

Anisotropic Gaussians & GDA

from last lecture: PDF $f(x) = n(q(x))$, $n(q) = \frac{1}{\sqrt{2\pi d} |\Sigma|} e^{-q/2}$, $q(x) = (x-\mu)^T \Sigma^{-1} (x-\mu)$

Cov. matrix

$$\text{eigendecomp: } \Sigma = V \Lambda V^T$$

Λ

evals of Σ are variances σ_i^2

each variance associated w/ an eigenvector $\Lambda_{ii} = \sigma_i^2$

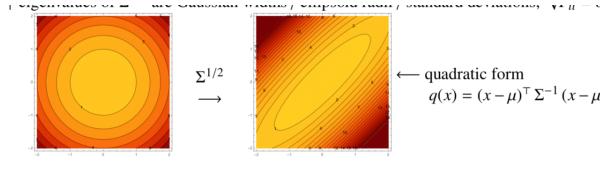
$$\rightarrow \Sigma^{1/2} = V \Lambda^{1/2} V^T \text{ maps spheres to ellipsoids}$$

more helpful
in terms of
STDs

(widths of Gaussians in
each i-th dimension)

$$\hookrightarrow \text{also radii of ellipsoids, } \sqrt{\Lambda_{ii}} = \sigma_i$$

Width of Univariate Gauss $\text{Std}, \sqrt{\lambda_1} \rightarrow \sqrt{\Lambda_{11}}$ of Σ



decomp Σ^{-1} : $\Sigma^{-1} = V \Lambda^{-1} V^T$ ← precision matrix/metric tensor
↳ defines isocontours of

Maximum Likelihood Estimation for Anisotropic Gaussians

given $(x_1, y_1), \dots, (x_n, y_n)$, find best fit normal dist. for each class

* x_i is a column vector

$$\text{For GDA: } \hat{\Sigma}_c = \frac{1}{n_c} \sum_{i: y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

mean rank 1 outer prod., $d \times d$

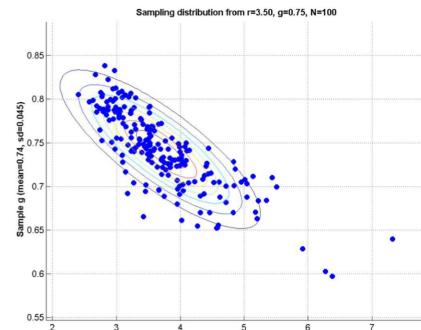
Sample mean

conditional covariance matrix for points in class c

- Prior $\hat{\pi}_c$ that a sample point is in class c

- mean $\hat{\mu}_c$ = avg of all points in c

Gaussian generated
from points



* recall \sum_c is PSD, but not always PD

What if I want to do LDA?

- In LDA, every class has same covariance

$$\text{LDA: } \Sigma = \frac{1}{n} \sum_c \sum_{i:y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

↑ contribution of each sample points within this class

← Pooled-within-class Covariance Matrix

Overall classes

Same variance,
different means



If not within-class, we'd get 1 big contour

QDA: choose class C that maximizes $f(x=x|Y=c)\pi_c$

Same as maximizing the quadratic discriminant function

$$Q_C(x) = \ln((\sqrt{2\pi})^d f_C(x)\pi_c) = -\frac{1}{2}(x - \mu_c)^T \Sigma_c^{-1}(x - \mu_c) - \frac{1}{2}\ln|\Sigma_c| + \ln\pi_c$$

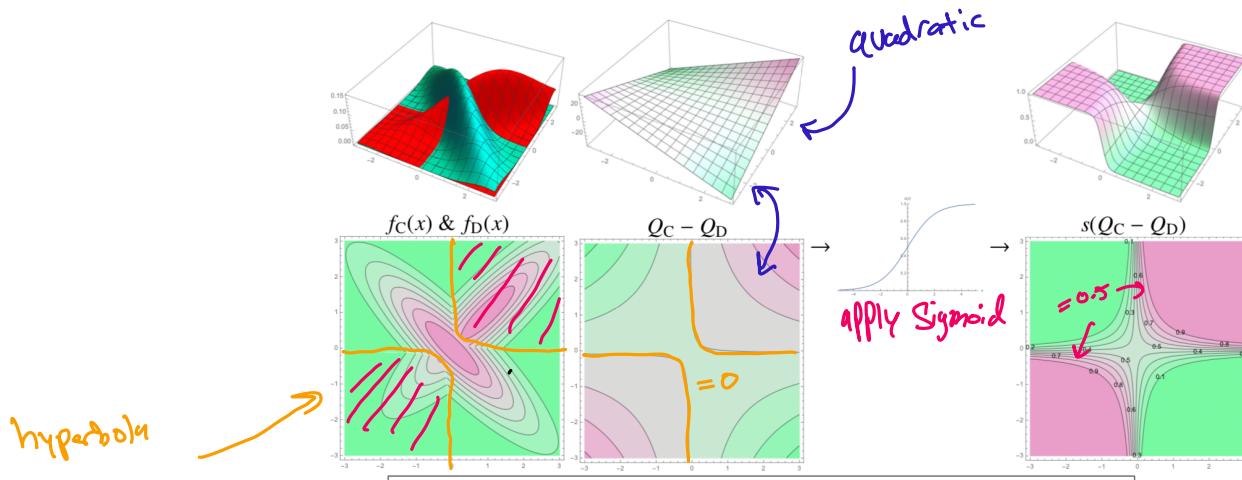
↑ not necessarily isotropic

for 2-classes: Decision fn $Q_C(x) - Q_D(x)$ is quadratic, might be indefinite

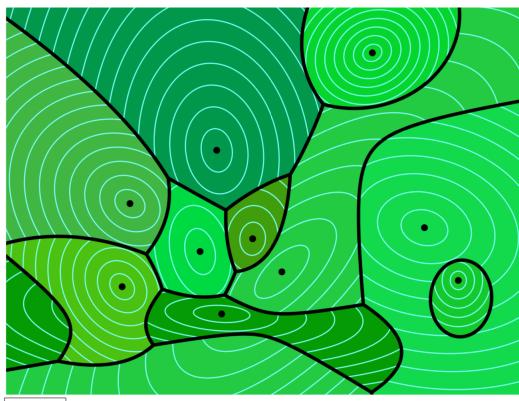
\Rightarrow Bayes decision boundary is quadratic

- Posterior is $P(Y=c|X=x) = s(Q_C(x) - Q_D(x))$, $s(\cdot)$ is logistic fn

Solution to quadratic in an arbitrary space



Anisotropic V. diagram



aniso.pdf [When you have many classes, their QDA decision boundaries form an anisotropic Voronoi diagram. Interestingly, a cell of this diagram might not be connected.]

w/ LDA_j quadratic terms
cancel out
 \downarrow
LDA: one \vec{z} for all classes

$$Q_C(x) - Q_D(x) = (\mu_C - \mu_D)^T \Sigma^{-1} x - \frac{\mu_C^T \Sigma^{-1} \mu_C - \mu_D^T \Sigma^{-1} \mu_D}{2} + \ln \pi_C - \ln \pi_D$$

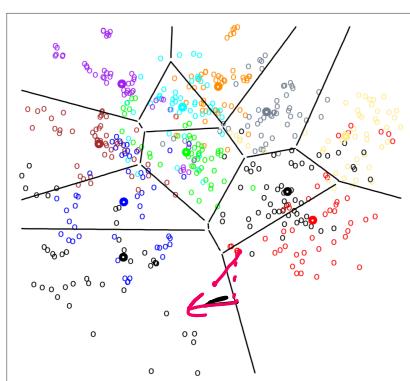
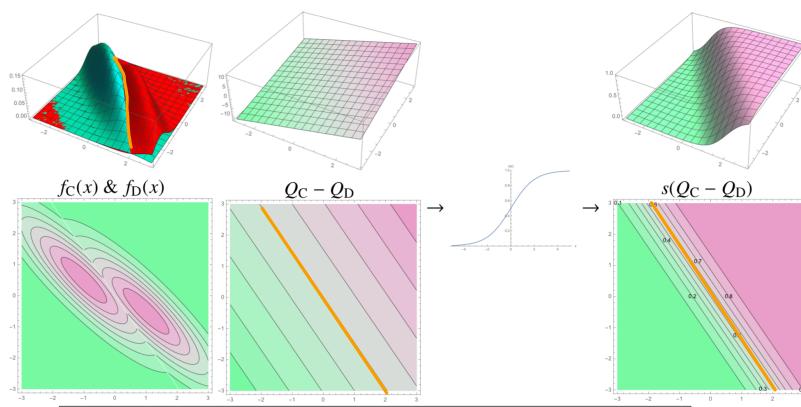
choose class that maximizes linear discriminant fn:

$$\boxed{\mu_C^T \Sigma^{-1} \vec{x}} \quad \boxed{-\frac{1}{2} \mu_C^T \Sigma^{-1} \mu_C + \ln \pi_C}$$

expensive, once for each class

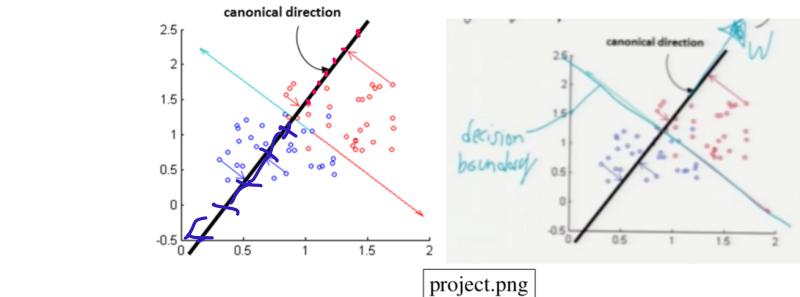
compute only 1 discriminant fn, use on all classes

- Decision boundary is $w^T x + d = 0$
- Posterior is $p(Y=c|X=x) = s(w^T x + d)$



Notes: LDA interpreted as projecting points onto normal vector, cut projection in half \rightarrow decision boundary

- LDA often interpreted as projecting points onto normal w ; cutting the line in half.



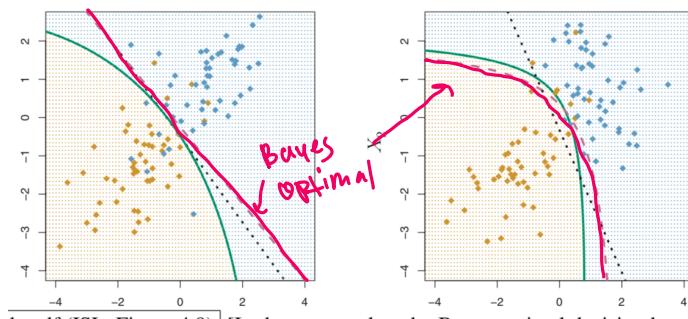
For 2 classes

for 2 classes, -LDA learns $d+1$ params (\vec{w} & λ)

- QDA has $\frac{d(d+3)}{2} + 1$ parameters (cross terms)

- QDA more likely to overfit to data

- LDA underfits



- w/ added features, LDA can be nonlinear

QDA can give boundaries not quadratic

- we don't get true Bayes classifier unless data is truly Gaussian

- Since we estimate Σ, μ from finite data

- most datasets aren't perfectly Gaussian

- changing priors or loss = adding constants to discriminant fns
 \hookrightarrow change isotopes

- posterior prob gives decision boundary for 10%, 50%, 90%.

\hookrightarrow can think choosing eigenvalue = prob. p same as choosing asymmetrical loss p for false positive, $1-p$ for false neg.

OR choosing prior to be $1-p, 1-p$

Some terms over just few weeks:

$X \in \mathbb{R}^{n \times d}$ is a design matrix of data points, each col. is feature

Centering X : Subtract μ from each row of X $\rightarrow \mu$ of each row becomes 0

Let R be uniform. Sample Covariance Matrix: $V_{\text{cov}}(R) = \frac{1}{n} \dot{X}^T \dot{X}$

Class specific cov. matrix: $\hat{\Sigma}_c = \frac{1}{n_c} \underset{\uparrow}{X_c^T X_c}$
only points of design matrix w/c

Decorrelating X : eigendecompose $\hat{\Sigma}$, rotate X until its axis-aligned

$$Z = \dot{X} V, V_{\text{cov}}(R) = V \Lambda V^T$$

$\rightarrow V_{\text{cov}}(Z) = \Lambda$, Z has diag. covariance