

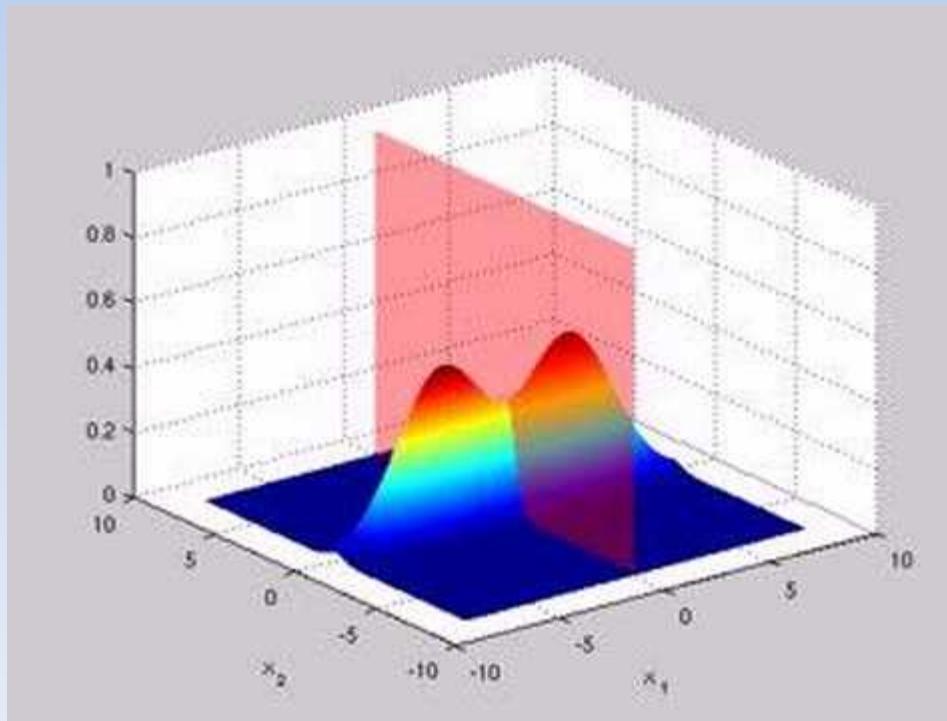
UCB - CS189
Introduction to Machine Learning
Fall 2015

Lecture 13: Linear Discriminant
Analysis (LDA) and QDA

Isabelle Guyon
ChaLearn

Come to my office hours...
Wed 2:30-4:30 Soda 329

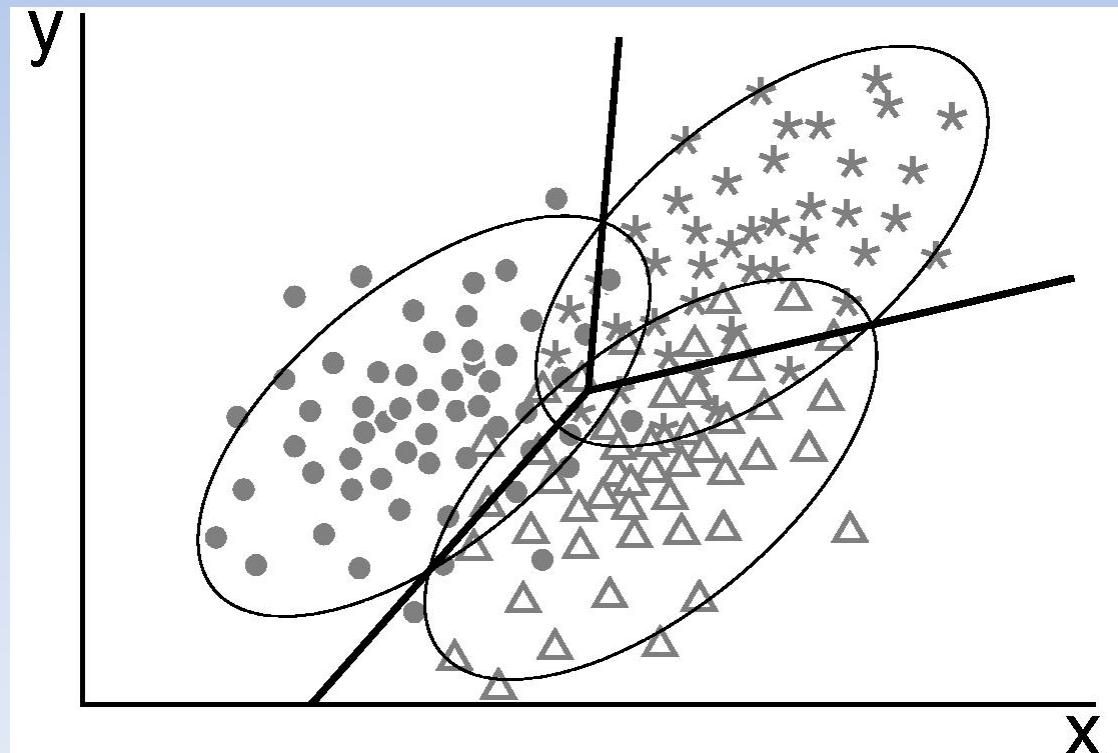
Last time: Gaussian classifier aka centroid method aka naïve Bayes aka Hebb's rule revised



Come to my office hours...

Wed 2:30-4:30 Soda 329

Today: LDA



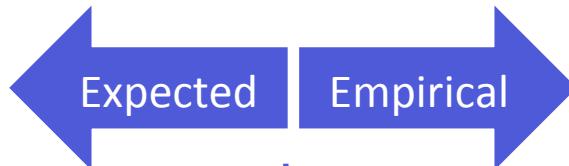
Math prerequisites

Independence, covariance, correlation, and dot product

Independence: $P(X, Y) = P(X)P(Y)$

Mean:

$$\mu_x = E(X)$$



$$\mu_x = (1/N) \sum_{k=1:N} x^k$$

Variance:

$$\sigma_x^2 = E[(X - \mu_x)^2]$$

$$\sigma_x^2 = (1/N) \sum_{k=1:N} (x^k - \mu_x)^2$$

Covariance:

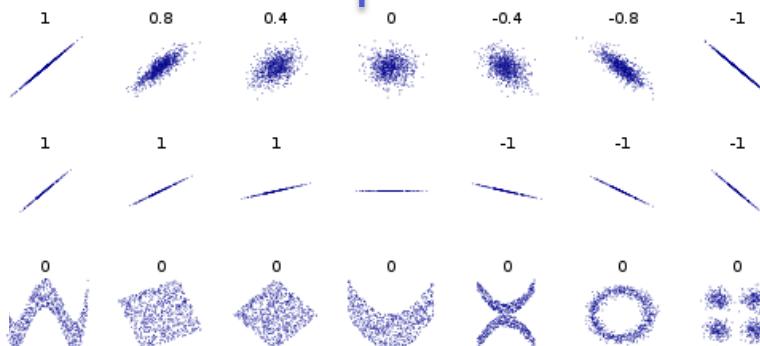
$$\sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)]$$

Covariance:

$$\sigma_{XY} = (1/N) \sum_{k=1:N} (x^k - \mu_x)(y^k - \mu_y)$$

Pearson correlation:

$$\rho_{X,Y} = \frac{\sigma_{xy}}{\sigma_X \sigma_Y}$$



Pearson correlation:

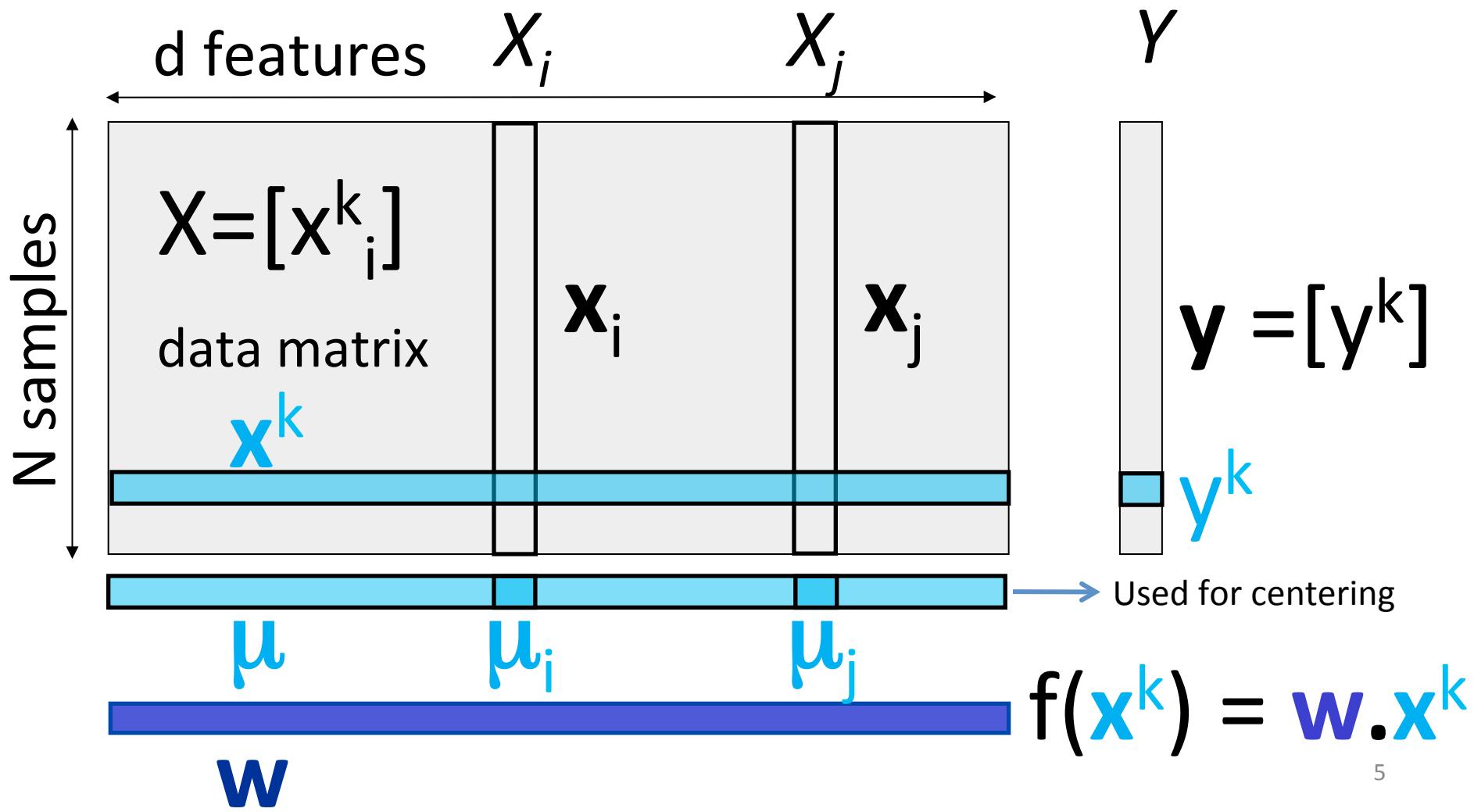
$$\rho_{X,Y} = \frac{\sigma_{xy}}{\sigma_X \sigma_Y}$$

Link to dot product:

$$E(X Y)$$

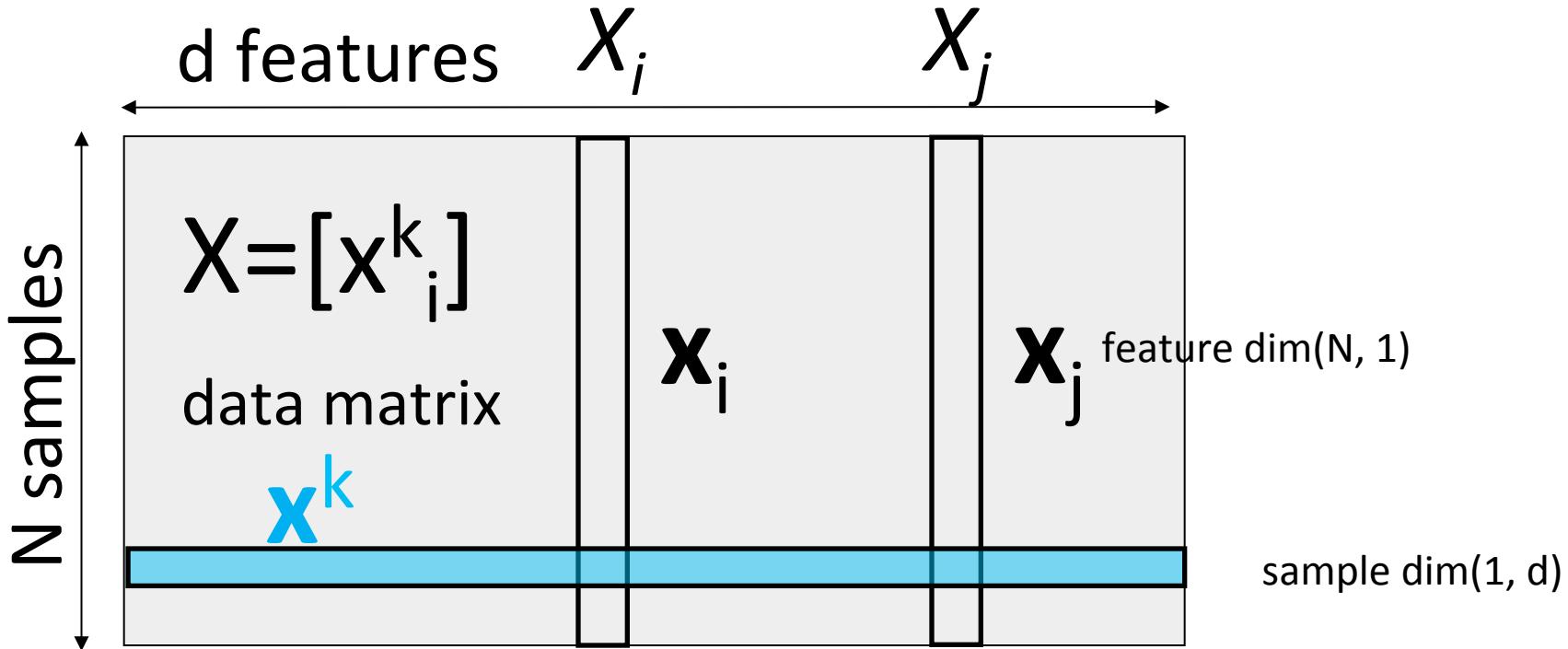
$$(1/N) \mathbf{x} \cdot \mathbf{y} = (1/N) \sum_{k=1:N} x^k y^k$$

What could be X and Y?



Outer product

X is a line random vector $X = [X_1 \ X_2 \ \dots \ X_i \ \dots \ X_d]$



$\Sigma = E(X^T X)$ covariance matrix (if X has 0 mean), dim (d,d)

$\Sigma = (1/N) X^T X$ outer product = empirical covariance matrix , dim (d,d)

Σ is positive semi-definite and has therefore a square root $R = \Sigma^{(1/2)}$

$\Phi = X R^{-1}$ is the “whitened data” $\Phi^T \Phi = I$

Data transforms

$\mathbf{x}^k = [x_1^k \ x_2^k \ \dots \ x_i^k \ \dots \ x_d^k]$ is a line of the data matrix \mathbf{X} .

Mean:

$$\boldsymbol{\mu} = (1/N) \sum_{k=1:N} \mathbf{x}^k$$

Variance:

$$\sigma_i^2 = (1/N) \sum_{k=1:N} (x_i^k - \mu_i)^2$$

Covariance:

$$\sigma_{ij} = (1/N) \sum_{k=1:N} (x_i^k - \mu_i)(x_j^k - \mu_j)$$

(d,d) covariance matrix:

$$\boldsymbol{\Sigma} = \boldsymbol{\sigma}_{ij} = (1/N) \boldsymbol{\Xi}^T \boldsymbol{\Xi}$$

(N, d) centered data matrix:

$$\boldsymbol{\Xi} = [\xi_i(\mathbf{x}^k)] = [x_i^k - \mu_i]$$

Centering: subtracting the mean of the features.

$$\boldsymbol{\xi}(\mathbf{x}^k) = \mathbf{x}^k - \boldsymbol{\mu}$$

Standardizing or “sphering”: subtracting the mean and dividing by the standard deviation (component-wise).

$$\boldsymbol{\phi}(\mathbf{x}^k) = (\mathbf{x}^k - \boldsymbol{\mu}) ./ \boldsymbol{\sigma}$$

$$\phi_i(\mathbf{x}^k) = (x_i^k - \mu_i) / \sigma_i$$

Whitening: multiplying by the square root of the inverse covariance matrix.

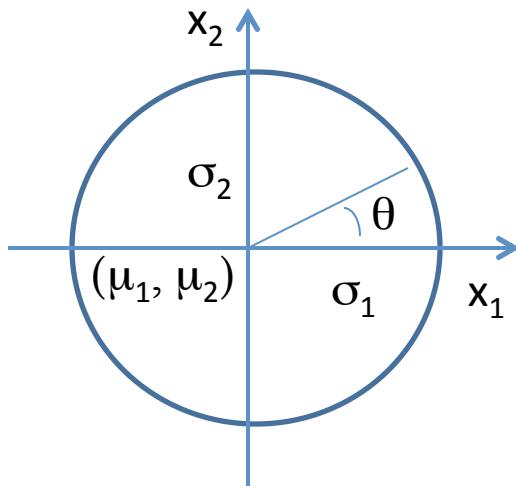
$$\boldsymbol{\Phi} = \boldsymbol{\Xi} \ \boldsymbol{\Sigma}^{-1/2}$$

$$(N, d) = (N, d) (d, d)$$

Math prerequisites

Ellipses

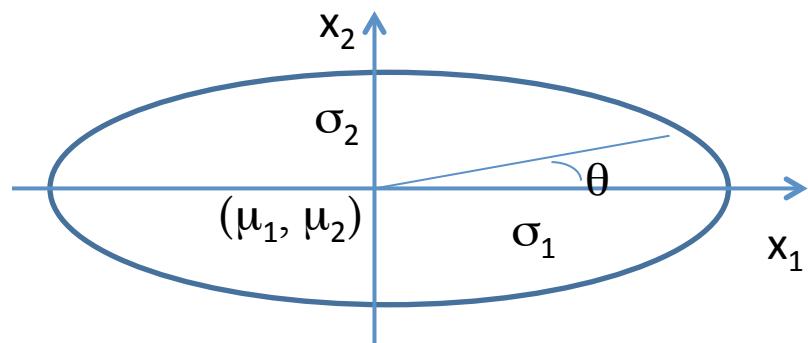
An ellipse in its principal axes is a flattened circle:



circle: $\sigma_1 = \sigma_2$

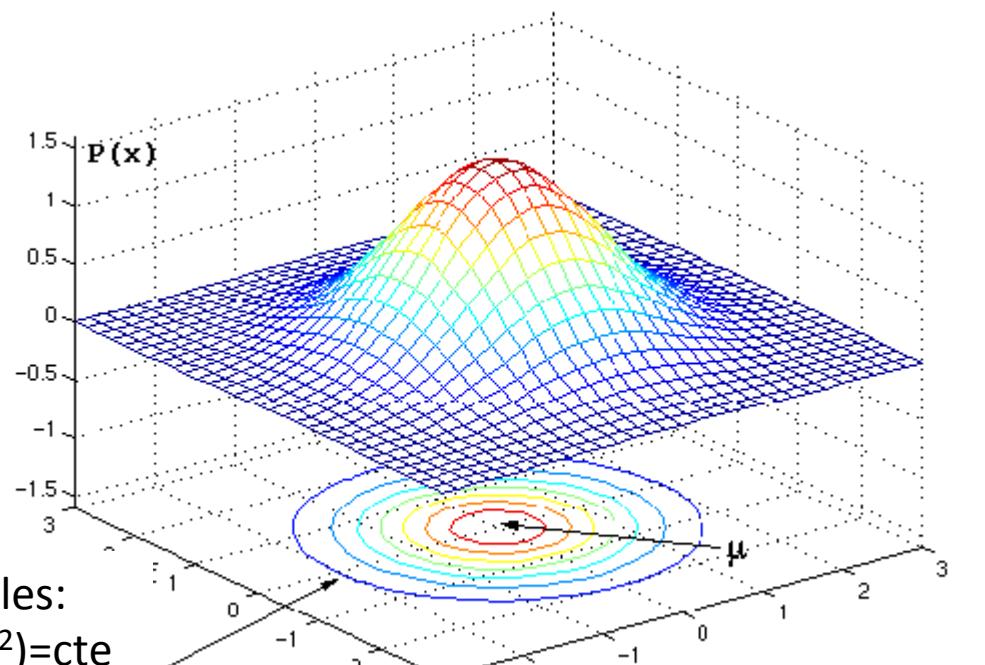
$$x_1 - \mu_1 = \sigma_1 \cos \theta$$
$$x_2 - \mu_2 = \sigma_2 \sin \theta$$

$$\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 = 1$$



ellipse: $\sigma_1 \neq \sigma_2$

Isotropic Gaussian $\sim \exp(-\|\mathbf{x}-\mu\|^2/2\sigma^2)$

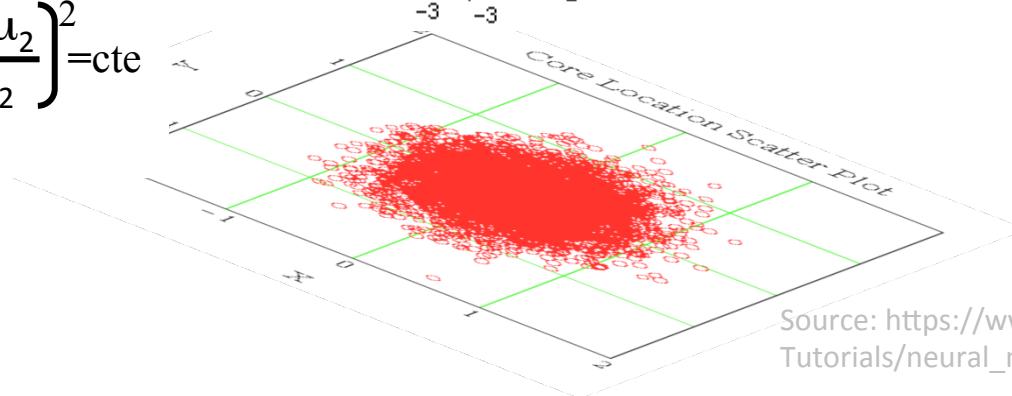


Contours = circles:

$$\exp(-\|\mathbf{x}-\mu\|^2/2\sigma^2)=\text{cte}$$

$$\left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2 = \text{cte}$$

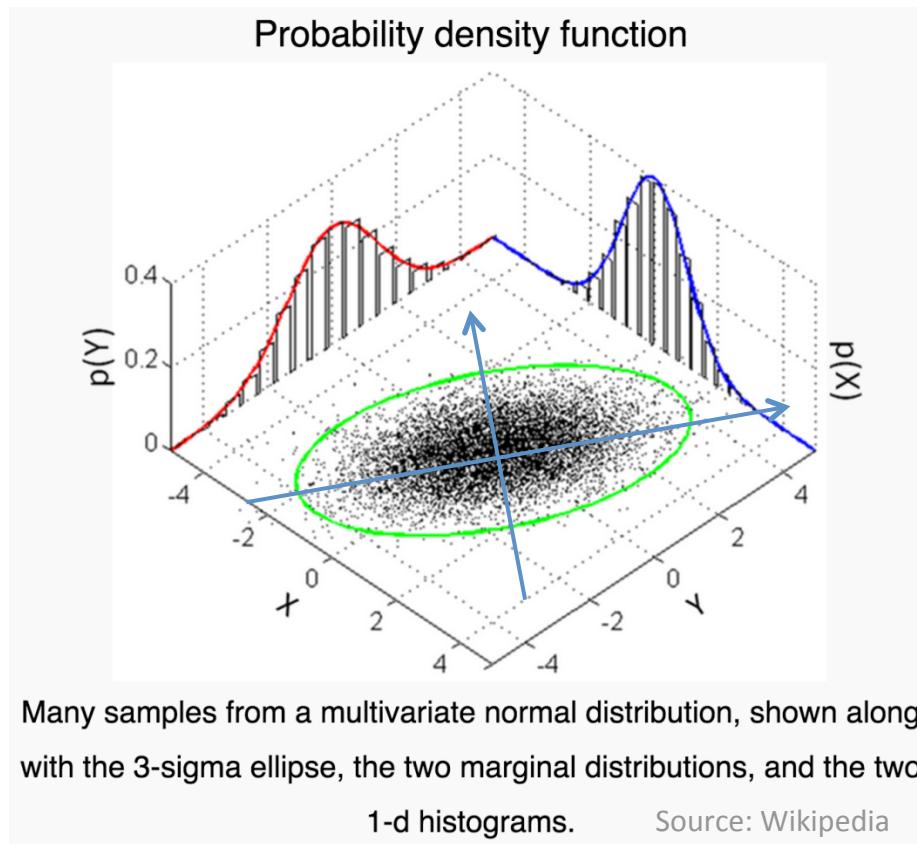
$$\sigma_1 = \sigma_2$$



Source: https://www.byclb.com/TR/Tutorials/neural_networks/ch4_1.htm

Multivariate Gaussian

$$\sim \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) = \exp(-(1/2)\phi\phi^T)$$



$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}$$

Contours in the principal axes:

$$\left(\frac{\phi_1 - \mu_{\phi_1}}{\sigma_{\phi_1}}\right)^2 + \left(\frac{\phi_2 - \mu_{\phi_2}}{\sigma_{\phi_2}}\right)^2 = \text{cte}$$

How to get the principal axes?

- $X = [x_i^k]$ is our data matrix (N lines = samples, d columns = feature).
- Assume the sample means $\mu_i = (1/N) \sum_k x_i^k$ were subtracted from columns x_i already.
- Then $X^T X$ is the (d, d) covariance matrix of the features (up to a factor N).
- We can diagonalize it as: $\Sigma = X^T X = US^2U^T$

With U and orthogonal matrix of column eigenvectors, such that $U^T U = I$.

S^2 is a diagonal matrix of positive eigenvalues s_i^2 .

Matrix $R = USU^T = \Sigma^{-1/2}$ is the square root of $X^T X$: $R^2 = X^T X$.

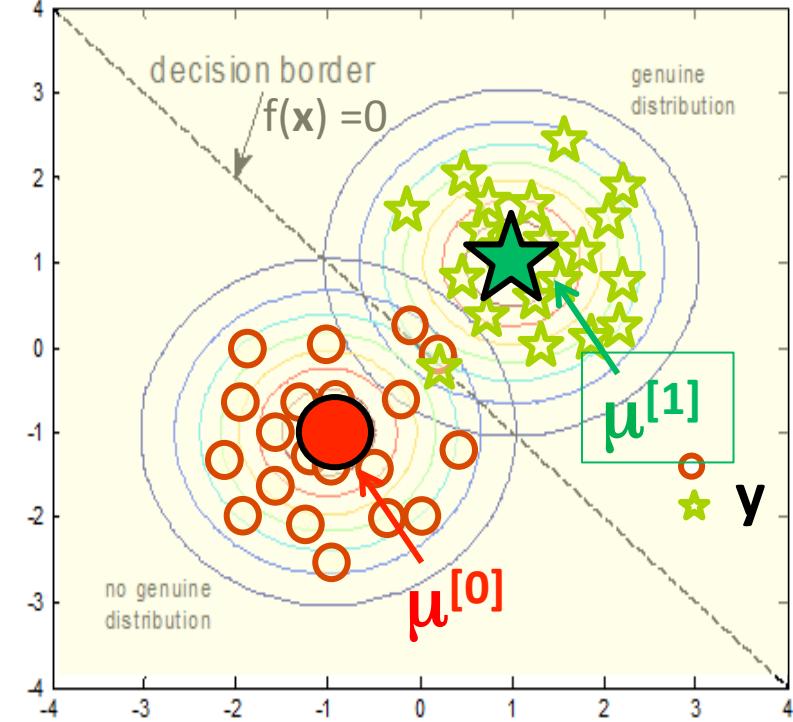
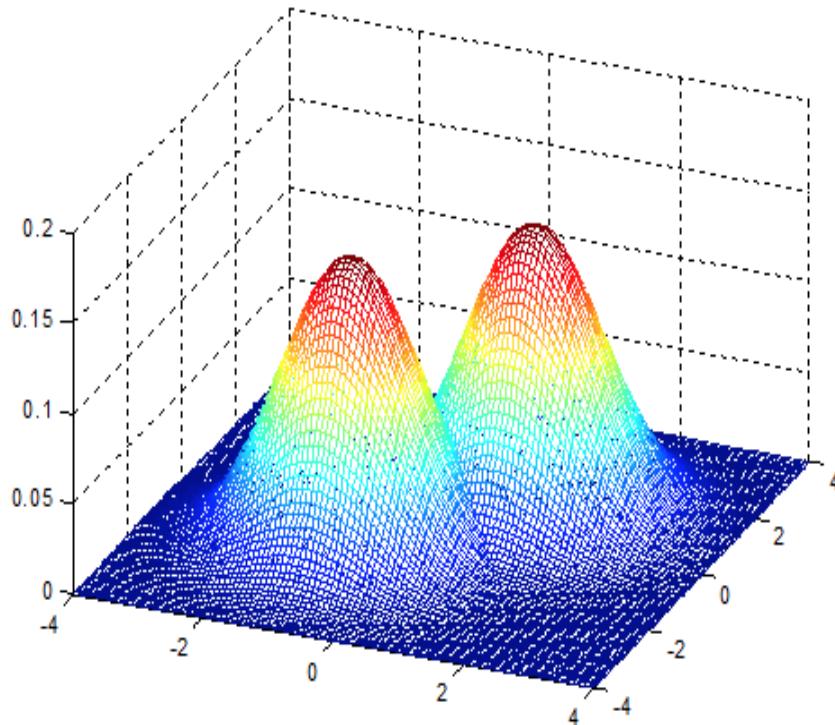
- So to obtain a diagonal covariance matrix S , one just needs to apply matrix U to rotate the axes (and get $\Phi^T \Phi = S^2$, that is a diagonal covariance matrix):
$$\Phi = XU$$
or apply R^{-1} for data **whitening** (and get $\Phi^T \Phi = I$, that is a diagonal covariance matrix, with all the feature variances equal to 1):
$$\Phi = X R^{-1} = X \Sigma^{-1/2}$$

Gaussian classifier

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

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

Source: <http://www.intechopen.com/source/html/17742/media/image25.png>



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

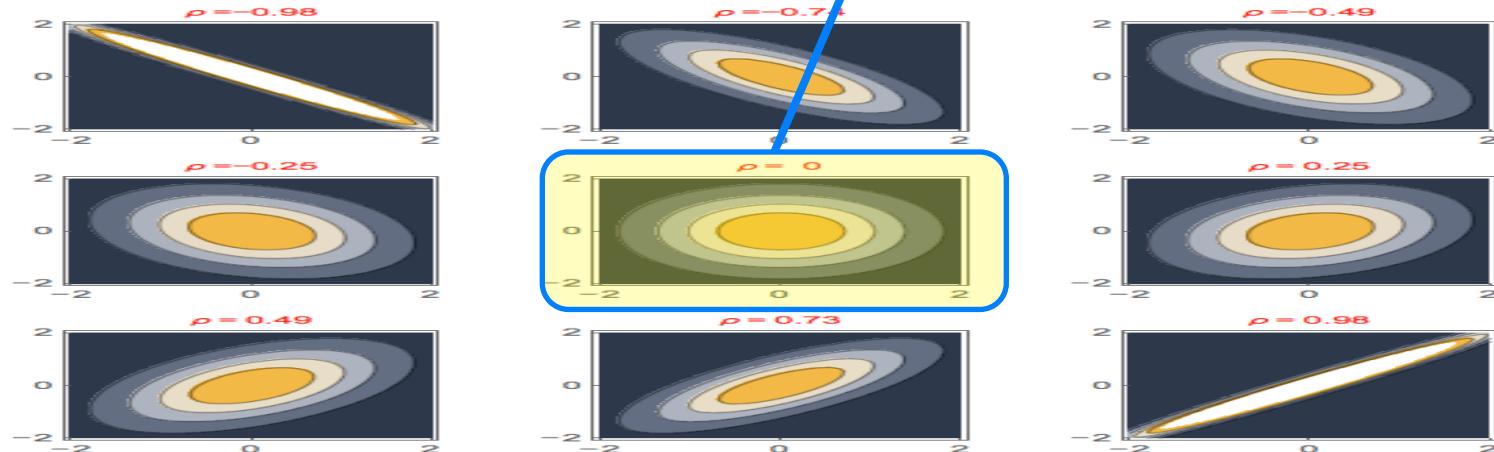
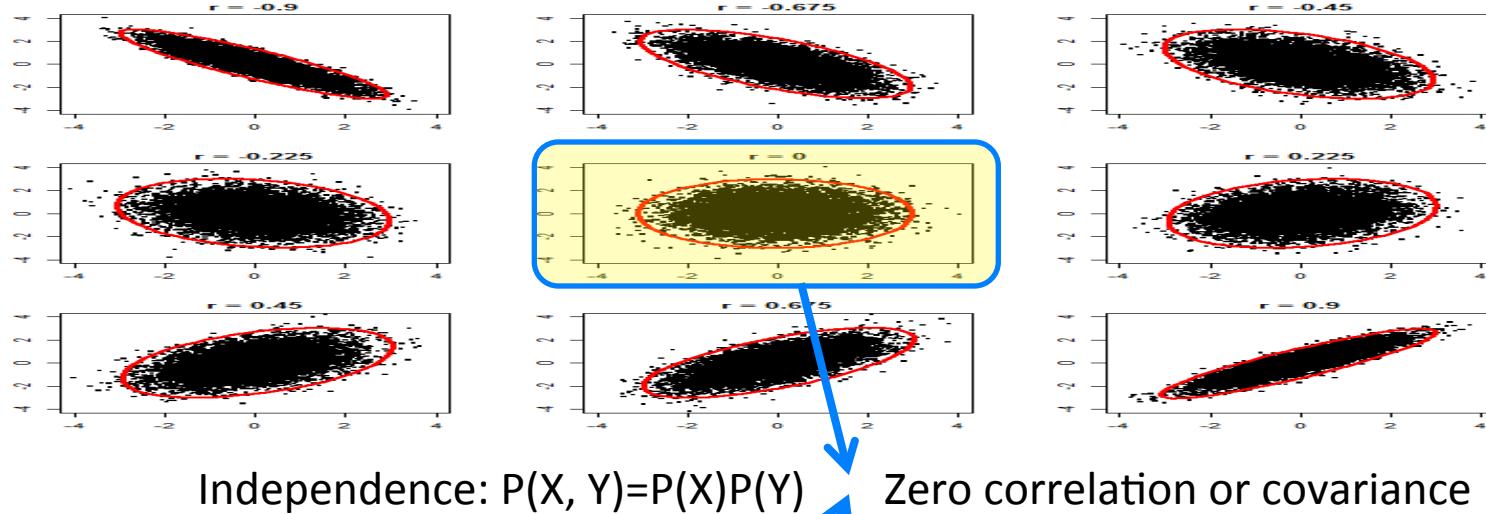
$$\underbrace{\phantom{(\mu^{[1]} - \mu^{[0]}) \cdot x}}_w$$

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

Multivariate Gaussian and covariance

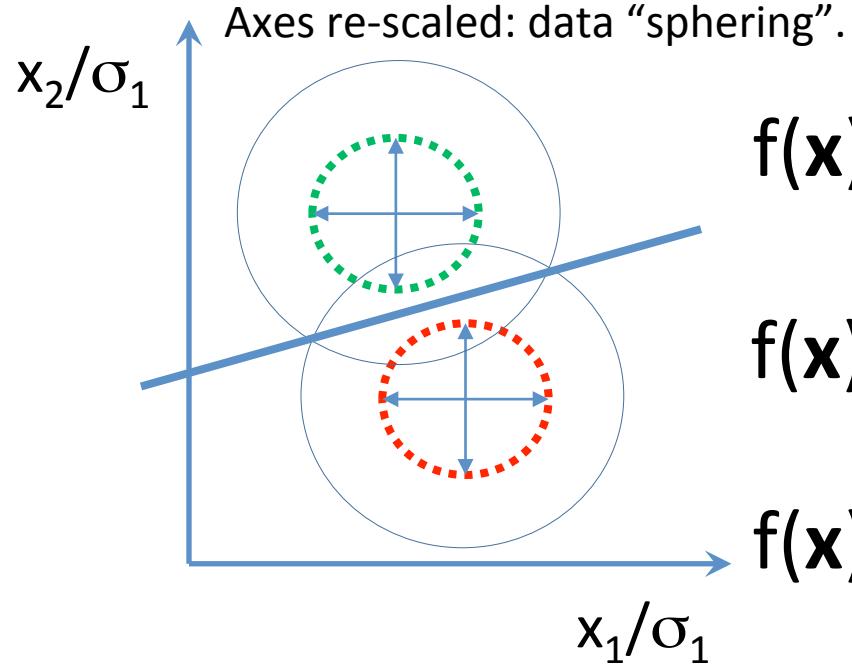
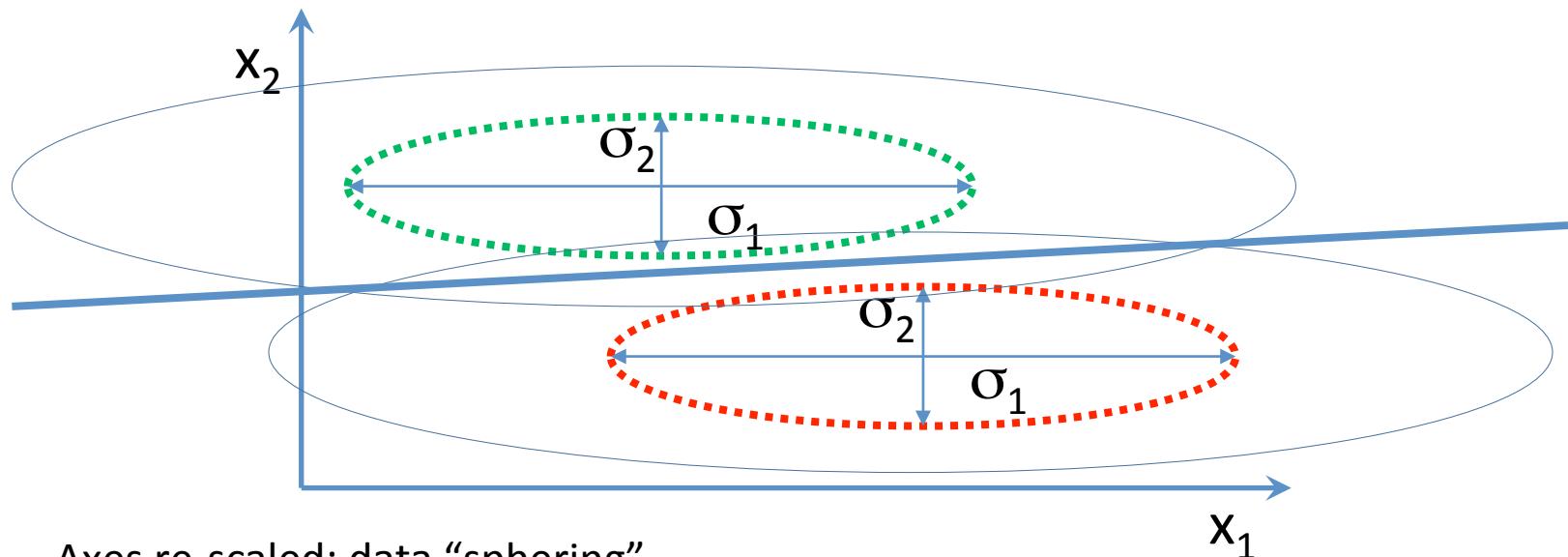
$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\text{cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$



► Fig. 12: Contour plots of the bivariate Normal pdf, for different values of ρ

Different σ_i

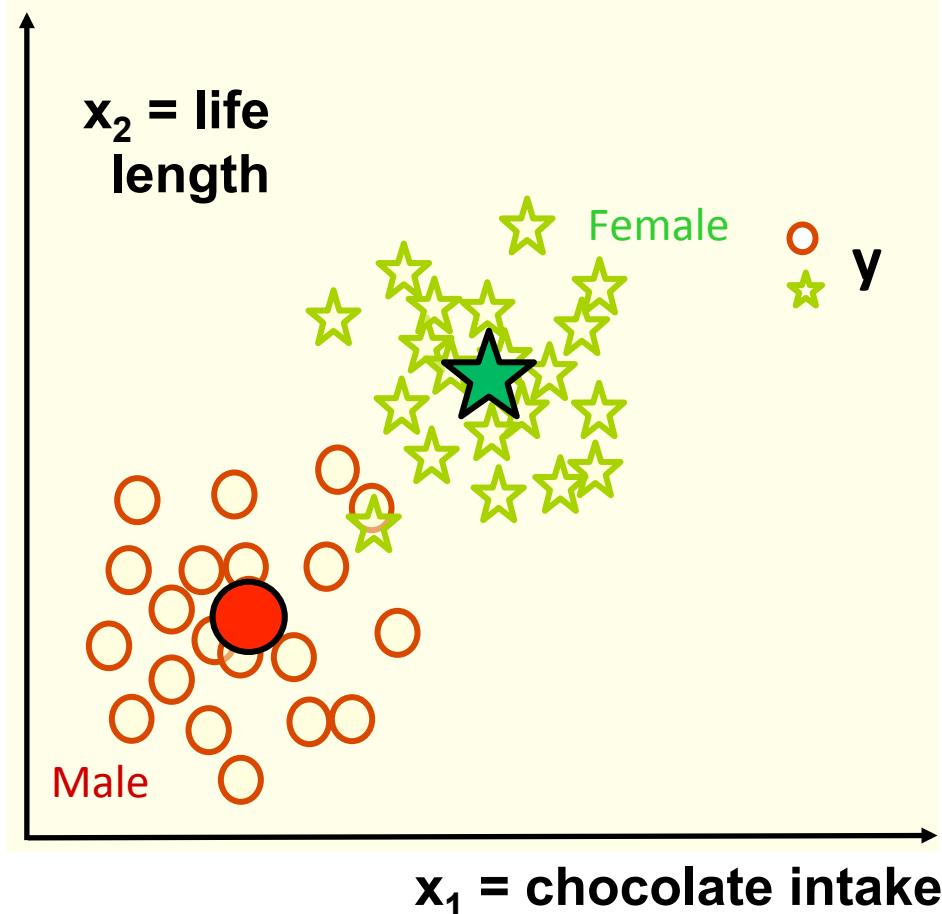


$$f(\mathbf{x}) = \sum_i [(\mu_i^{[1]} - \mu_i^{[0]})/\sigma_i^2] x_i + b$$

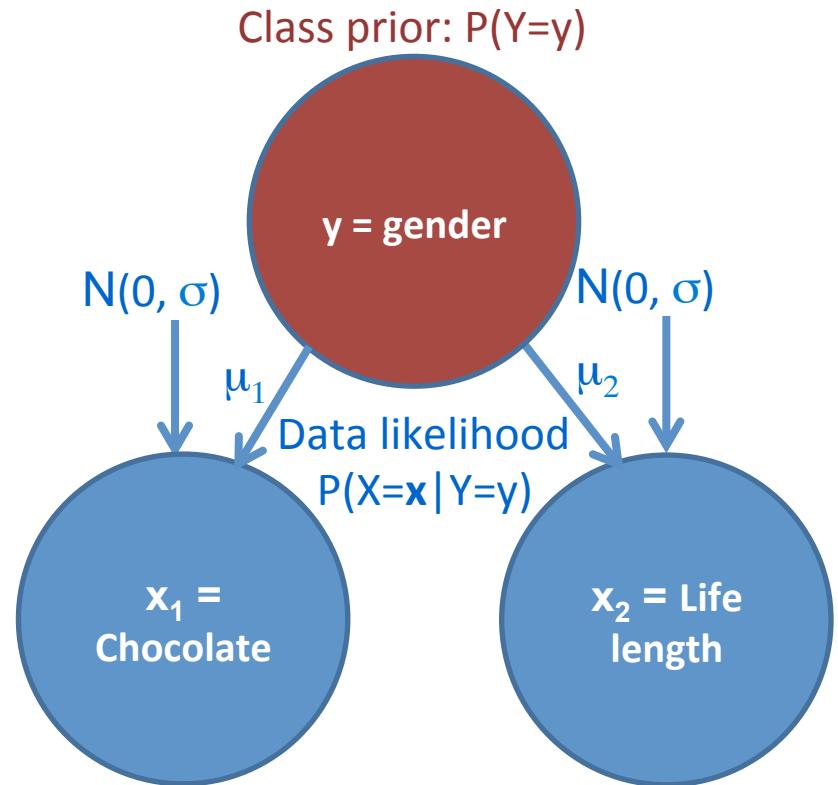
$$f(\mathbf{x}) = \sum_i [\mu_i^{[1]}/\sigma_i - \mu_i^{[0]}/\sigma_i] x_i/\sigma_i + b$$

$$f(\mathbf{x}) = \sum_i [\mu_i^{[1]'} - \mu_i^{[0]'}] x_i' + b$$

Class conditional independence

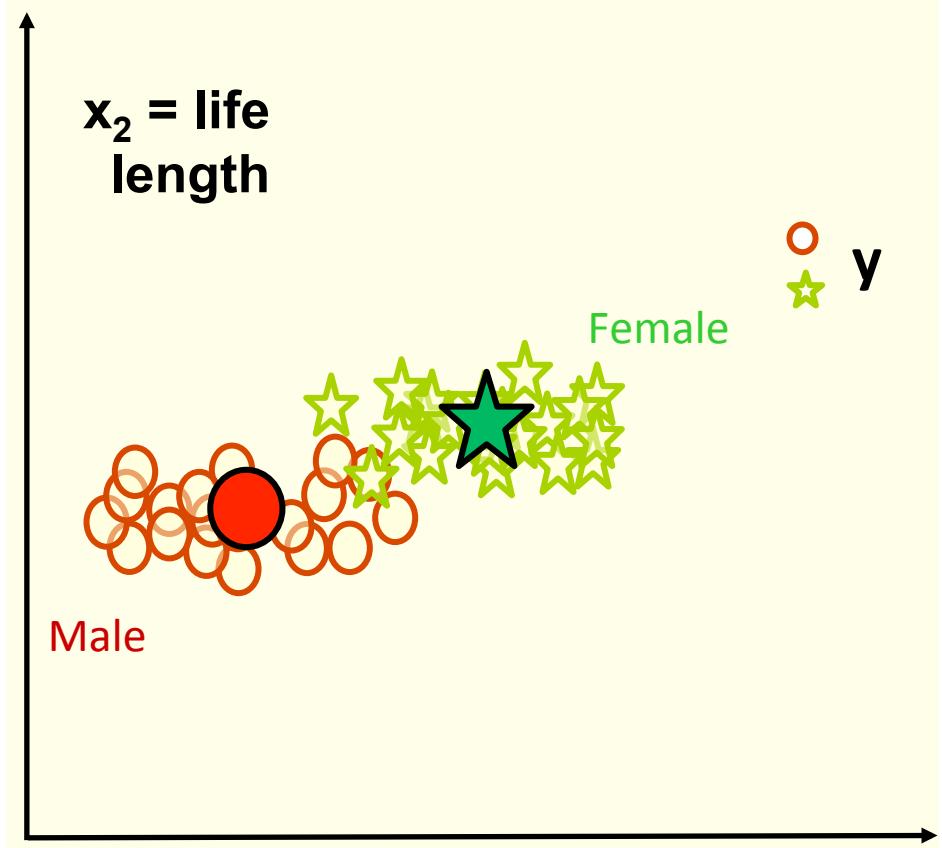


$$P(Y=y \mid X=x) \sim P(Y=y) P(X=x \mid Y=y)$$

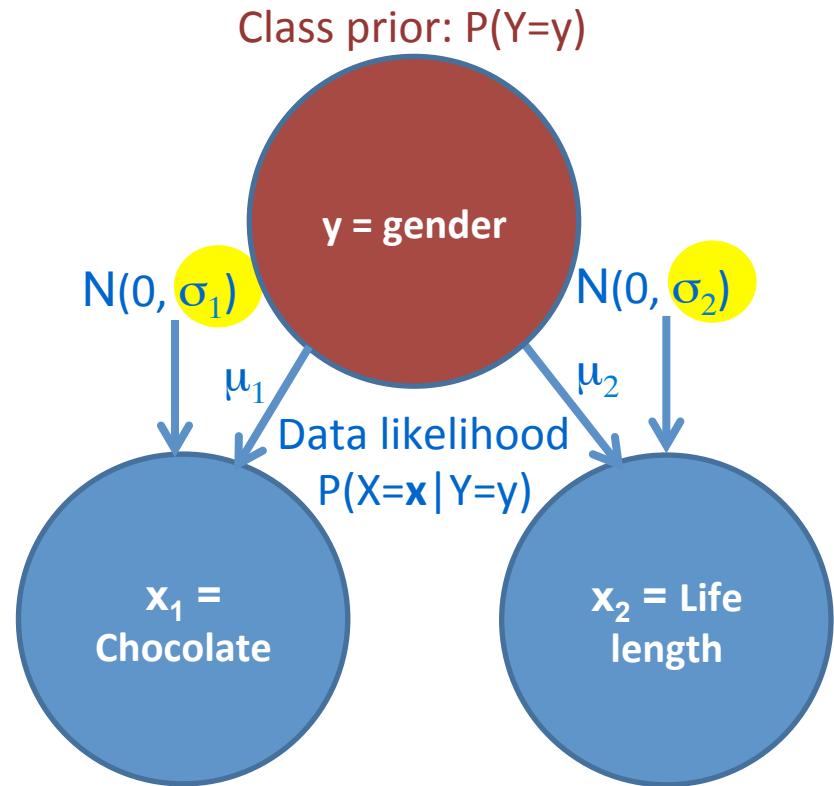


$$\begin{aligned} P(X=x \mid Y=y) &\sim \exp(-\|x-\mu\|^2/2\sigma^2) \\ &= \exp(-\sum_{i=1:N} (x_i - \mu_i)^2 / 2\sigma^2) \\ &= \prod_{i=1:N} \exp(-(x_i - \mu_i)^2 / 2\sigma^2) \\ &\sim \prod_{i=1:N} P(X_i \mid Y) \end{aligned}$$

Class conditional independence

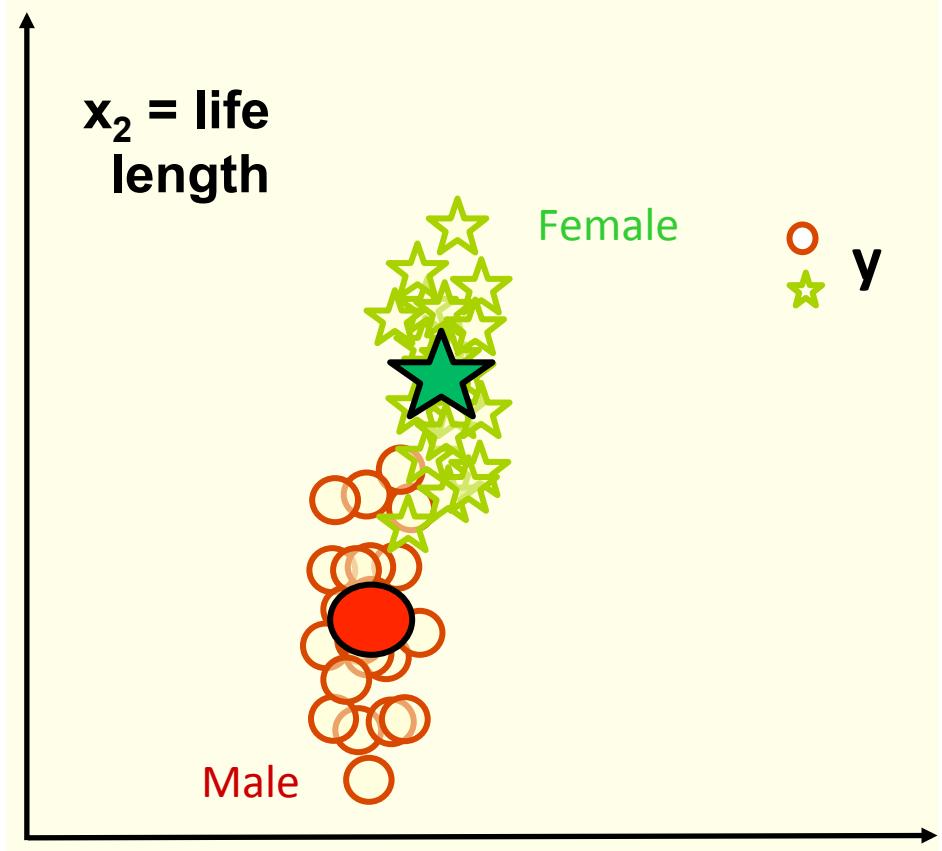


$$P(Y=y \mid X=x) \sim P(Y=y) P(X=x \mid Y=y)$$



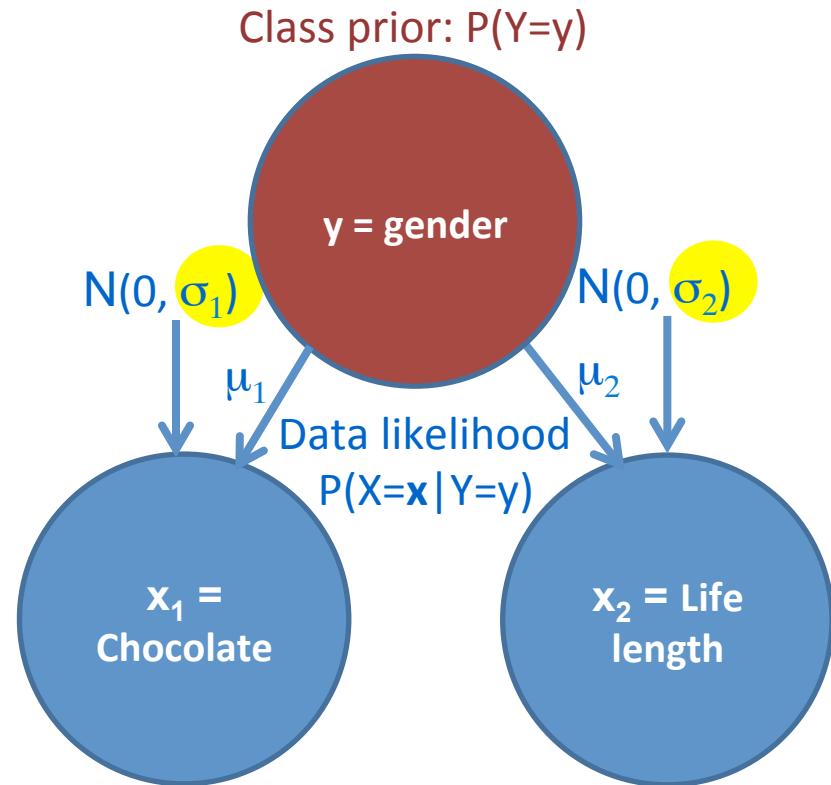
$$\begin{aligned} P(X=x \mid Y=y) &\sim \exp\left(-\sum_{i=1:N} (x_i - \mu_i)^2 / 2\sigma_i^2\right) \\ &= \prod_{i=1:N} \exp(-(x_i - \mu)^2 / 2\sigma_i^2) \\ &\sim \prod_{i=1:N} P(X_i \mid Y) \end{aligned}$$

Class conditional independence



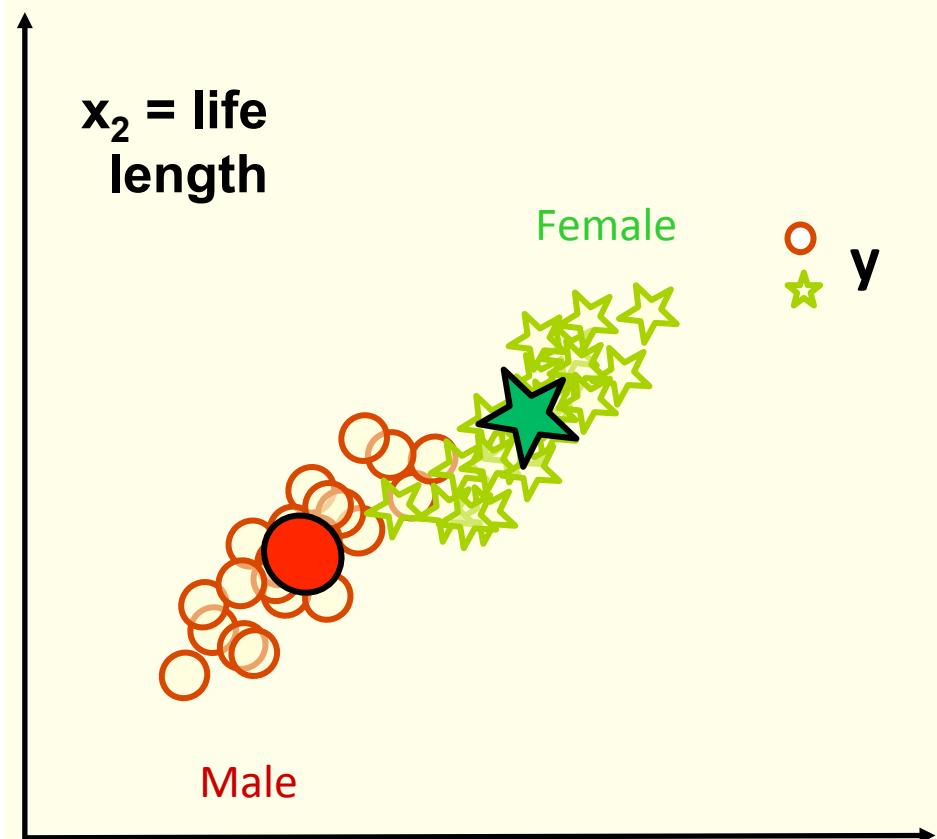
$$x_1 \perp x_2 \mid y$$

$$P(Y=y \mid X=x) \sim P(Y=y) P(X=x \mid Y=y)$$

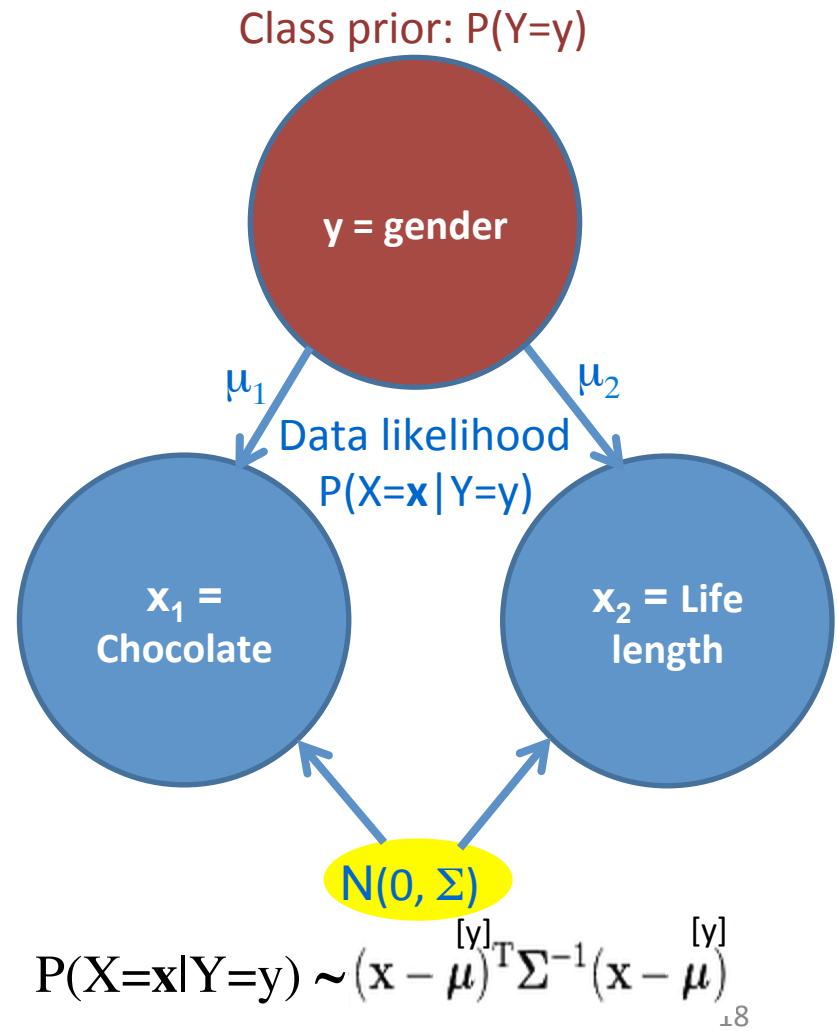


$$\begin{aligned} P(X=x \mid Y=y) &\sim \exp\left(-\sum_{i=1:N} (x_i - \mu_i)^2 / 2\sigma_i^2\right) \\ &= \prod_{i=1:N} \exp(-(x_i - \mu)^2 / 2\sigma_i^2) \\ &\sim \prod_{i=1:N} P(X_i \mid Y) \end{aligned}$$

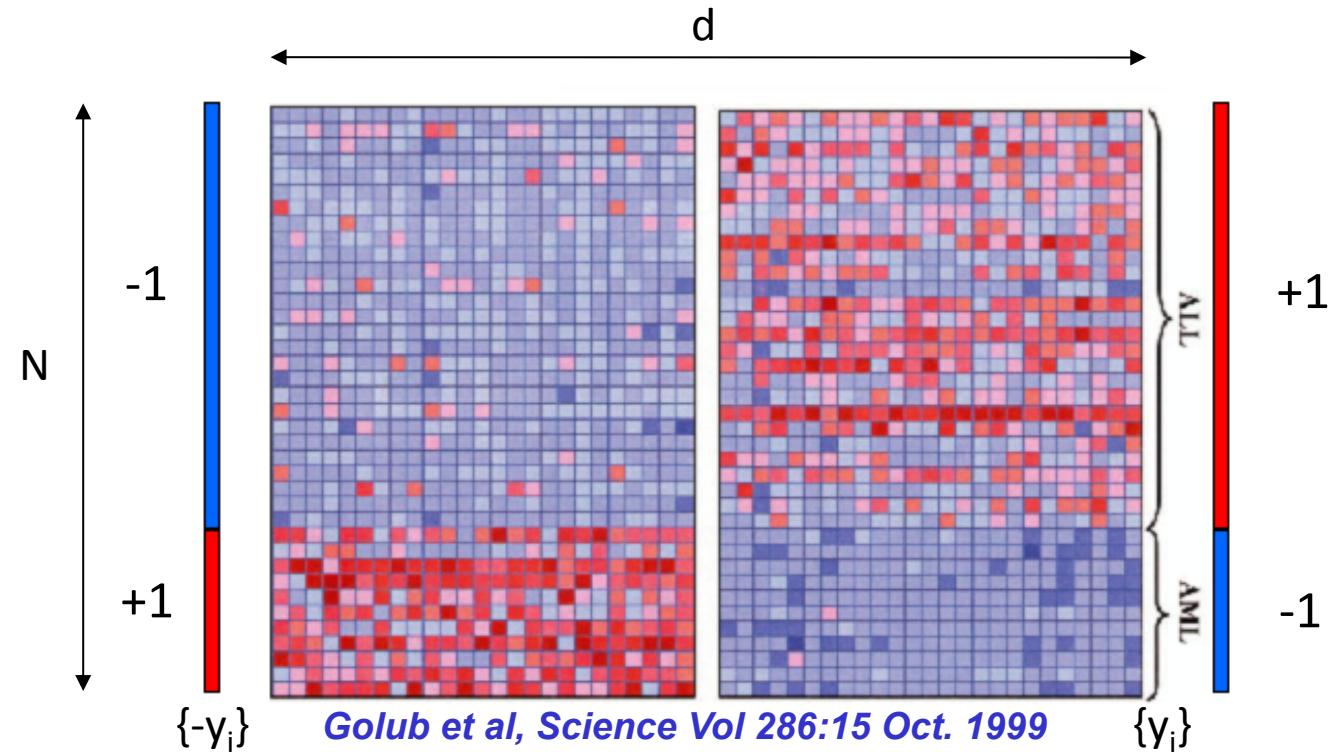
No independence



$$P(Y=y \mid X=x) \sim P(Y=y) P(X=x \mid Y=y)$$



Data sphering

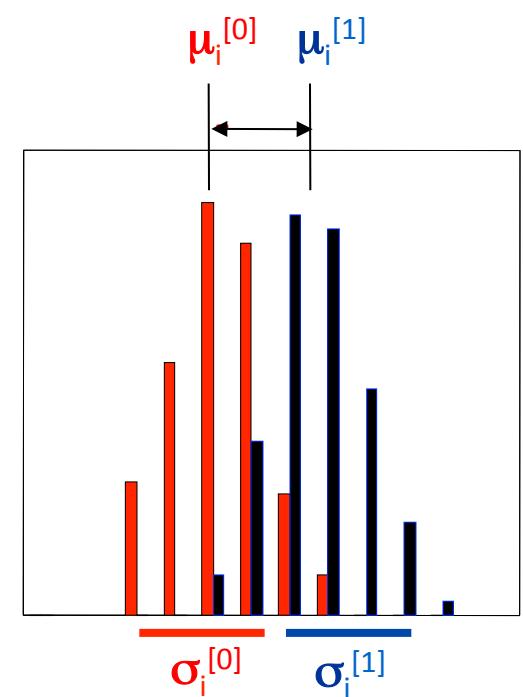


$$w_i = \frac{\mu_i^{[1]} - \mu_i^{[0]}}{\sigma_i^2} \quad \sigma_i \approx \frac{1}{2} (\sigma_i^{[0]} + \sigma_i^{[1]})$$

Rank features according to “signal-to-noise” ratio:

$$S2N = | \sigma_i w_i | \approx | \text{corrcoef} | \sim | x \cdot y |$$

For “balanced classes and after “standardization” $x \leftarrow (x - \mu_x)/\sigma_x$



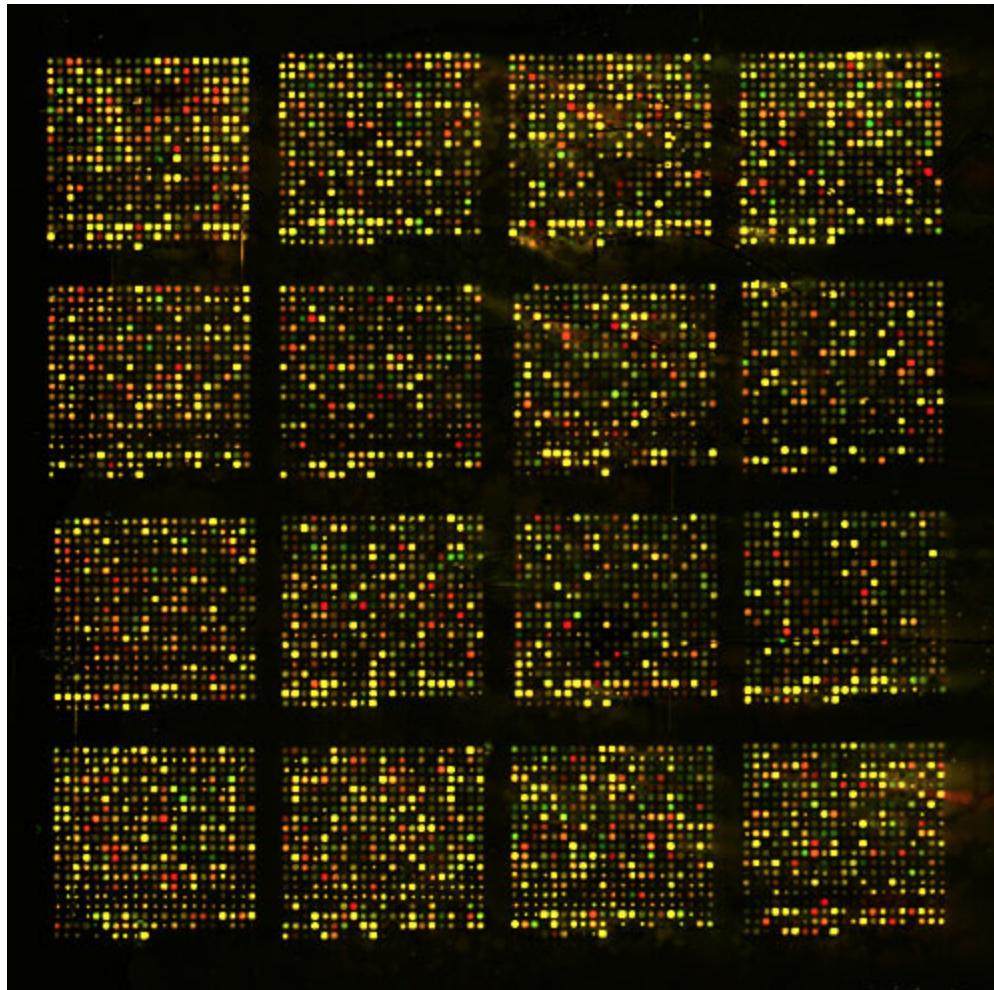
Matrix notations

- For the examples of a given class y , *within class variance*: $\sigma_i^2 = E[(x_i - \mu_i^{[y]})^2]$ (same for all classes).
- Independence assumption: $x_i \perp x_j \mid y \Rightarrow$ no *within class covariance*: $E[(x_i - \mu_i^{[y]})(x_j - \mu_j^{[y]})] = 0$
- Diagonal covariance matrix: $\Sigma = \begin{pmatrix} \sigma_1^2 & & & & 0 \\ & \sigma_2^2 & & & \\ & & \sigma_3^2 & & \\ 0 & & & \ddots & \\ & & & & \sigma_d^2 \end{pmatrix}$
- Trivial to invert, for each class y :

$$(\mathbf{x} - \boldsymbol{\mu})_{(1, d)}^{[y]} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})_{(d, 1)}^{[y] T} = \sum_{i=1:d} (x_i - \mu_i^{[y]})^2 / \sigma_i^2$$

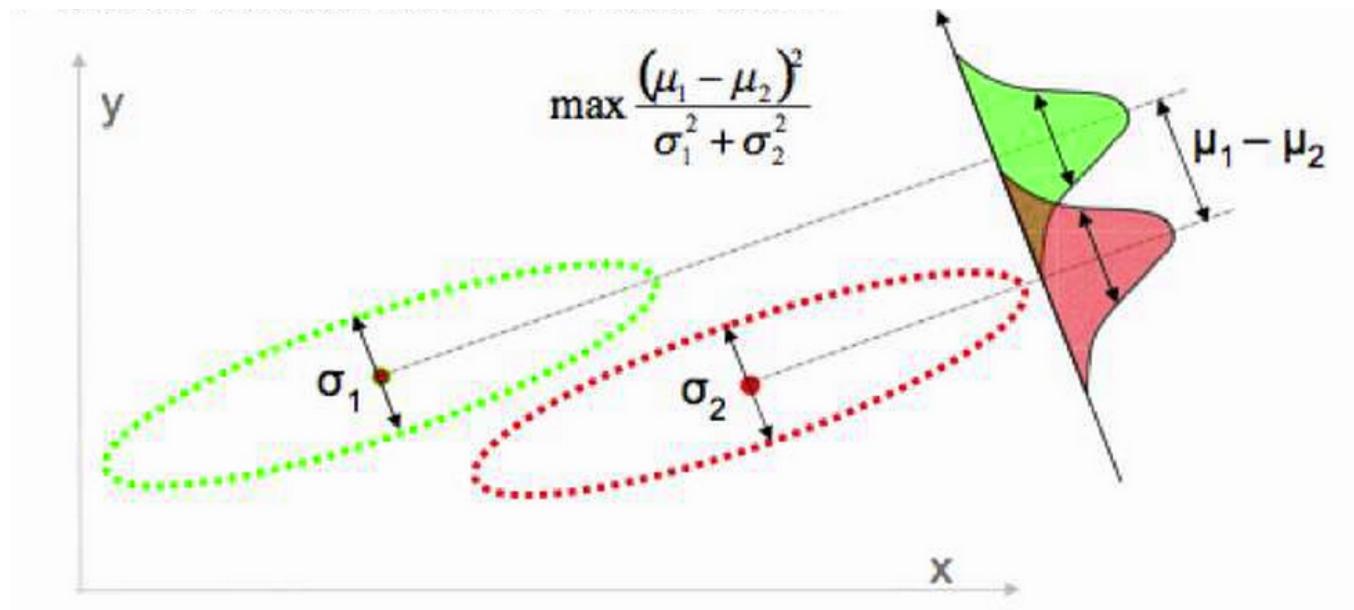
$$P(X=x \mid Y=y) \sim \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{[y]} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})_T^{[y]}\right)$$

Correlated noise

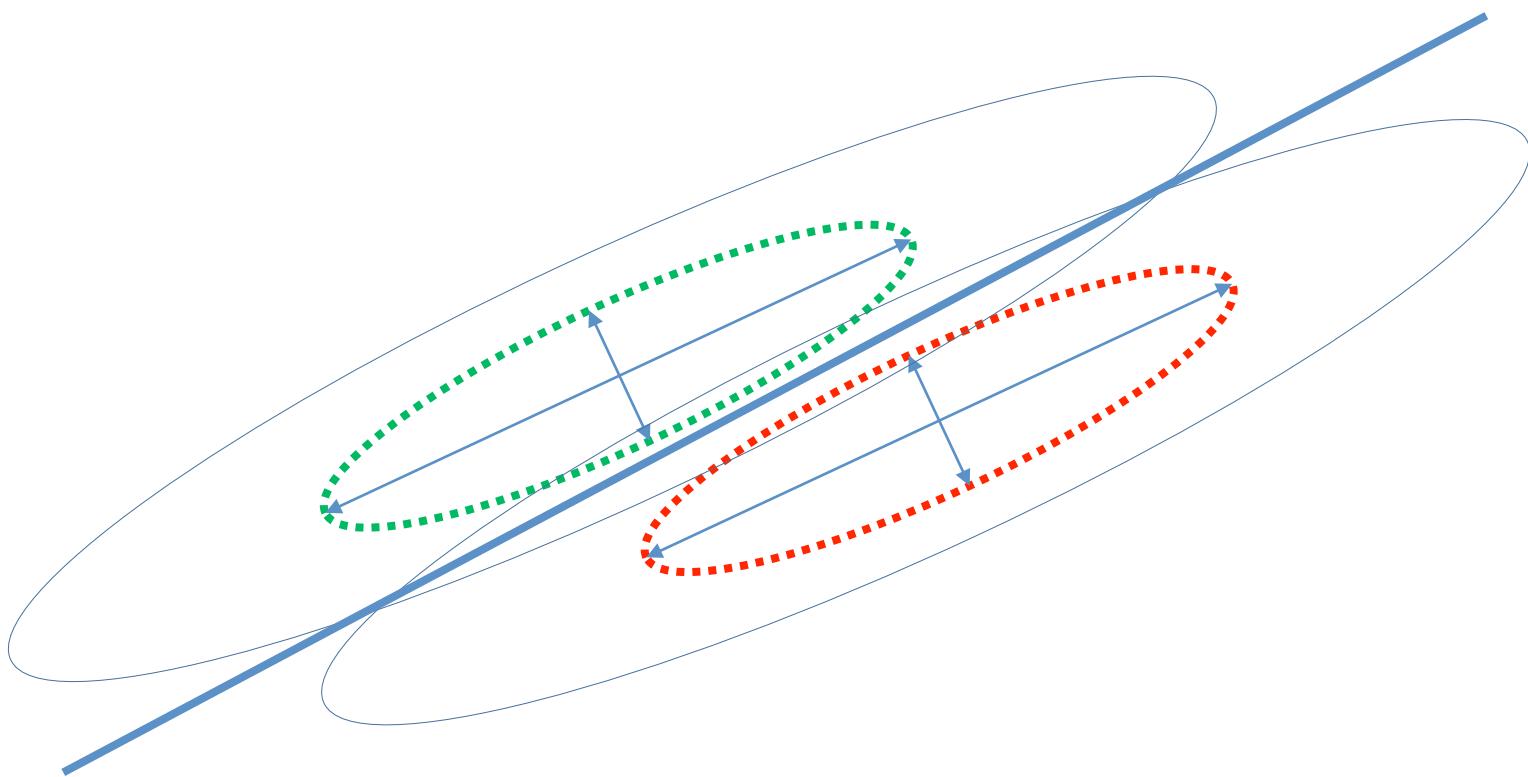


Linear Discriminant Analysis (LDA)

- Rotate input space.
- Find the direction that maximizes the ratio of the between class variance over the within class variance.



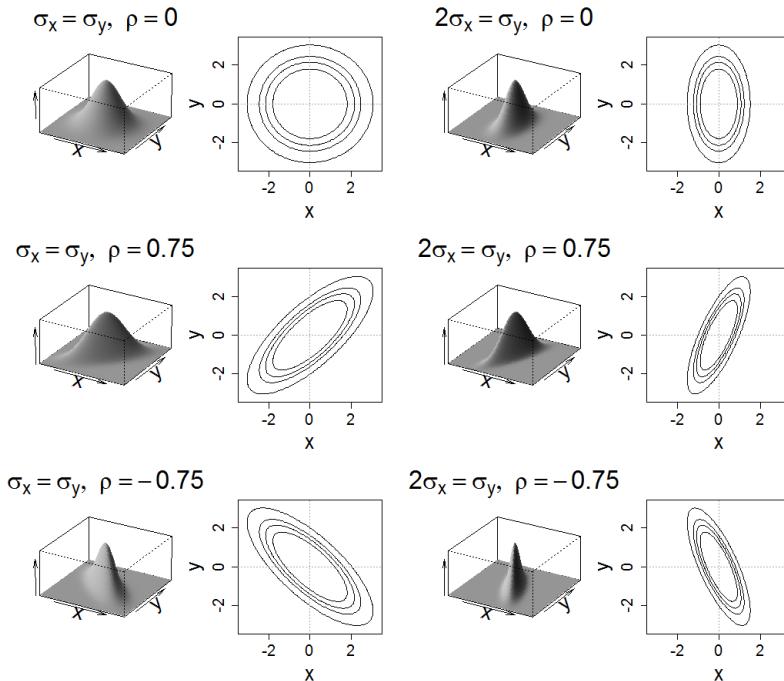
This is still a linear classifier
(assumes identical covariance matrices)



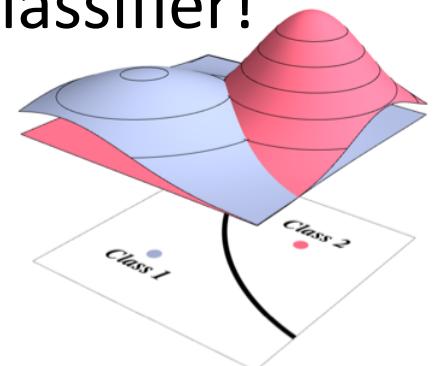
General case

$$P(X=x | Y=y) \sim \exp \left(-\frac{1}{2} (x - \boldsymbol{\mu})_{(1,d)}^{[y]} \Sigma_{(d,d)}^{-1} (x - \boldsymbol{\mu})_{(d,1)}^{[y]T} \right)$$

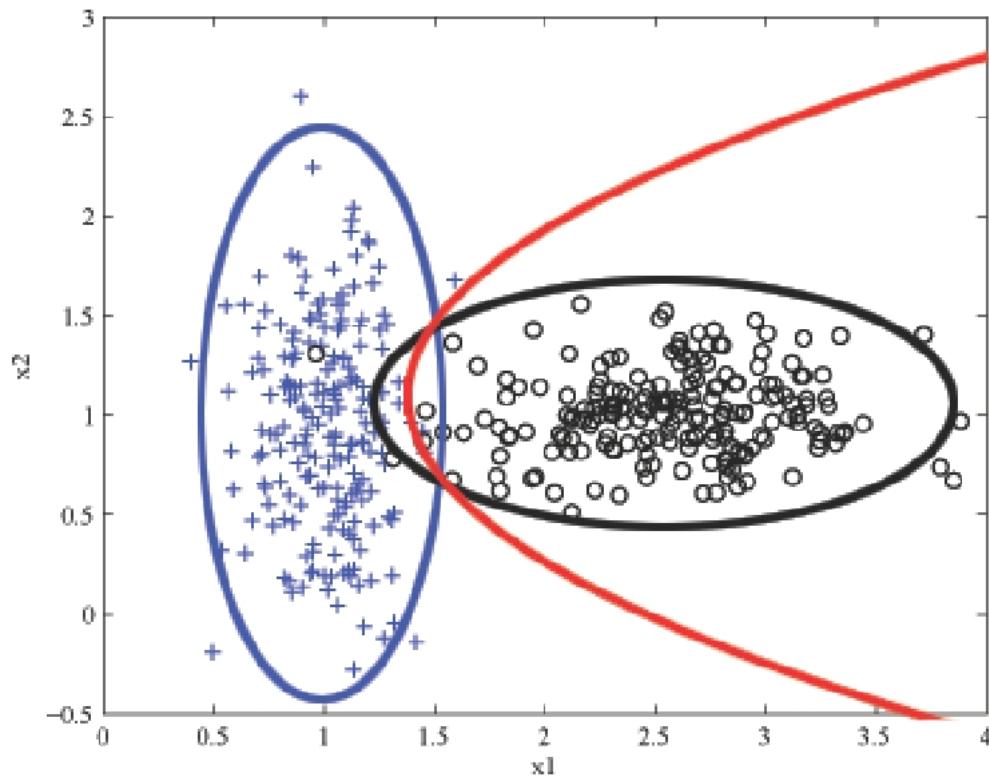
$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} & \cdots & \sigma_{2d} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 & \cdots & \sigma_{3d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \sigma_{d3} & \cdots & \sigma_d^2 \end{pmatrix}$$



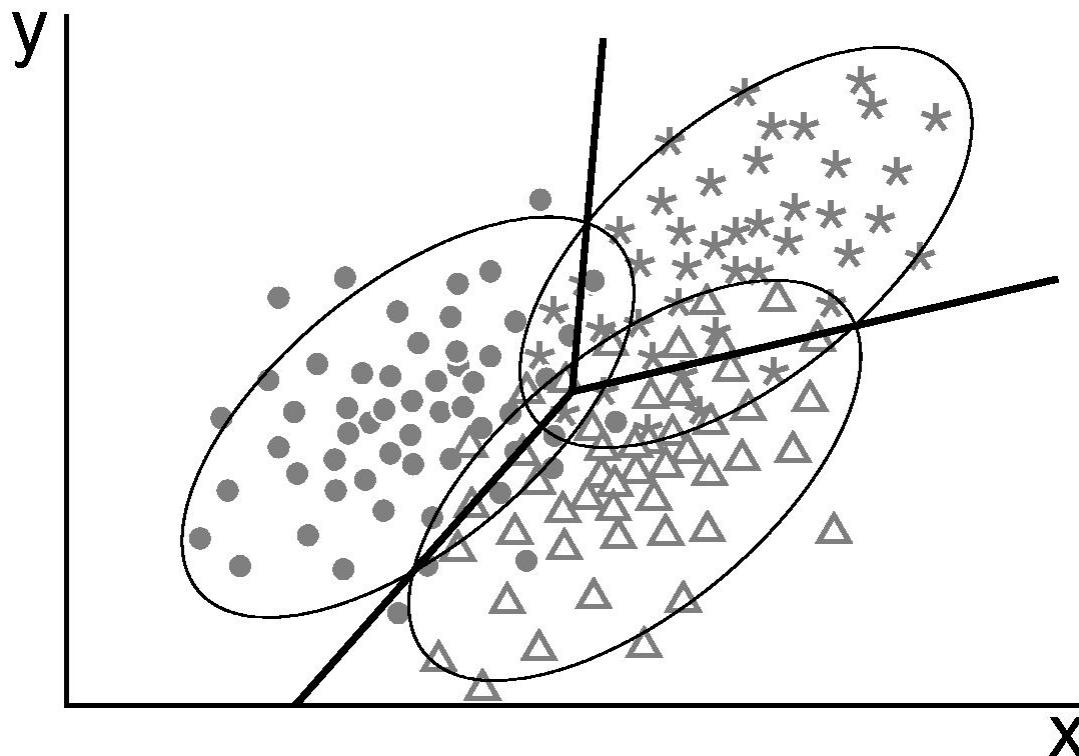
Not necessarily
a linear classifier!



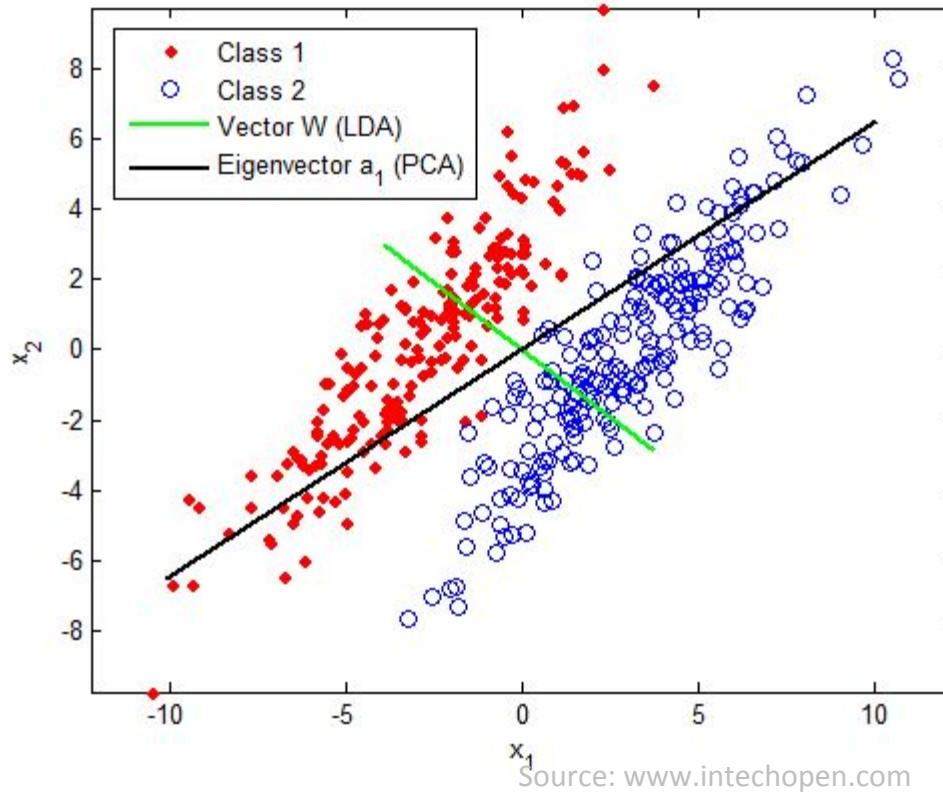
QDA



Multiclass LDA



LDA and PCA



Apply whitening in both case with $\Phi = X \Sigma^{-1/2}$

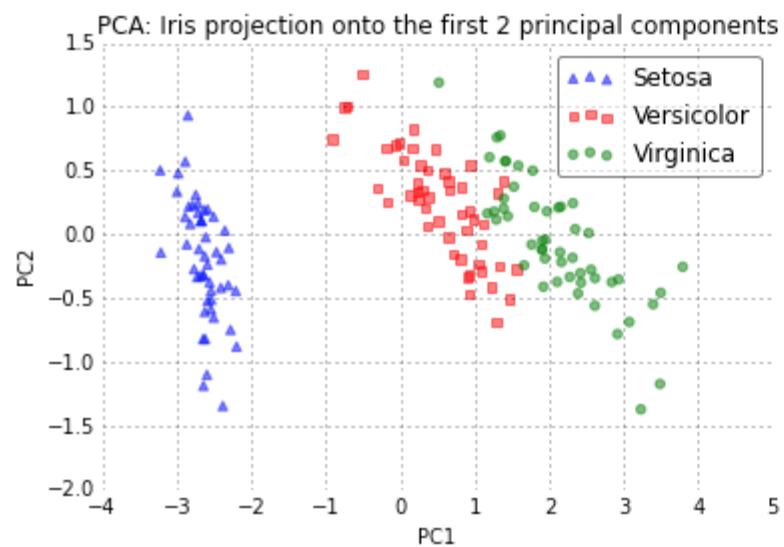
PCA: Σ is the TOTAL covariance matrix.

LDA: Σ is the POOLED WITHIN CLASS covariance matrix.

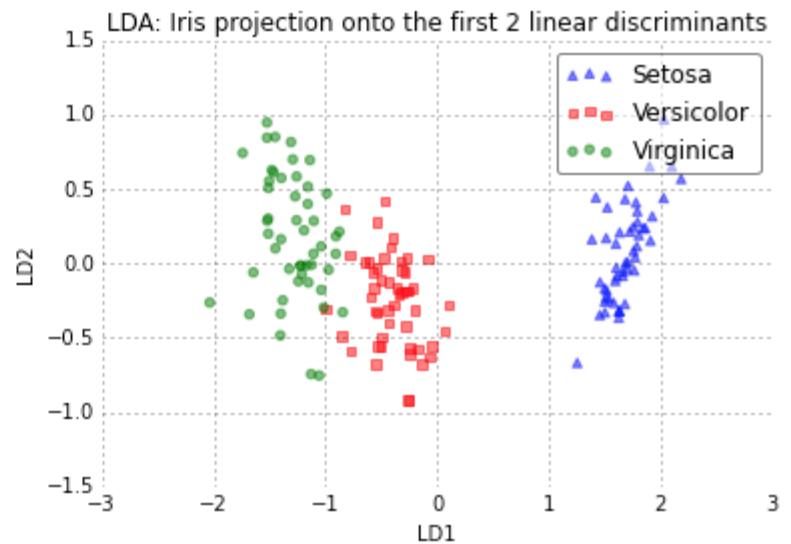
Then use the centroid method.

LDA and PCA (continued)

PCA



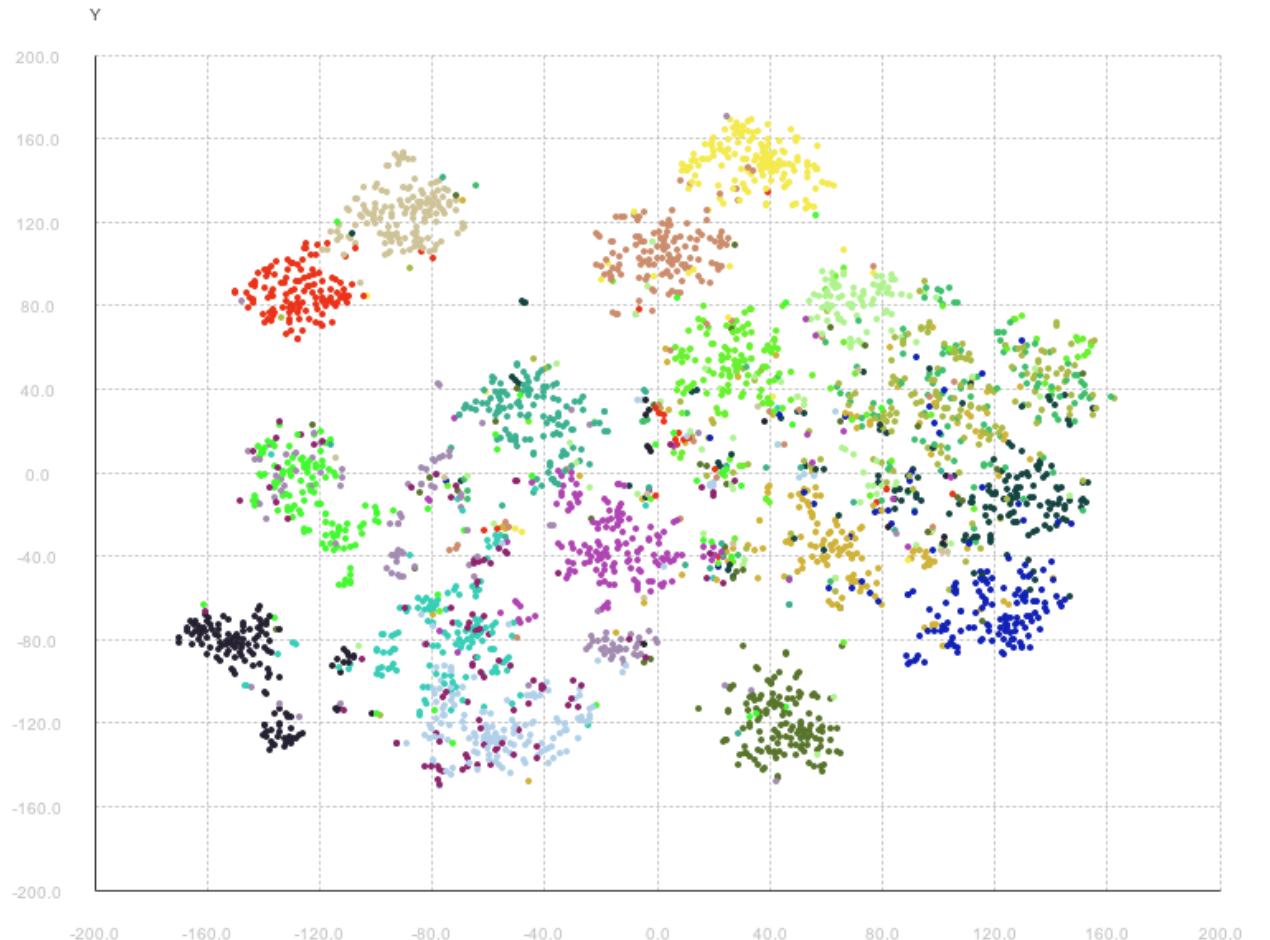
LDA



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

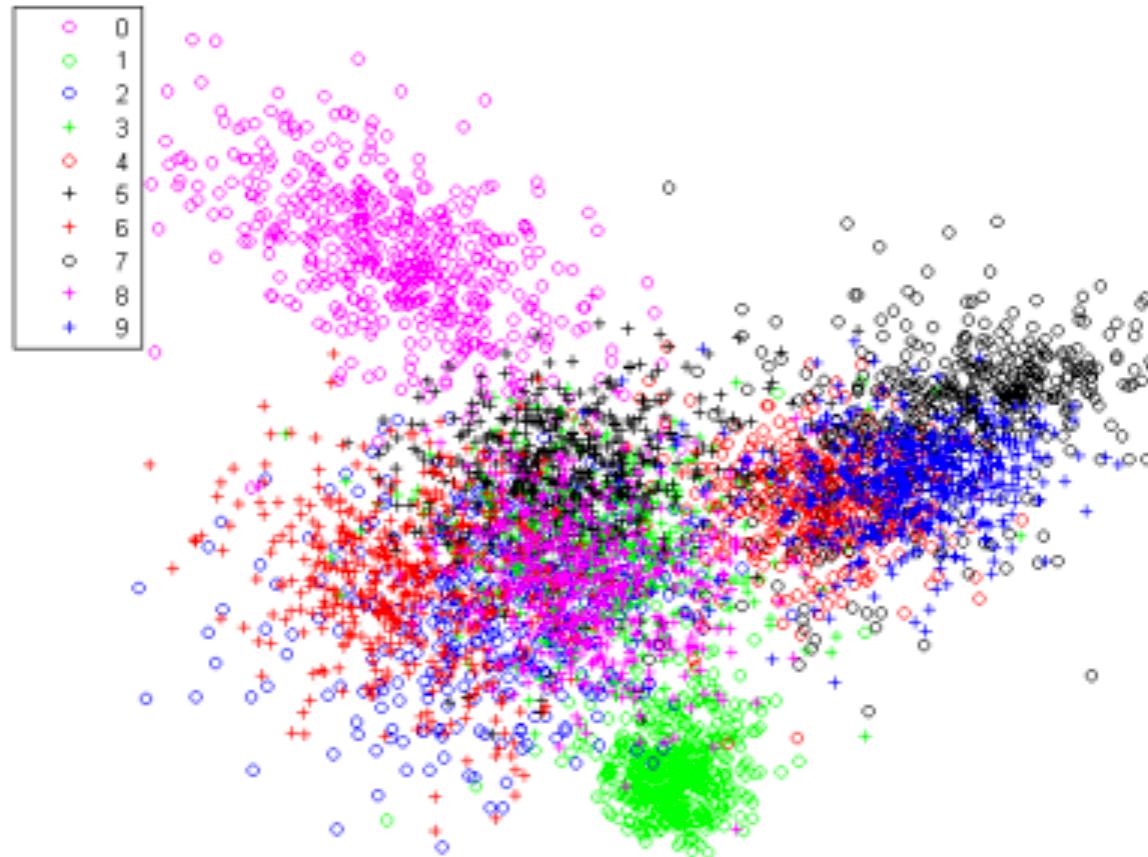
LDA for data visualization

20 Newsgroups



LDA for data visualization

MNIST



Source: Wang et al. 2014 Generalized Autoencoder

Shrinkage (comparison with ridge regression)

- Remember capacity control: Monitor C/N. $C = d_{\text{eff}}$

Model family	m models	Linear (d feat.)	Kernel method	SVM
Capacity	$\leq \log m$	$\leq d+1$	$\leq N$	$\leq N_{\text{sv}}$

- Minimize $\|\mathbf{w}\|^2$: this boils down for regression and SVM to adding a small value $\lambda > 0$ to $\mathbf{X}^T \mathbf{X}$!
- Do the same for LDA before inverting $\Sigma = \mathbf{X}^T \mathbf{X}$.

$$\Sigma \leftarrow \Sigma + \lambda \mathbf{I}$$
- Then whitening: $\Phi = \mathbf{X} \Sigma^{-1/2}$
- Compare with ridge regression: $f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w}^T$

$$\mathbf{w}^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^+ \mathbf{y}$$
- Whitening: $\Phi = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1/2}$ is the whitened data. $\Phi^+ = \Phi^T$ ($\Phi^T \Phi = \mathbf{I}$)

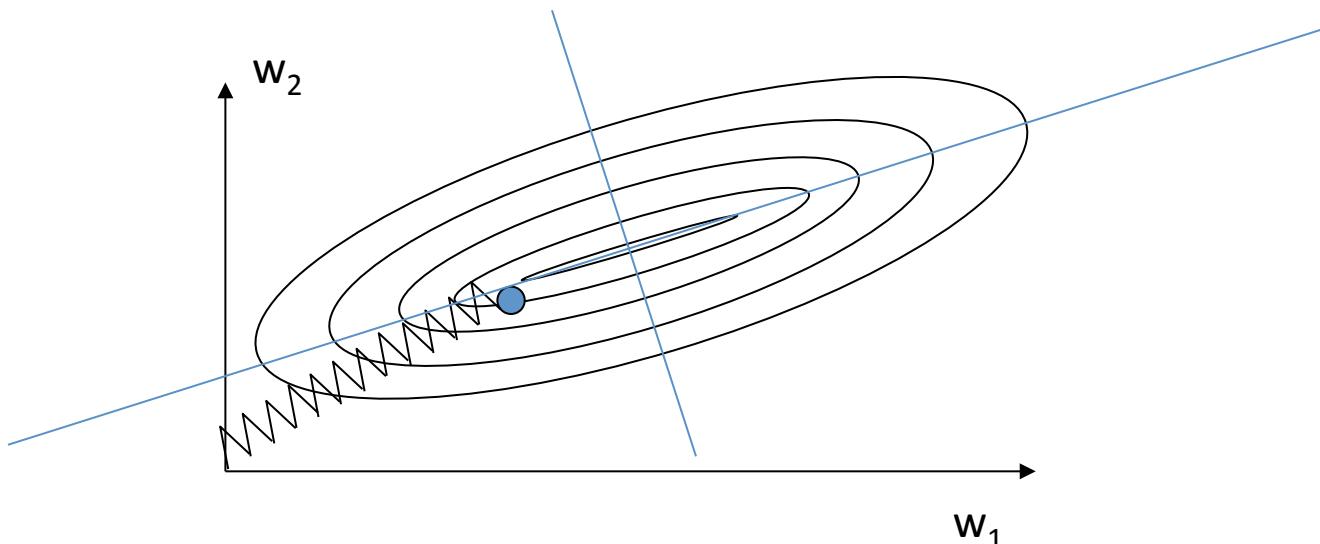
$$\phi = \mathbf{x} (\mathbf{X}^T \mathbf{X})^{-1/2}$$

 In Φ -space, $\mathbf{w}^T = \Phi^T \mathbf{y}$

$$f(\mathbf{x}) = \phi \cdot \mathbf{w}^T$$
- LDA does the same thing, with target values $1/N_1$ and $-1/N_0$ and using the pooled within class covariance matrix.

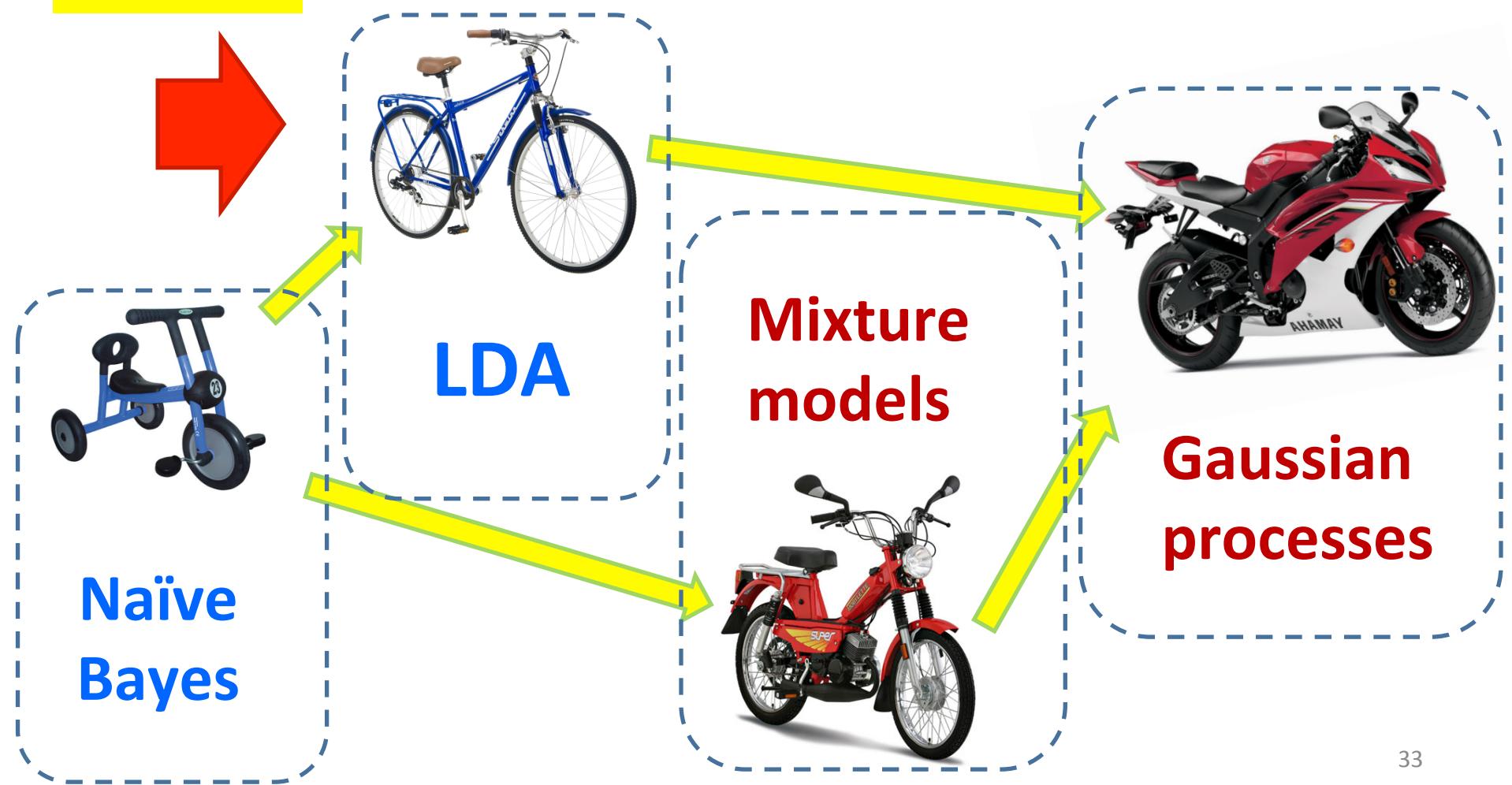
How do ridge regression and regularized LDA control capacity?

- For $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$ the VC dimension is $C = d$ (the dimension of input space).
- $R_{\text{reg}} = \text{RSS} + \lambda \|\mathbf{w}\|^2$ $\text{RSS} = \|\mathbf{X}\mathbf{w}^T - \mathbf{y}\|^2$
 $= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$
- Gradient: $\nabla_{\mathbf{w}} R = 2 (\mathbf{X}^T \mathbf{X} \mathbf{w}^T - \mathbf{X}^T \mathbf{y})$
- Hessian: $H = 2\mathbf{X}^T \mathbf{X}$



You are
here

Generative models



Naïve
Bayes

LDA

Mixture
models

Gaussian
processes

Summary

- LDA is a generalization of the Gaussian classifier for cases in which the input variables are not statistically independent, but all classes have the **same covariance matrix Σ** .
- Once we rotate input space into the “principal axes” of Σ and rescale by the eigen values, LDA is like the isotropic Gaussian classifier, a.k.a. **centroid method**.
- PCA and ridge regression use the covariance matrix of all the data. LDA uses the **“pooled” within class covariance**.
- Two-class LDA is also called Fisher’s linear discriminant and is similar to least square regression with $1/N_1$ and $-1/N_0$ target values.
- LDA is useful for multi-class classification and data visualization. Other classifiers have multi-class versions.
- LDA and logistic regression estimate $P(Y=y | X=x)$ in different ways.

Come to my office hours...
Wed 2:30-4:30 Soda 329

Next time: Mixture models

