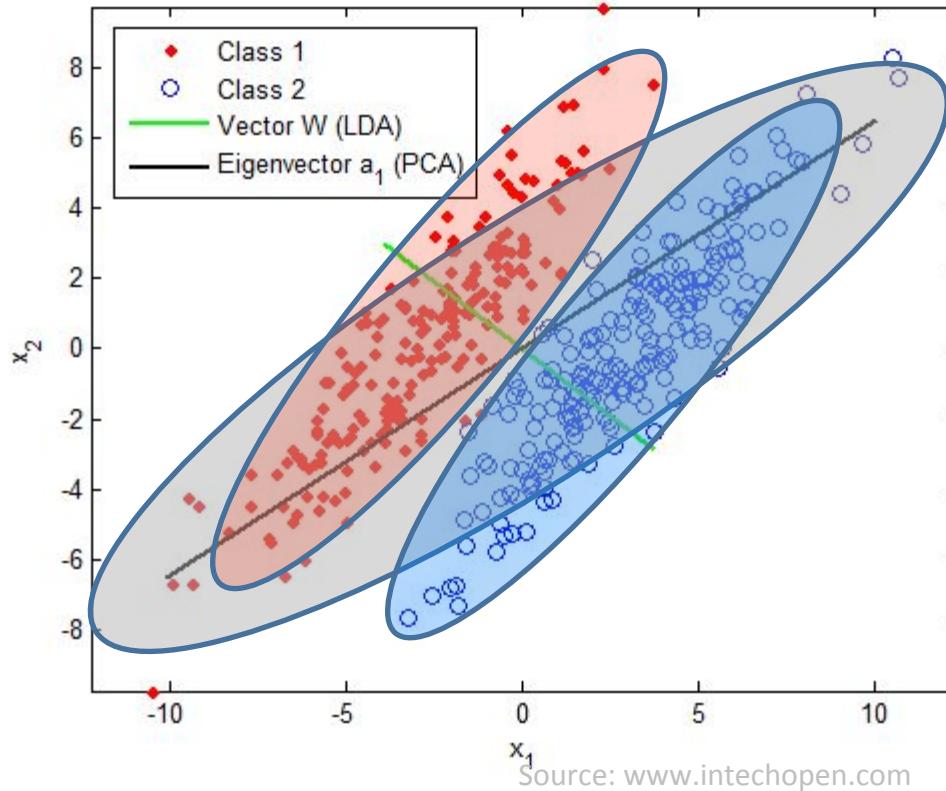
Hypothesize no covariance (even if there is some).
Rescale all axes by the standard deviation or the samples.

Sphering



Benefit: The weights of any linear method are on the same scale.
But we still have some covariance between features...

LDA and PCA difference

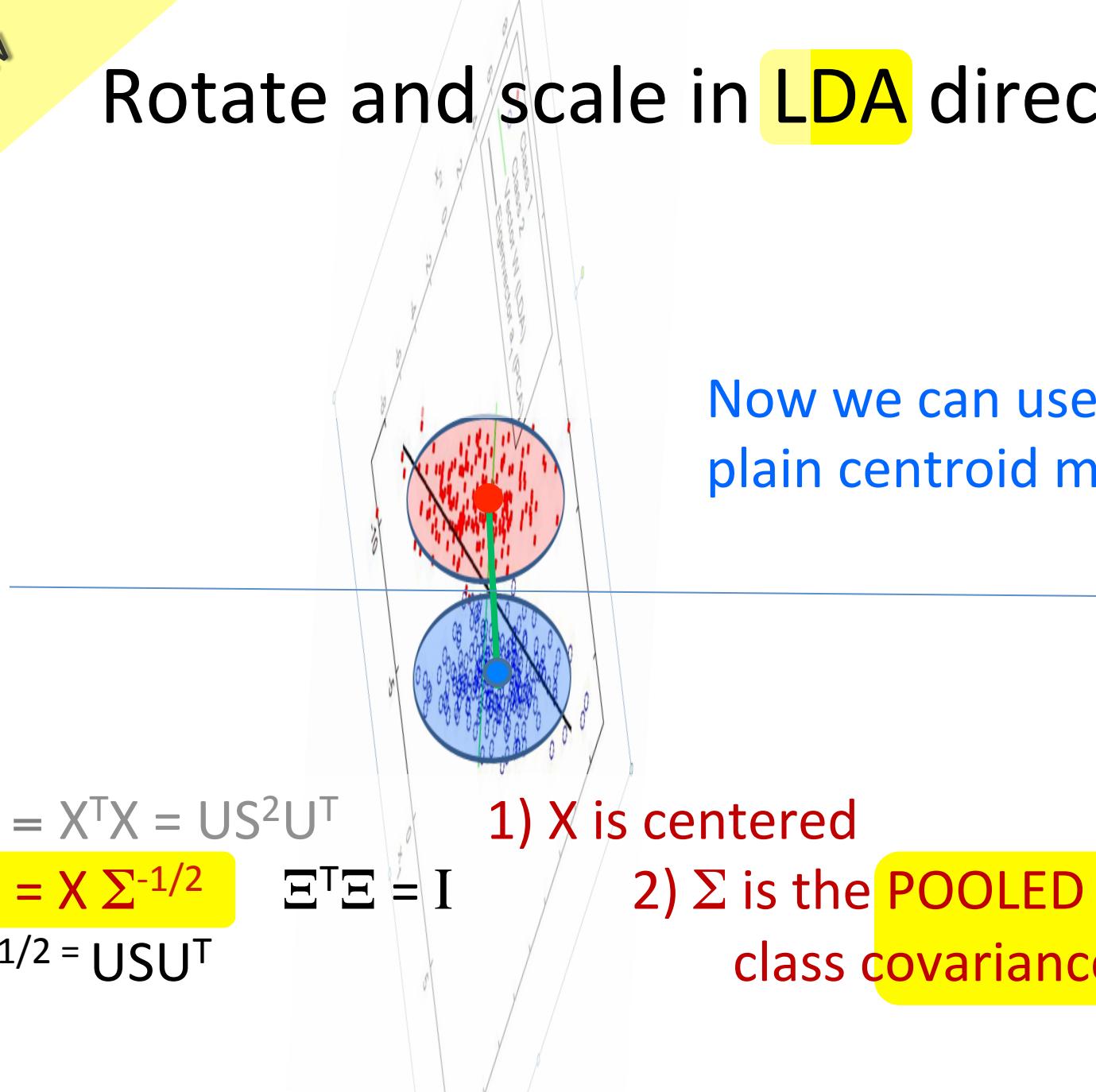


Apply whitening in both case to get $\Xi = X \Sigma^{-1/2}$

PCA: Σ is the TOTAL covariance matrix.

LDA: Σ is the POOLED WITHIN CLASS covariance matrix.

Rotate and scale in LDA directions



Now we can use the plain centroid method.

$$\Sigma = X^T X = U S^2 U^T$$

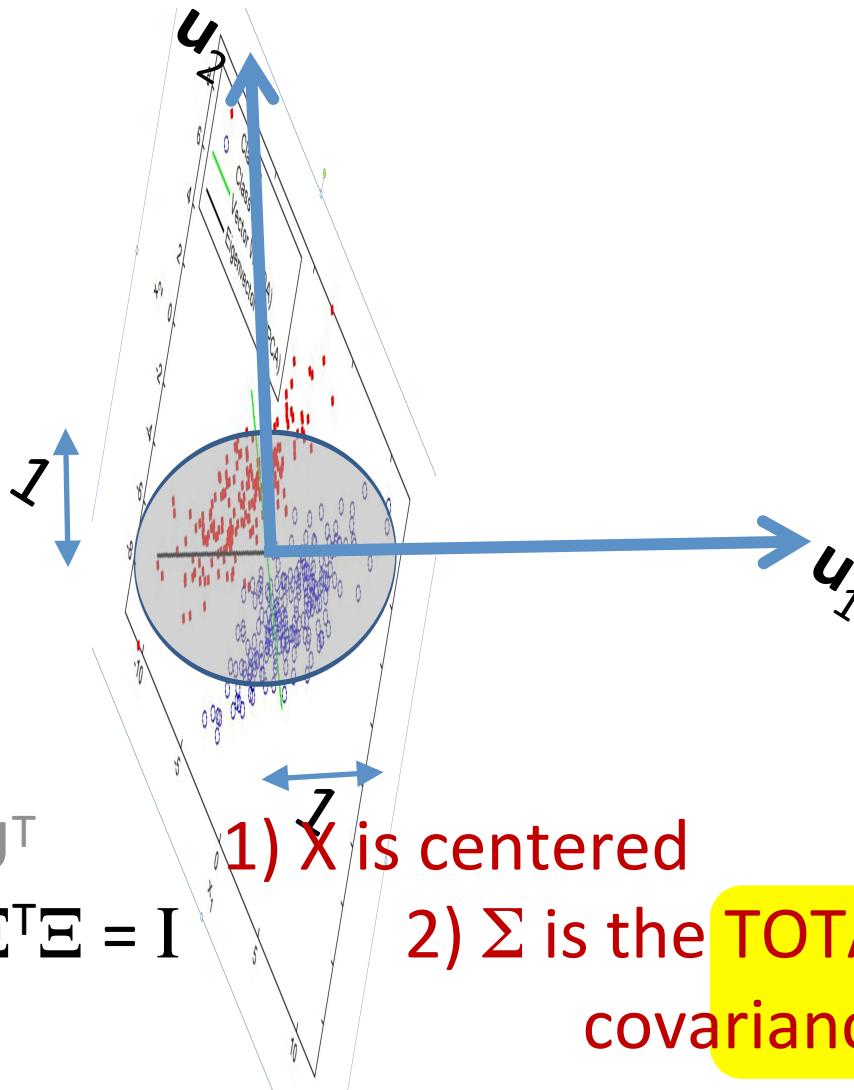
$$\Xi = X \Sigma^{-1/2}$$

$$\Sigma^{1/2} = U S U^T$$

1) X is centered

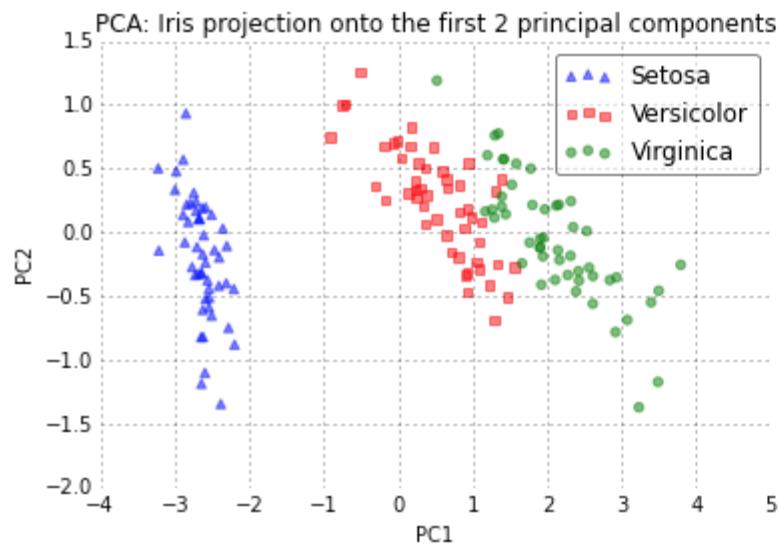
2) Σ is the POOLED within class covariance

Rotate and scale in PCA directions

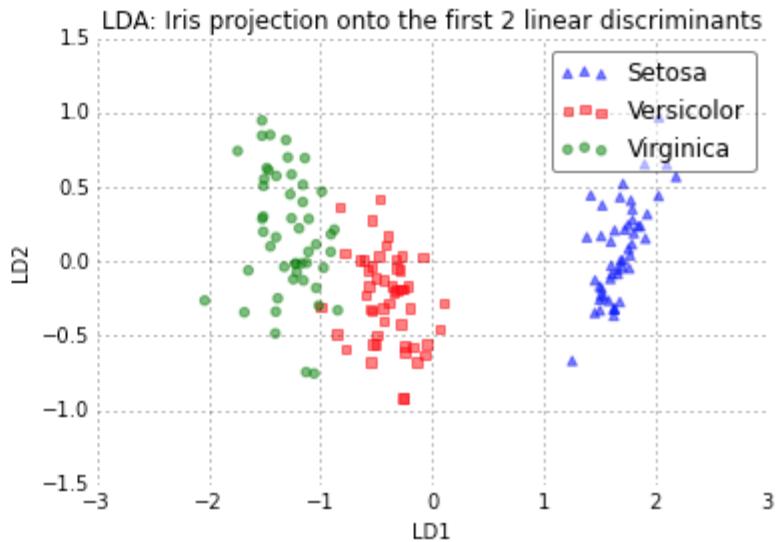


LDA and PCA visualization

PCA



LDA



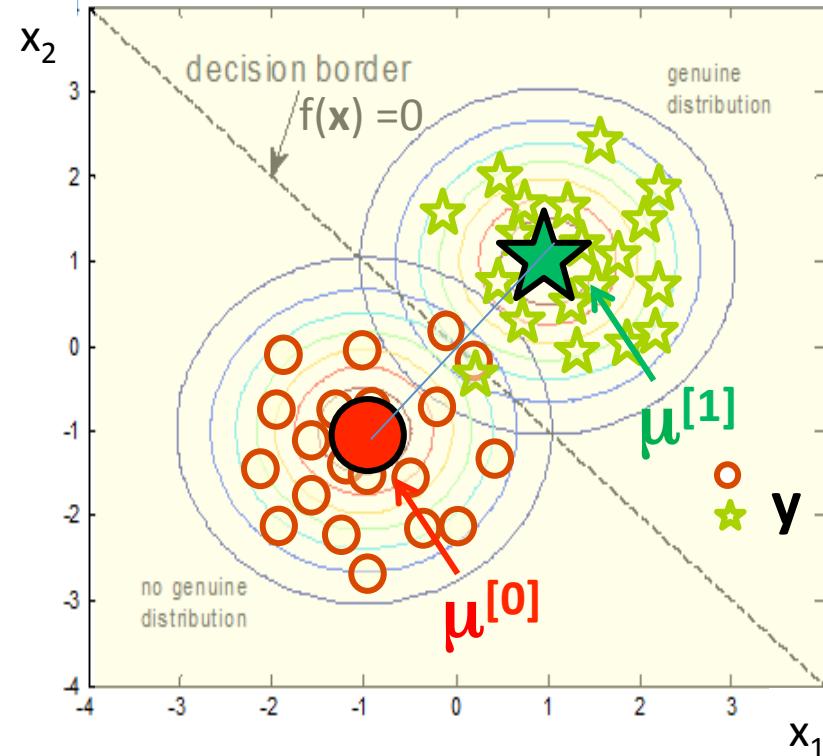
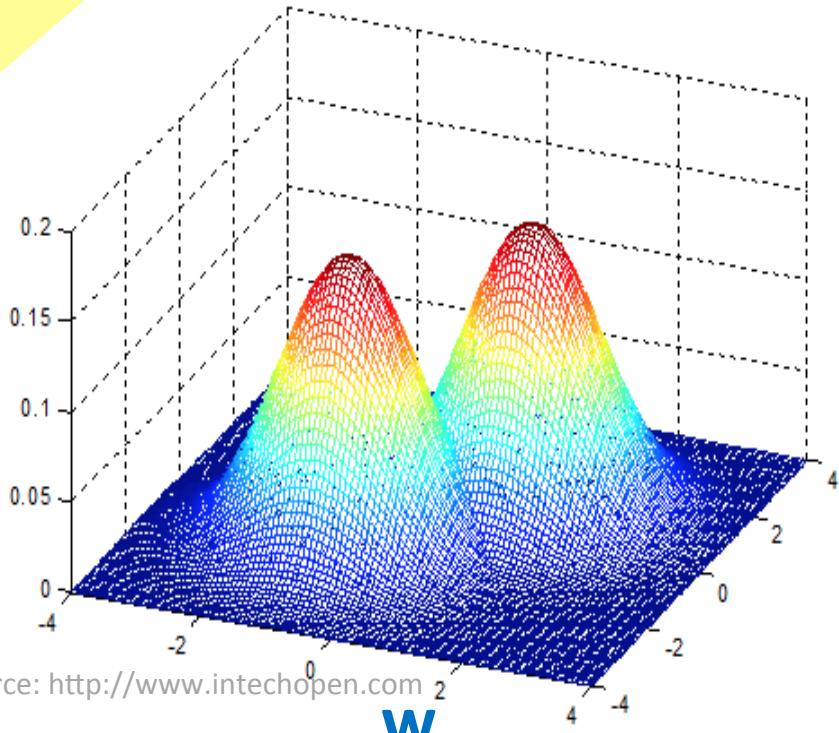
http://sebastianraschka.com/Images_old/2014_intro_supervised_learning/iris_pca_lda.png

Review

Gaussian classifier

$$\hat{y} = \operatorname{argmax}_y P(Y=y | X=x) \sim P(Y=y) P(X=x | Y=y)$$

$$\sim P(Y=y) \exp(-\|x - \mu^{[y]}\|^2 / 2\sigma^2)$$



$$f(x) = (\mu^{[1]} - \mu^{[0]}) \cdot x + b$$

$$f(x) = (\mu^{[1]}/\sigma^{[1]2} - \mu^{[0]}/\sigma^{[0]2}) \cdot x + b$$

$$f(x) = \Sigma^{-1}(\mu^{[1]} - \mu^{[0]}) \cdot x + b$$

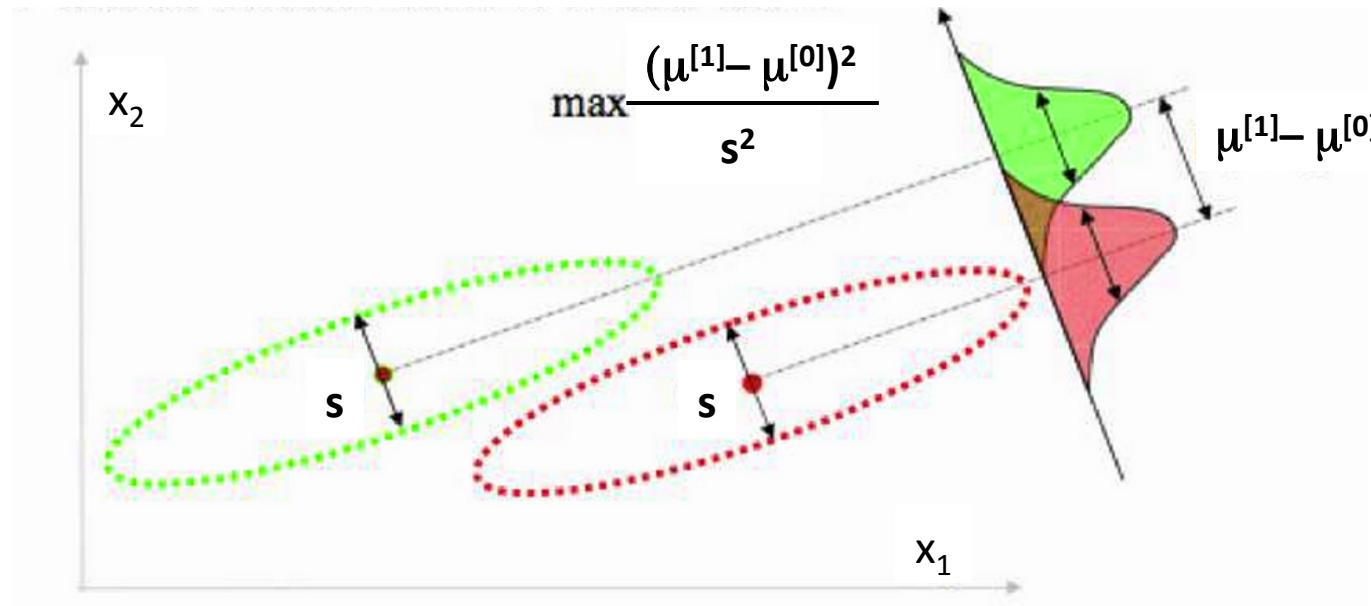
$$b = (\mu^{[0]2} - \mu^{[1]2})/2 + \log(N_1/N_0)$$

$$b = [(\mu^{[0]}/\sigma^{[0]2})^2 - (\mu^{[1]}/\sigma^{[1]2})^2]/2 + \log(N_1/N_0)$$

$$b = [(\mu^{[0]\top} \Sigma^{-1} \mu^{[0]} - \mu^{[1]\top} \Sigma^{-1} \mu^{[1]})/2 + \log(N_1/N_0)]$$

LDA is a “glorified” Gaussian classifier

$$P(X=x | Y=y) \sim \exp \left(-\frac{1}{2} \begin{matrix} [y] \\ (1, d) \end{matrix} \Sigma^{-1} \begin{matrix} [y] \\ (d, d) \end{matrix} \begin{matrix} [y] \\ (d, 1) \end{matrix}^T \right)$$



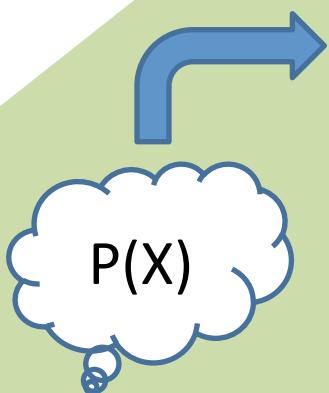
\mathbf{W}

$$f(\mathbf{x}) = \sum^{-1}(\boldsymbol{\mu}^{[1]} - \boldsymbol{\mu}^{[0]}) \cdot \mathbf{x} + b$$

$$b = [(\boldsymbol{\mu}^{[0]T} \Sigma^{-1} \boldsymbol{\mu}^{[0]} - \boldsymbol{\mu}^{[1]T} \Sigma^{-1} \boldsymbol{\mu}^{[1]})/2 + \log(N_1/N_0)]$$

Which assumption is right for you?

$$P(X, Y) = P(X)P(Y|X) = P(Y)P(X|Y)$$

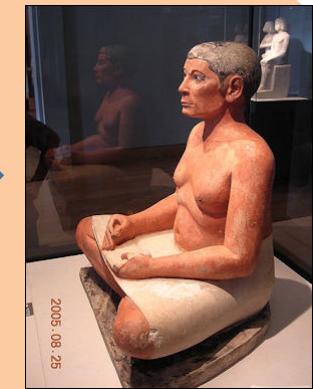
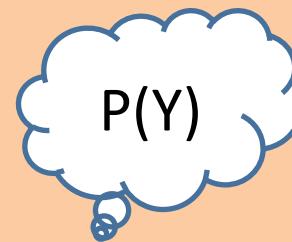


86	98	78	6	63	78	78	14	16	1
0	16	12	74	8	85	85	30	30	87
2	86	1	0	23	80	22	67	75	75
40	40	76	76	60	29	29	1	81	30
32	32	99	69	76	56	6	56	52	2
8	23	83	57	57	21	21	70	1	8
76	3	72	82	80	2	52	4	54	4
91	16	00	10	86	98	78	63	78	34



$P(Y|X)$

96 98 78 6 63 79...



$P(X|Y)$



5

It does not really matter!

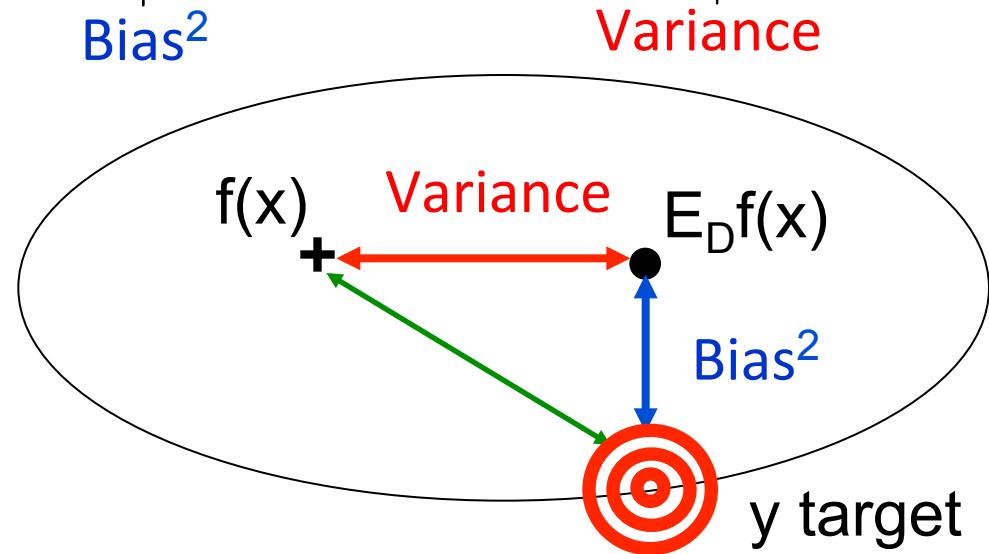
Bias vs. variance tradeoff



- f trained on a training set D of size m (m fixed)
- For the square loss:

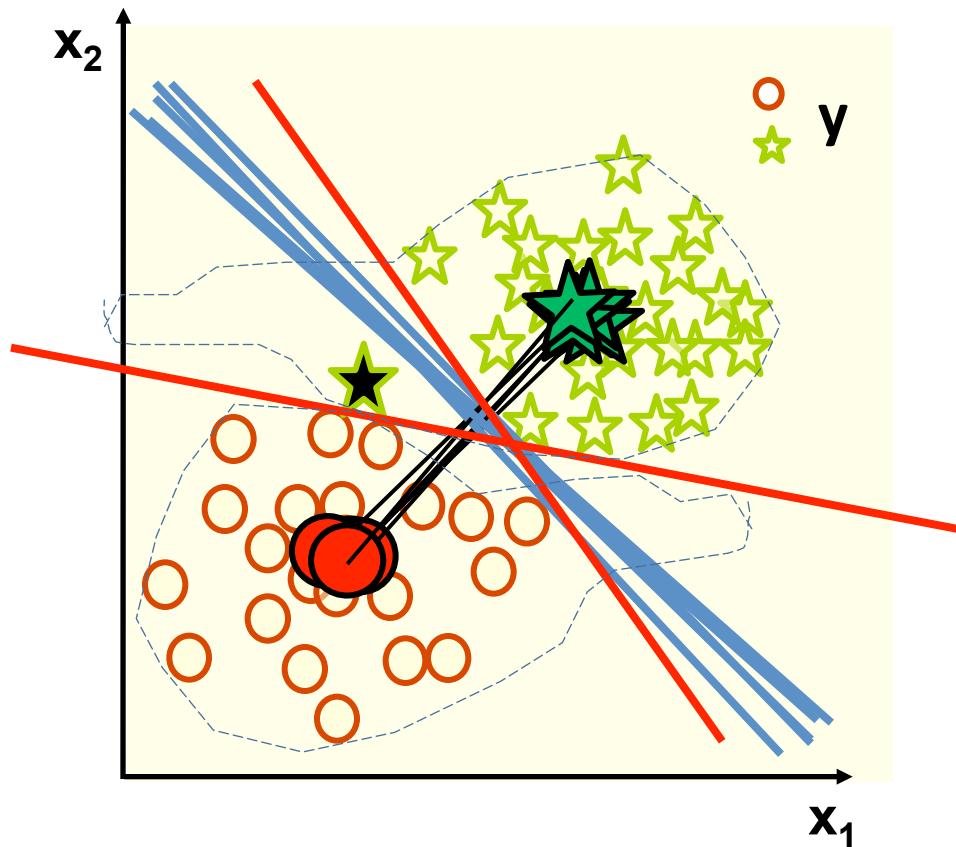
$$\underbrace{E_D[f(x)-y]^2}_{\text{Expected value of the loss over datasets } D \text{ of the same size}} = \underbrace{[E_D f(x) - y]^2}_{\text{Bias}^2} + \underbrace{E_D[f(x) - E_D f(x)]^2}_{\text{Variance}}$$

Expected value
of the loss over
datasets D of
the same size



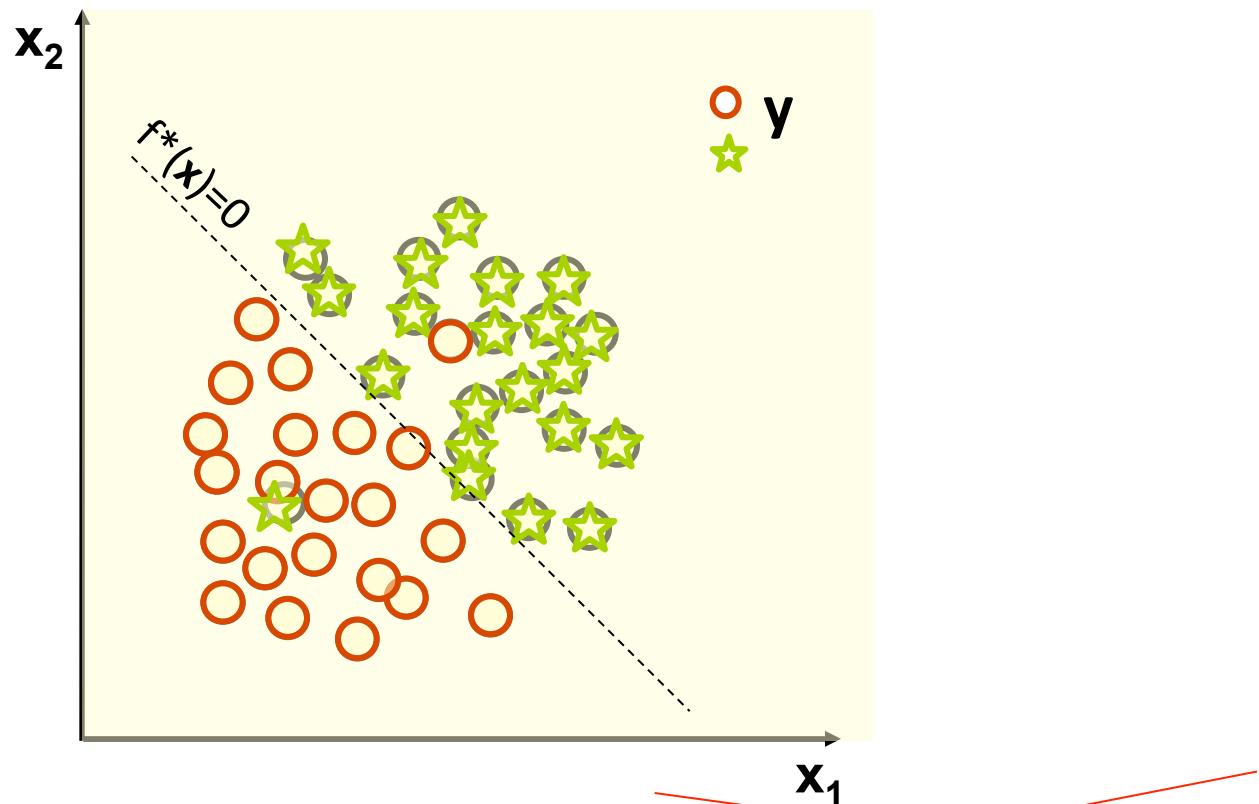
Violated assumption (1): Gaussianity

Centroid vs. SVM



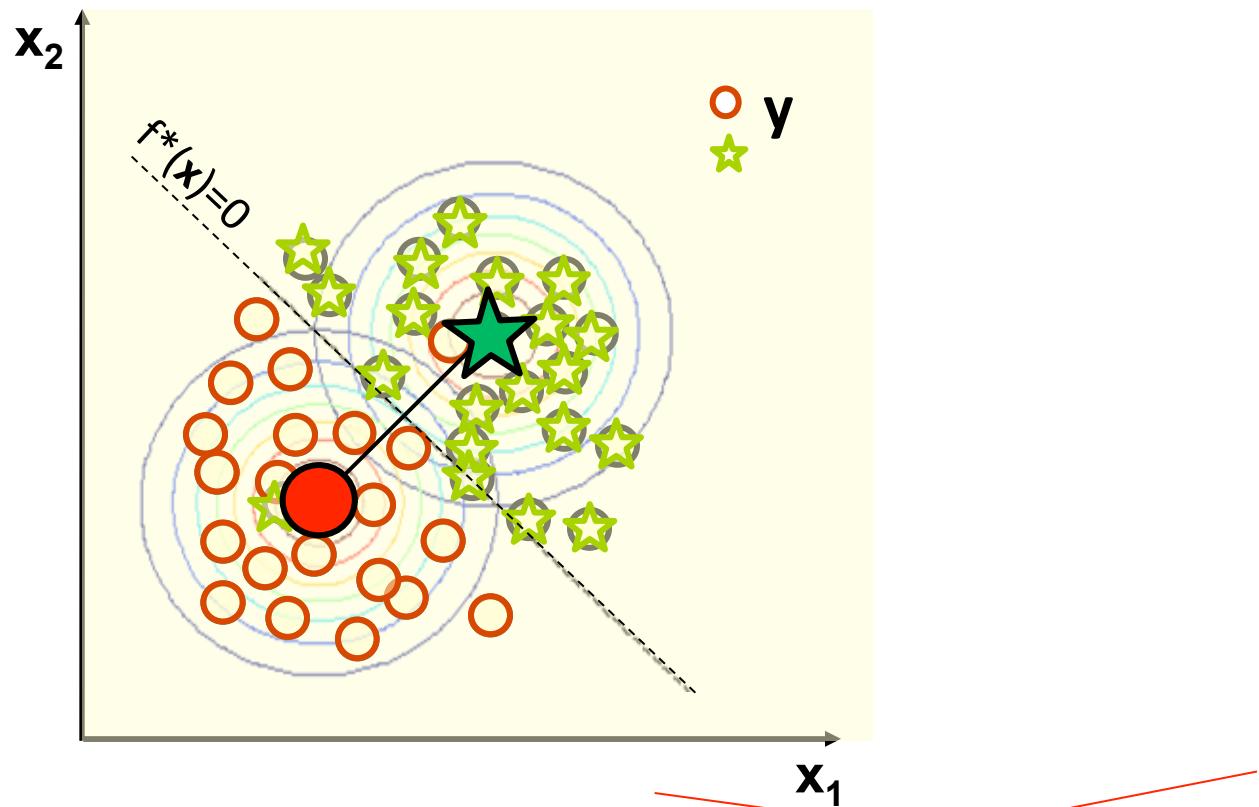
Noisy data or lack of training data: Even if the class clusters have a tailed distribution, you might be better off making Gaussian assumptions.

Violated assumption (2): wrong data generating “direction”



“Wrong” data generating assumption: $P(Y=y)$ then $\cancel{P(X=x|Y=y)}$
In reality: $P(X=x)$ then $P(Y=y|X=x)=1$ if $f^*(\mathbf{x})>0$, or 0 if $f^*(\mathbf{x})<0$.

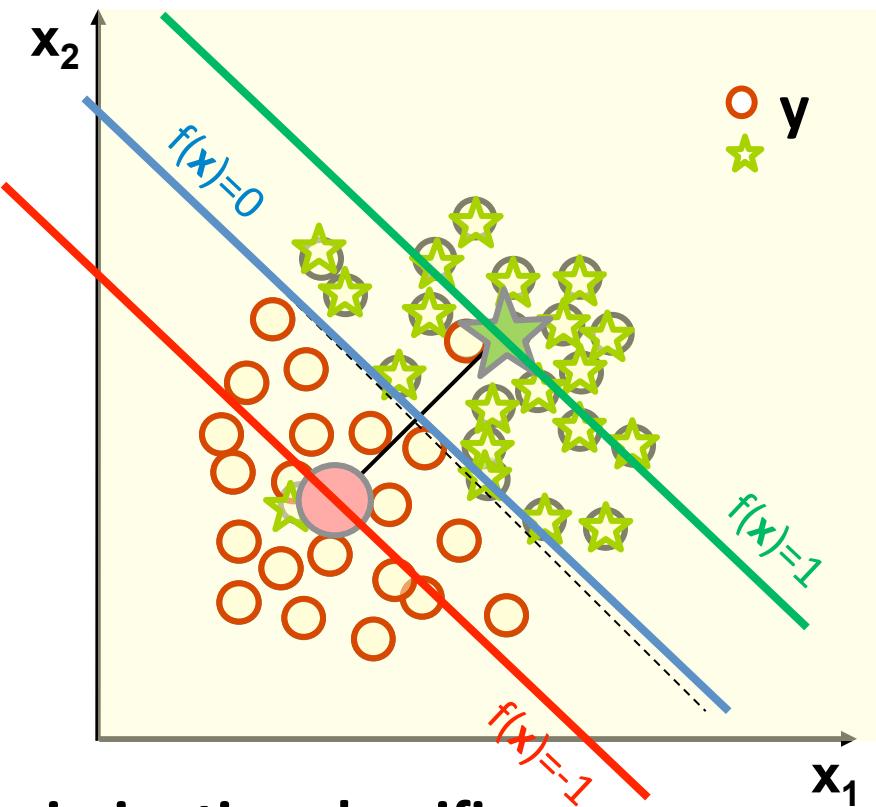
Violated assumption (2): $P(X=x | Y=y)$ not Gaussian



“Wrong” data generating assumption: $P(Y=y)$ then $P(X=x | Y=y)$
In reality: $P(X=x)$ then $P(Y=y | X=x)=1$ if $f^*(x)>0$, or 0 if $f^*(x)<0$.
The Gaussian model assumptions are violated.

Would the centroid method still be OK?

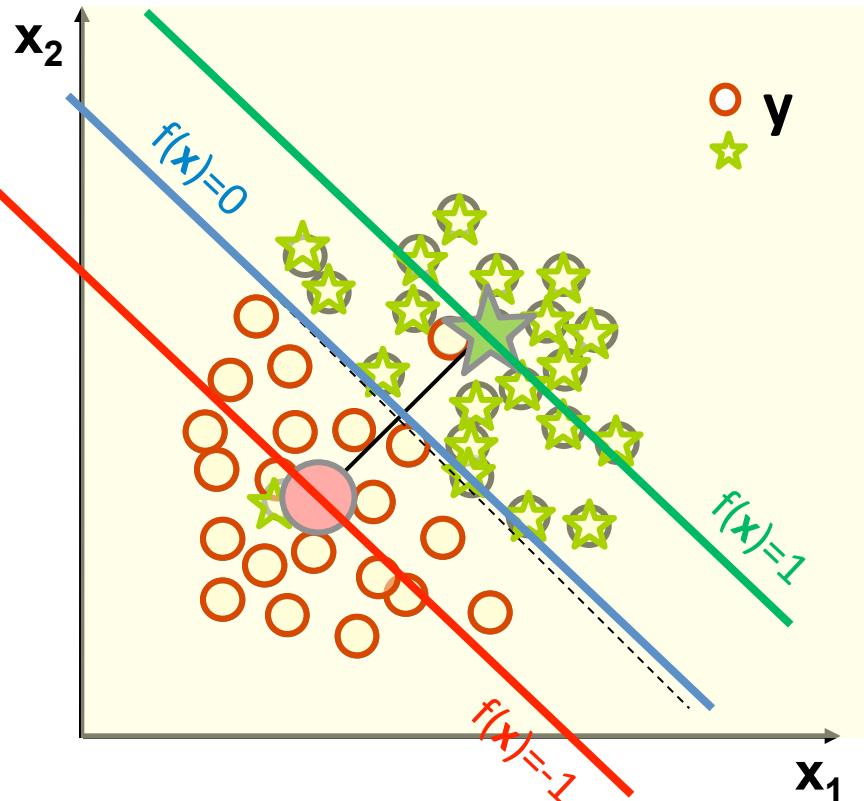
Let's compare with the regression solution



Regression = **discriminative classifier**:

- “Correct” data generating assumption: $P(X=x)$ then $P(Y=y | X=x)$
- but “wrong” loss function $(f(\mathbf{x}) - y)^2$ for classification? ($y=\pm 1$) 27

Explanation: Regression for classification ≡ LDA / Fisher linear discriminant



If:

- covariance matrix \equiv pooled within class covariance, and
- $y \in \{+1/N_1, -1/N_0\}$.

LDA:

$$\mathbf{w} = \Sigma^{-1}(\mu^{[1]} - \mu^{[0]})$$

RIDGE REGRESSION:

$$\mathbf{w} = \Sigma^{-1} \mathbf{X}^T \mathbf{y}$$

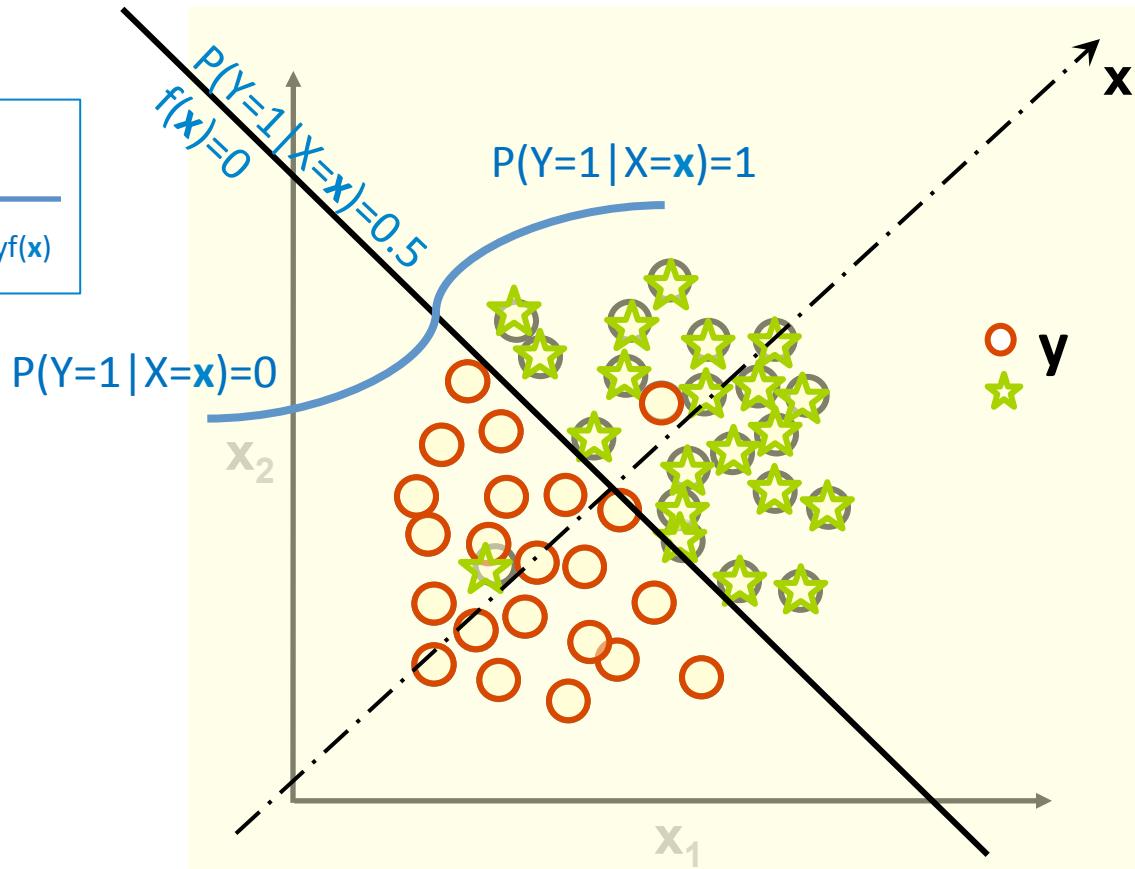
Regularized Σ^{-1}

$$\Sigma^{-1} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$$

What about logistic regression?



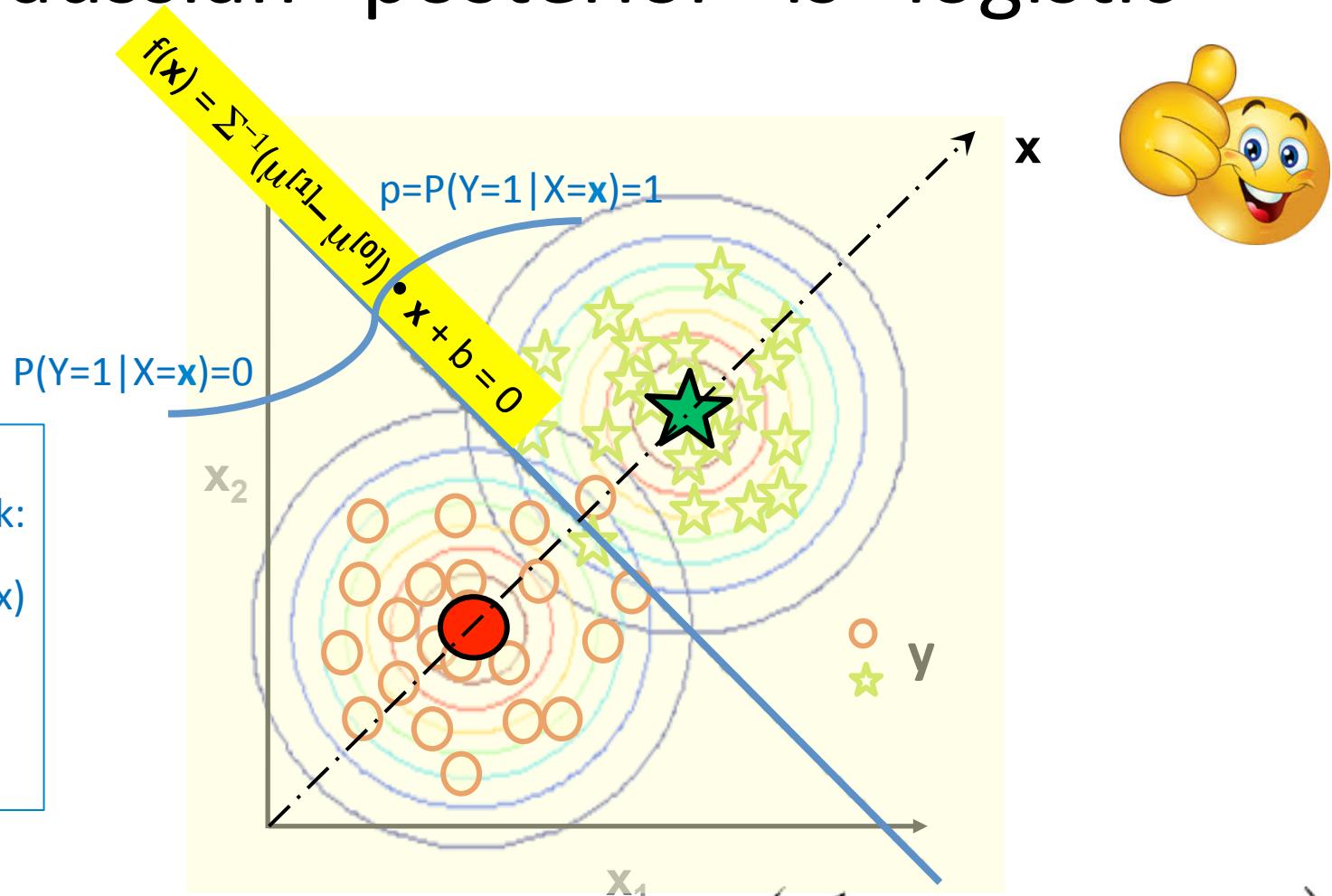
$$P(Y=y|X=x) = \frac{1}{1 + e^{-yf(x)}}$$



Logistic regression = **discriminative classifier**:

- “Correct” data generating assumption: $P(X=x)$ then $P(Y=y|X=x)$
- “Correct” loss function $-\log P(Y=y|X=x) = \log(1 + e^{-yf(x)})$ ($y=\pm 1$) ²⁹

The Gaussian “posterior” is “logistic”

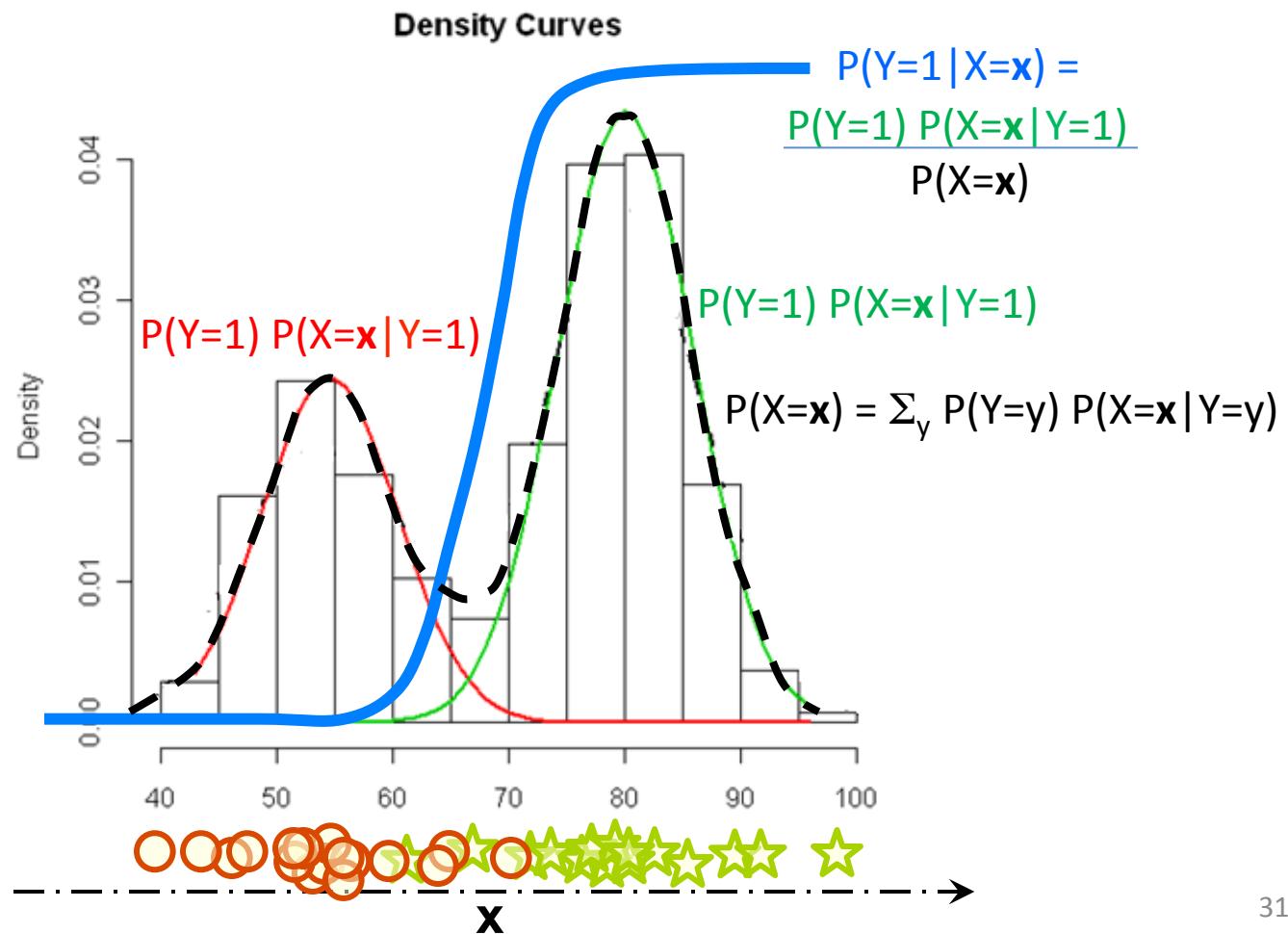


Generative model: $P(Y=y)$ then $P(X=x|Y=y) \sim \exp\left(-\frac{1}{2}(x - \mu^{[y]})^\top \Sigma^{-1}(x - \mu^{[y]})\right)$

We have seen that the log odds-ratio is a linear model.

So... The posterior $P(Y=1 | X=x)$ of the Gaussian model is logistic.

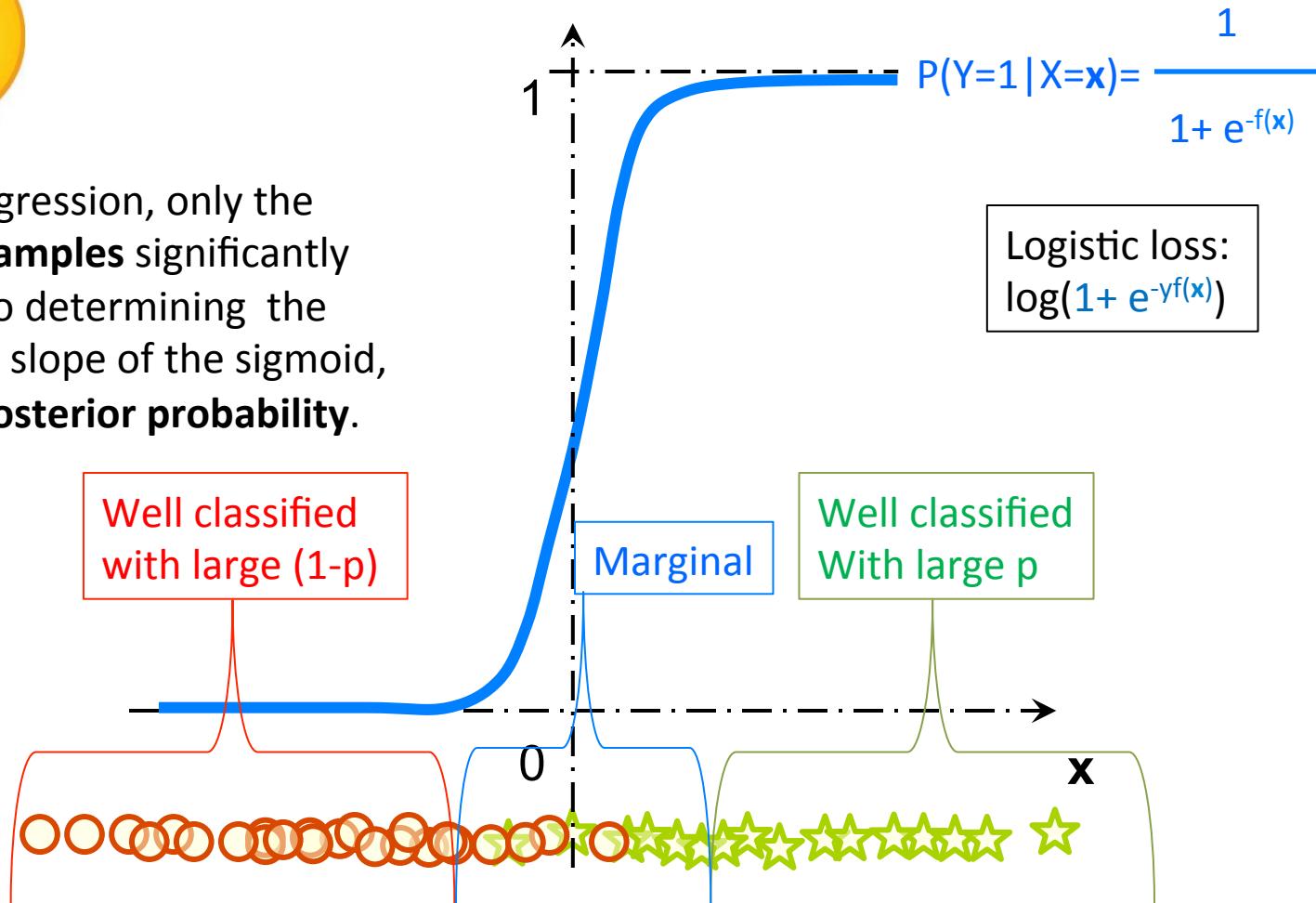
The posterior $P(Y=1 | X=x)$ of the Gaussian model is logistic



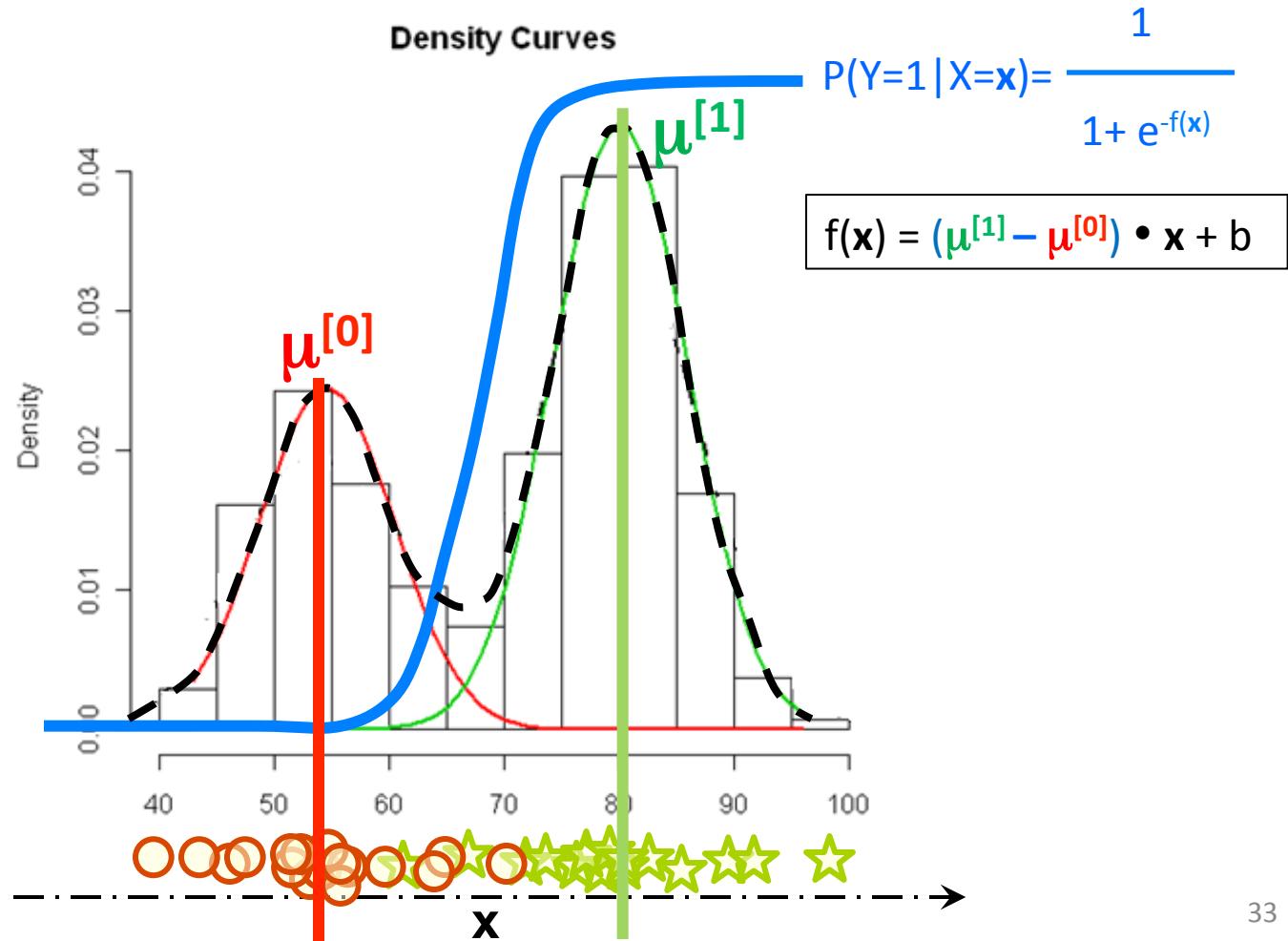
But the logistic model does not make any Gaussian assumption!



In logistic regression, only the **marginal examples** significantly contribute to determining the position and slope of the sigmoid, hence the **posterior probability**.



The Gaussian model posterior estimation uses the examples



Do we want to bother with LDA?

My own cookbook:



1) Preprocessing:

- Pattern normalization: $\mathbf{x}^k \leftarrow \mathbf{x}^k / \|\mathbf{x}^k\|$ (optional)
- Feature centering: $\mathbf{x}_i \leftarrow \mathbf{x}_i - \mu_i$
- Singular value decomposition: $\mathbf{X} = \mathbf{V}\mathbf{S}\mathbf{U}^T$
 $\mathbf{U}^T\mathbf{U} = \mathbf{I}$, $\mathbf{V}^T\mathbf{V} = \mathbf{I}$, \mathbf{S} diagonal $\text{dim}(\mathbf{r}, \mathbf{r})$, $\mathbf{r} = \text{rank}(\mathbf{X}) \leq \min(d, N)$
- Visualize data in 2 dim (top 2 singular values) $\Xi = \mathbf{X} [\mathbf{u}^1, \mathbf{u}^2]$
 - Further preprocess (e.g. non-linear scaling) $(N, d) \rightarrow (d, 2)$
 - Eventually replace Σ by pooled Σ
 - Reduce space dimension with PCA or LDA $\Xi = \mathbf{X} [\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^{\max}]$

2) Linear methods:

1. Centroid (with targets $y \in \{+1/N_1, -1/N_0\}$ or $y = \pm 1$)
2. Soft margin SVM (\sim ridge regression for large λ)
3. Post-fit SVM output to sigmoid (skip logistic regression)

3) Non-linear methods

Double random process

$P(Y)$



5



5

One style
 $P(X|Y)$

5



$P(Y)$

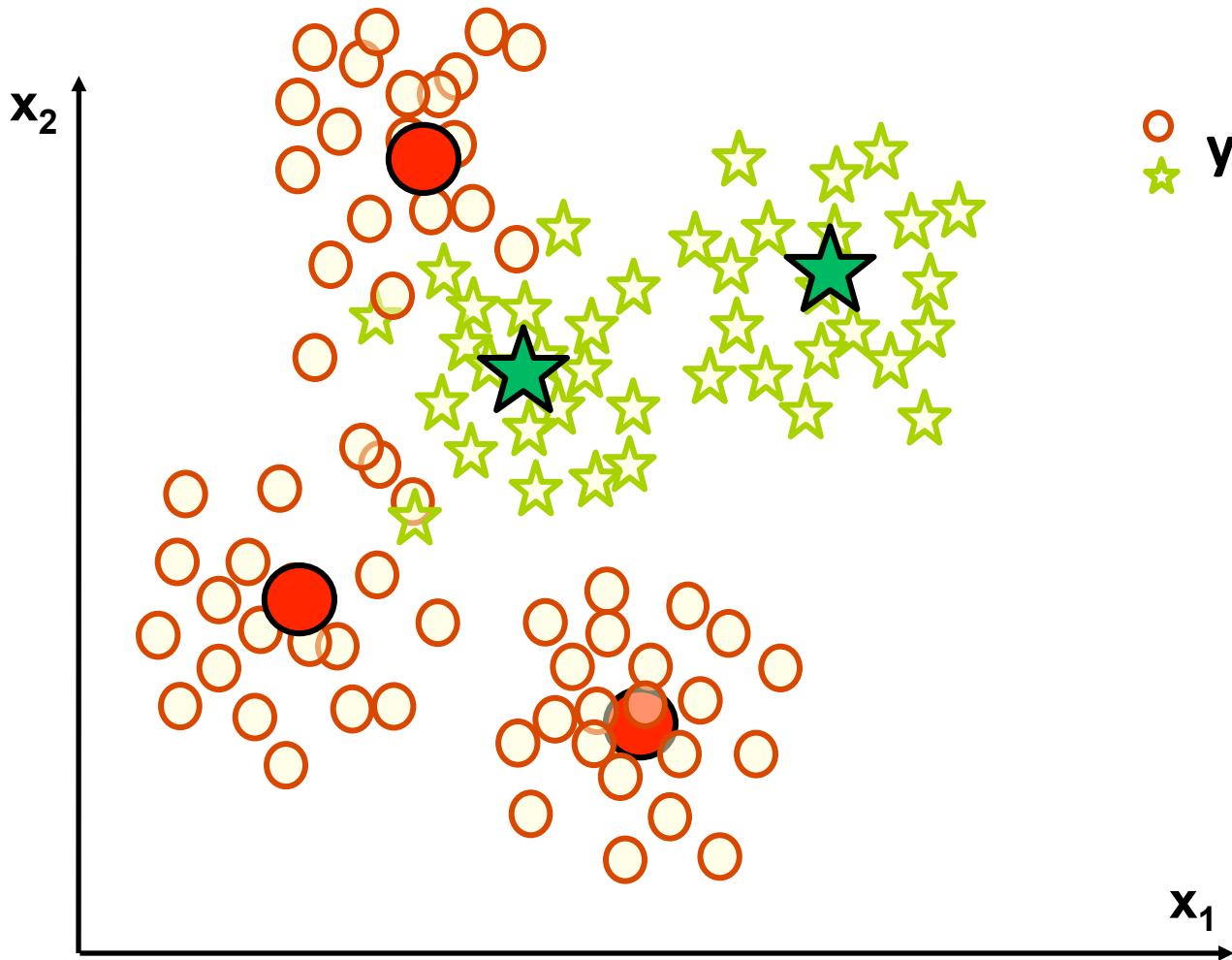
The image cannot be displayed. Your computer may not have enough memory to open the image, or the image may have been corrupted. Restart your computer, and then open the file again. If the red x still appears, you may have to delete the image and then insert it again.



5

$P(X|Y)$ Several styles²⁵

Multiple clusters per class



Mixture models

$$P(Y=1 | X=x) \sim$$

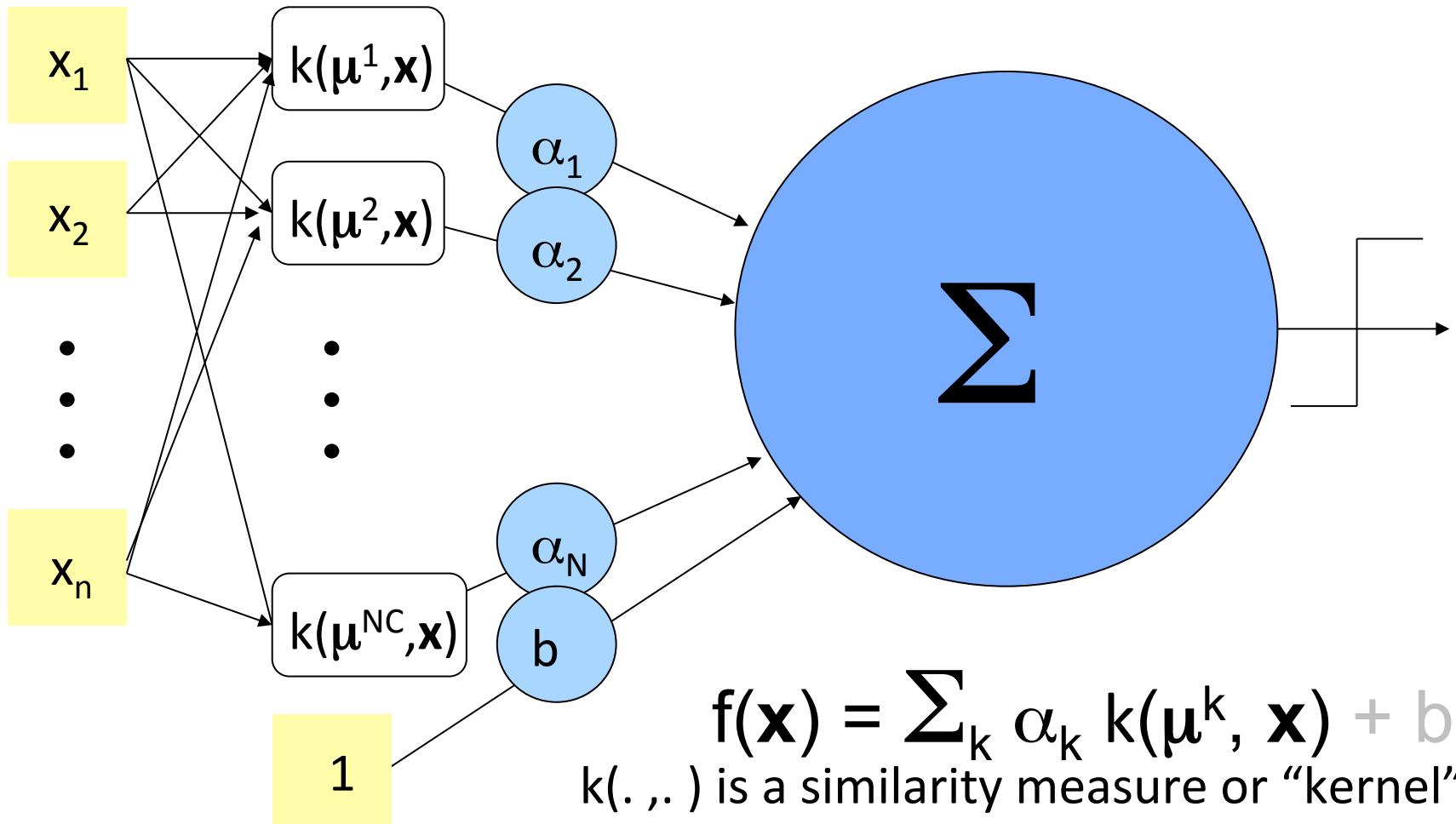
$$\begin{aligned} & \underbrace{P(X=x | Y=y)}_{\text{likelihood}} \underbrace{P(Y=y)}_{\text{prior}} = \sum_k P(X=x, S=s_k | Y=y) P(Y=y) \\ & = \sum_k \underbrace{P(X=x | S=s_k, Y=y)}_{\sim \exp(-\|x-\mu_k\|^2/2\sigma^2)} \underbrace{P(S=s_k | Y=y)}_{\sim \alpha_k} P(Y=y) \end{aligned}$$

$$f(x) = P(Y=1 | X=x) - P(Y=-1 | X=x)$$

$$\sim \sum_{k=1:N_c} \alpha_k \exp(-\|x-\mu_k\|^2/2\sigma^2) + b$$

RBF network

RBF = radial basis function



Parameter estimation

- $f(\mathbf{x}) = \sum_{k=1:N_c} \alpha_k k(\mu^k, \mathbf{x}) + b$

$$k(\mu^k, \mathbf{x}) = \exp(-\|\mathbf{x}-\mu_k\|^2/2\sigma^2) \quad \text{Gaussian kernel}$$

- **Parameters:**
 - N_c = number of clusters
 - μ_k = cluster centers
 - α_k = cluster weights
 - σ = kernel width
- **Simple way:**
 - Fix σ .
 - Fix N_c and take $N_c/2$ in each class.
 - In each class run k-means clustering to get the cluster centers μ_k .
 - Treat $f(\mathbf{x})$ as a model linear in its parameters and fit the α_k by gradient descent.
 - Optimize σ and N_c by cross-validation.

Or, yet simpler: Parzen windows

- $f(\mathbf{x}) = \sum_{k=1:N} y_k k(\mathbf{x}^k, \mathbf{x}) + b$

$$k(\mathbf{x}^k, \mathbf{x}) = \exp(-\|\mathbf{x}-\mathbf{x}_k\|^2/2\sigma^2) \quad \text{Gaussian kernel}$$

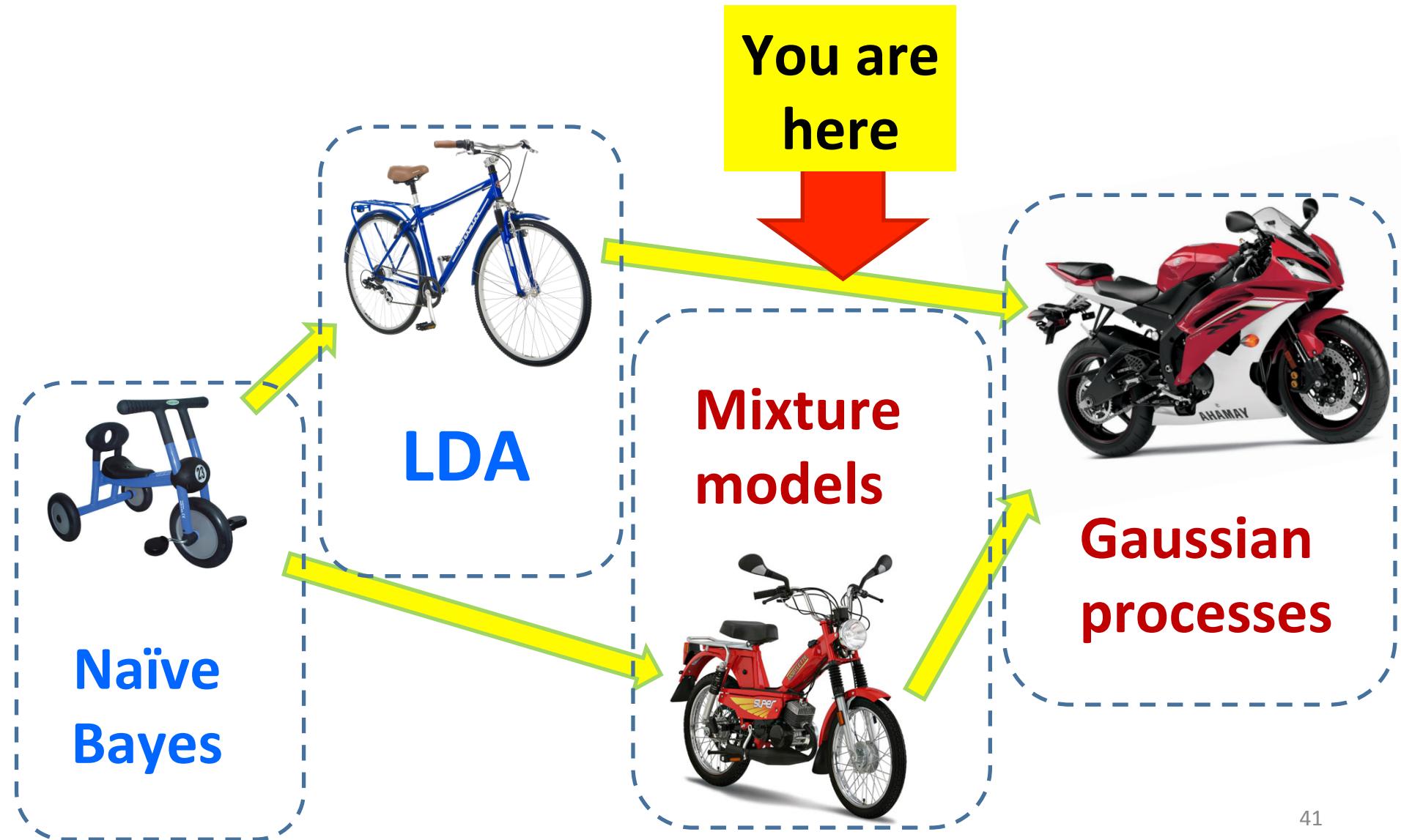
- **Parameters:**

- N_c = number of clusters
- μ_k = cluster centers
- α_k = cluster weights
- σ = kernel width

- **Simpler way:**

- There is one “cluster” per example ($N_c = N$, $\mu_k = \mathbf{x}^k$).
- The α_k are also fixed: $\alpha_k = y_k$
- Optimize σ by cross-validation.

Generative models



Summary

- We are trying to make sense of the models we are working with from a **data generating process** point of view. Useful for gaining insight, and potentially gaining performance advantages (by injecting “prior knowledge”):
 - **Visualization** allows us to evaluate our hypotheses.
 - Even if the hypotheses are wrong, we may get good results if the **data are noisy** or if we have **few training samples**, because of the “**bias/variance**” **tradeoff**.
 - SVM and logistic regression do not make data generating assumptions. But, ridge regression applied to classification problems is similar to LDA.
- **Gaussian mixtures:** What if the single cluster per class hypothesis is violated? We can introduce multiple clusters per class; this is the idea of Gaussian mixtures and several “Radial Basis Function” (RBF) algorithms. This bridges the gap with kernel methods.

Come to my office hours...

Wed 1:30-3:30 Soda 329

Next time: Gaussian processes

