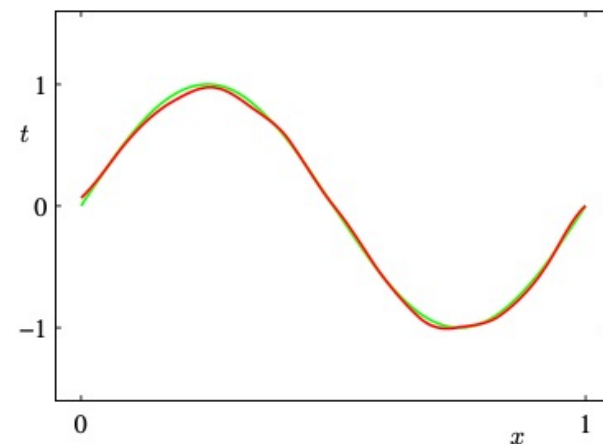
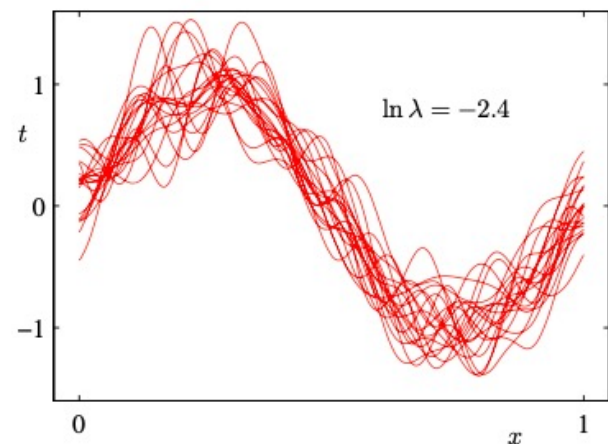
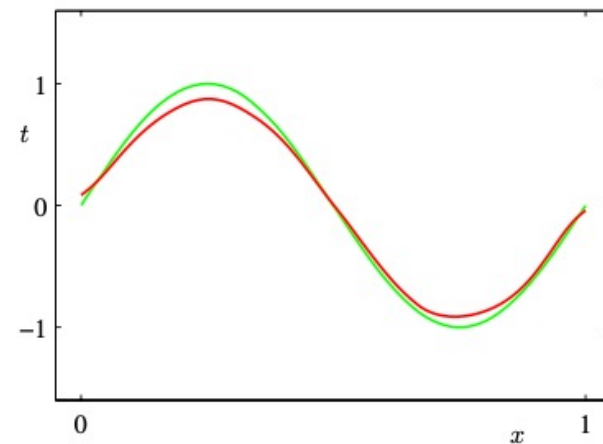
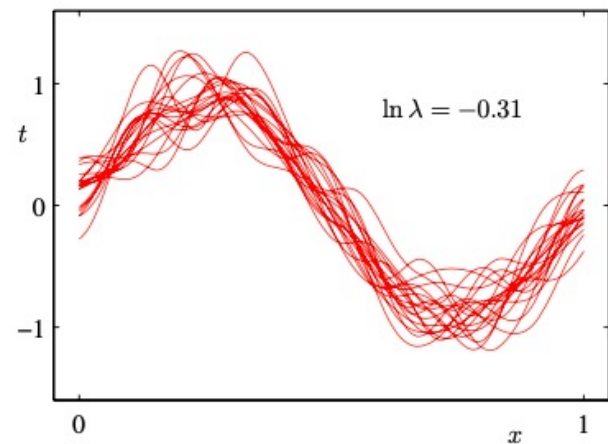
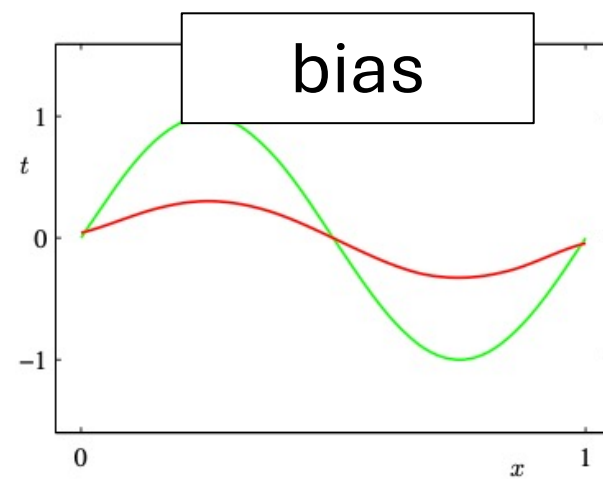
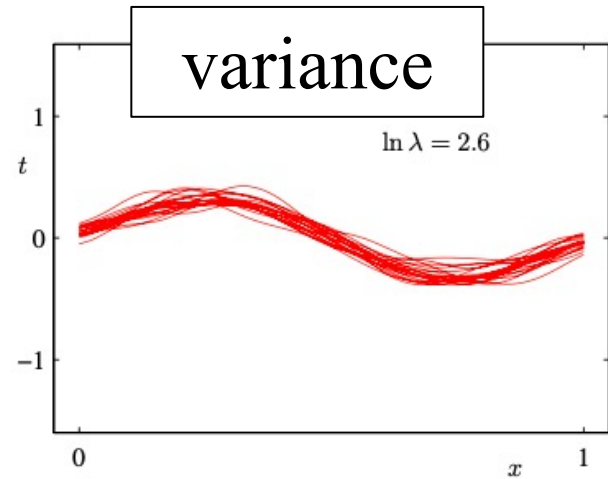


CS 461: Machine Learning Principles

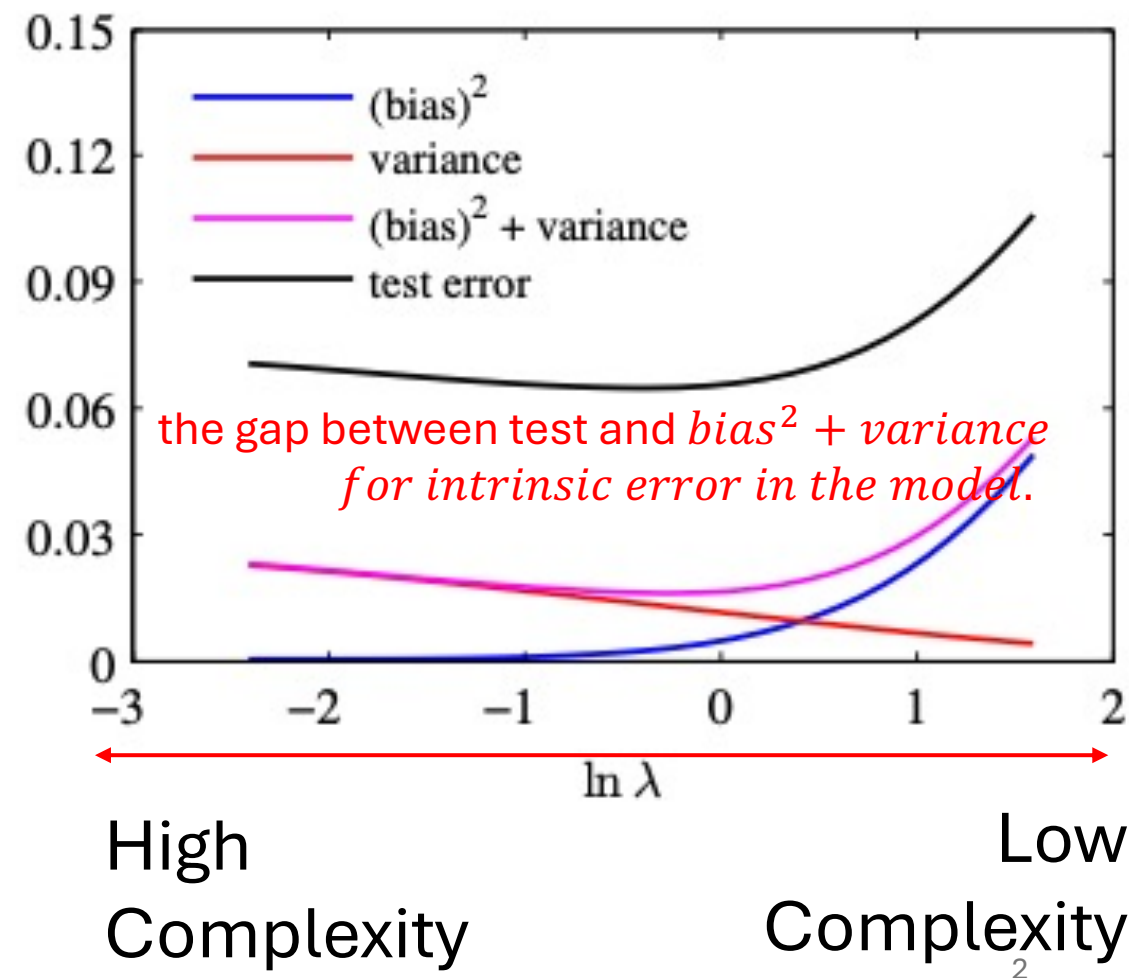
Class 8: Sept. 30

Binary Classification: **Linear Discriminant Analysis (LDA)**

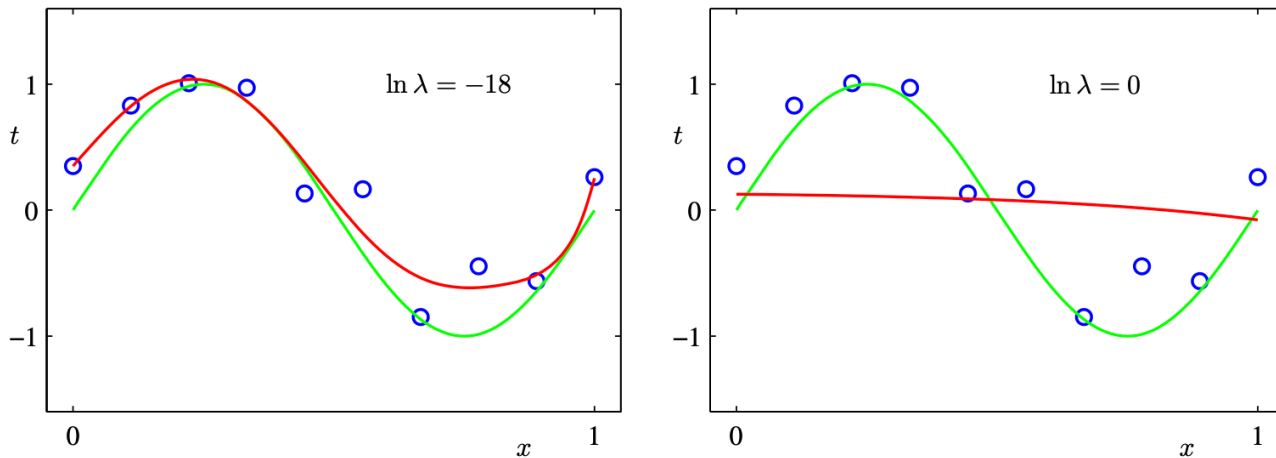
Instructor: Diana Kim



Trade off between
Variance and Bias according to λ
in Ridge Regression (Complexity)



Overfitting can be avoided by adopting Bayesian Approach



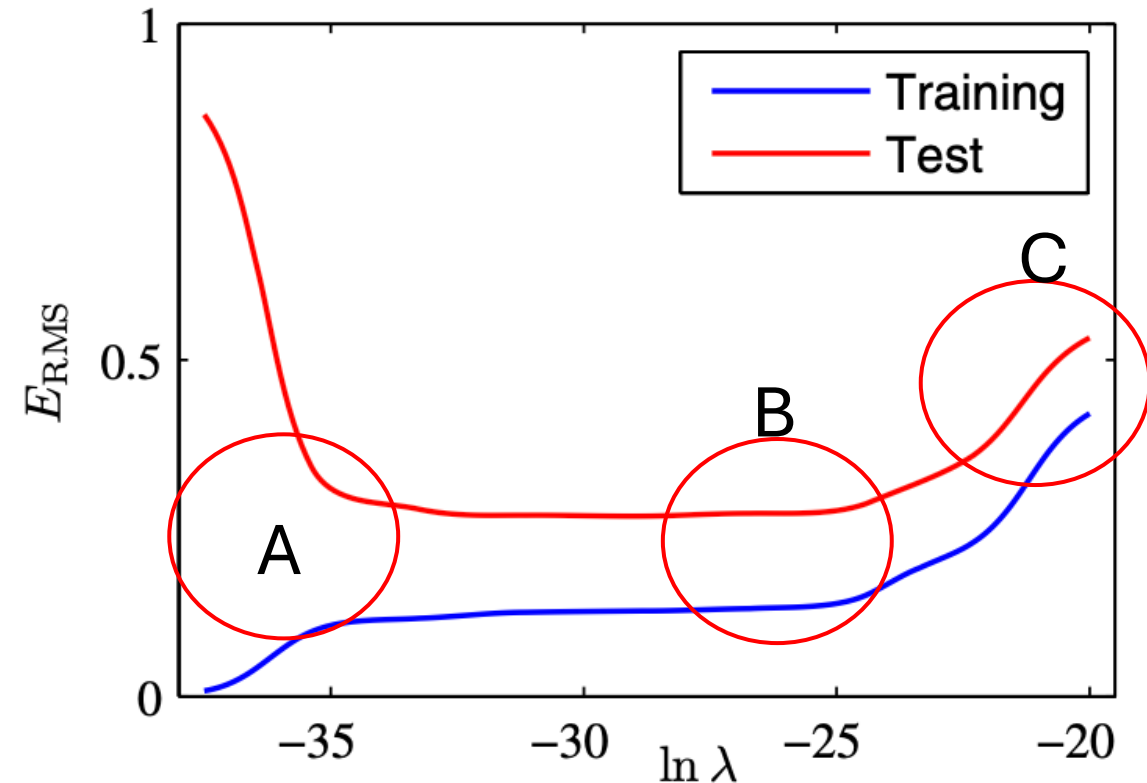
- Ridge regression regulates complexity.
- We need to choose regularization parameter:
- Q: which λ would you choose A or B?

+ In the lecture, I mentioned that the performance gap between train & test is small, so the answer was B instead of A (both test performance are similar)

+ But I realized that this may not be true (we only consider test error).

Let me revisit this problem again and open for discussion in the next lecture.

We need to choose the lambda showing the lowest test performance.



Classification Problem

How can we solve a classification problem?

How can we assign an input x to a certain class C_k ?

- Regression: Learning $y = f(x)$ “*the functional relation between x and y* ”
- Classification: Learning the discriminant functions $f_k(x)$

$C_k = \arg \max_k f_k(x)$: by using discriminant functions,
we can assign a class to x

To classify inputs to three groups,
How many discriminant functions do we need?

Example) three discriminant functions in the feature space of \mathbb{R}^2

1. $f_1(x_1, x_2) = x_2 - x_1 - 1$

2. $f_2(x_1, x_2) = x_2 + x_1 - 1$

3. $f_3(x_1, x_2) = x_2$

$$C_k = \arg \max_k f_k(x)$$

Q: Based on the discriminant functions above, assign a class for the points?

	f_1	f_2	f_3	class inference
$(-2,0)$	1	-3	0	1
$(0,0)$	-1	-1	0	3
$(2,0)$	-3	1	0	2

Example) computing the decision boundaries/ Splitting the feature space

1. $f_1(x_1, x_2) = x_2 - x_1 - 1$

2. $f_2(x_1, x_2) = x_2 + x_1 - 1$

3. $f_3(x_1, x_2) = x_2$

$$C_k = \arg \max_k f_k(x)$$

Q: Based on the discriminant functions above, assign a class for the points?

	f_1	f_2	f_3
$(-2,0)$	1	-3	0
$(0,0)$	-1	-1	0
$(2,0)$	-3	1	0

Example) computing the decision regions

1. $f_1(x_1, x_2) = x_2 - x_1 - 1$

2. $f_2(x_1, x_2) = x_2 + x_1 - 1$

3. $f_3(x_1, x_2) = x_2$

- $R_1: x_1 < -1$

- $R_2: x_1 \geq 1$

- $R_3: -1 \leq x_1 < 1$

Today, we are going to focus on **binary classification**. (+ / -)
To classify inputs by two groups,
How many discriminant functions do we need?

Today, we are going to focus on binary classification.

To classify inputs by two groups,

How many discriminant functions do we need?

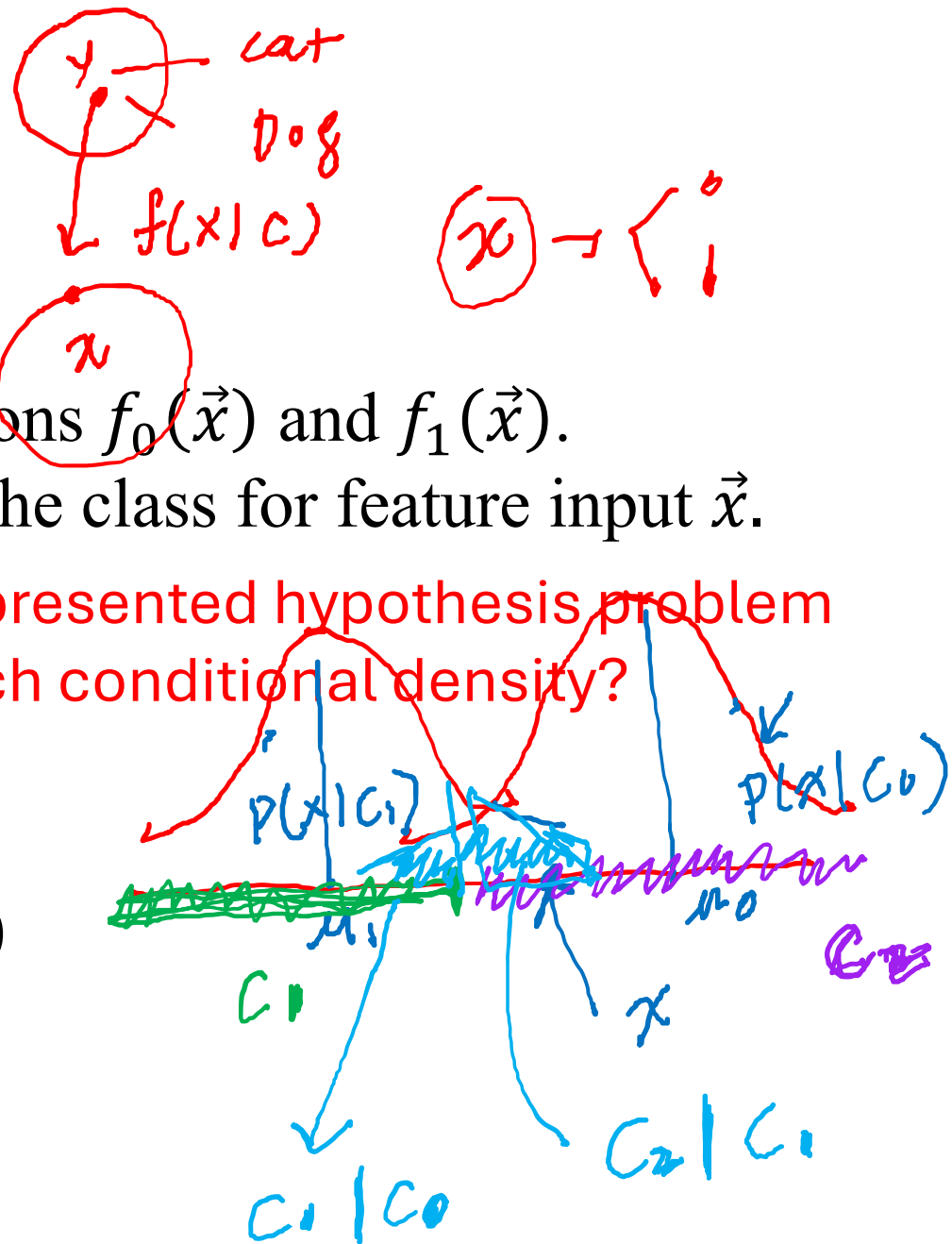
We want to design the two discriminant functions $f_0(\vec{x})$ and $f_1(\vec{x})$.

Once we have them, we can decide/infer the class for feature input \vec{x} .

classification can be represented hypothesis problem
 X is generated from which conditional density?

$$\begin{cases} \mathcal{H}_0 : & \vec{x} \sim f(\vec{x}|C_0) \\ \mathcal{H}_1 : & \vec{x} \sim f(\vec{x}|C_1) \end{cases}$$

$$f_0(\vec{x}) \underset{H_1}{\overset{H_0}{\geq}} f_1(\vec{x})$$



Binary Classification Problem & Decision Regions / Boundaries

Binary Classification Problem

$$\begin{cases} \mathcal{H}_0 : & \vec{x} \sim f(\vec{x}|C_0) \\ \mathcal{H}_1 : & \vec{x} \sim f(\vec{x}|C_1) \end{cases}$$

we need to learn
these discriminant functions.

$$\begin{matrix} H_0 \\ f_0(\vec{x}) \geq f_1(\vec{x}) \\ H_1 \end{matrix}$$

Decision Regions and Decision Boundaries

- $R_0 : f_0(\vec{x}) \geq f_1(\vec{x})$
- $R_1 : f_0(\vec{x}) < f_1(\vec{x})$

$$\begin{aligned} f_1(\vec{x}) &\stackrel{H}{=} f_0(\vec{x}) \\ f_1(\vec{x}) - f_0(\vec{x}) &= 0 \end{aligned}$$

Hyperplane
in the feature space
Q: dimension?

- MAP is a possible way forming the discriminant functions.

Posterior $P(H_0 | x)$

$$p(x|H_0)\pi_0 \stackrel{?}{\geq} p(x|H_1)\pi_1$$

$P(H_0 | x)$


$P(H_1 | x)$

$f_0(x)$ H_1

likelihood prior

$f_1(x)$

- MAP Rule Minimizes Expected Classification Error $E[R]$



$$\begin{aligned}
 E[R] &= \pi_0 \cdot E[R|H_0] + \pi_1 E[R|H_1] = \pi_1 \cdot \int_{y_0} f(y|H_1) dy + \pi_0 \cdot \int_{y_1} f(y|H_0) dy \quad P[H_0 = \text{true}] \\
 &= \pi_1 \cdot \int_{y_0} f(y|H_1) dy + \pi_0 \cdot \left(1 - \int_{y_0} f(y|H_0) dy\right) \\
 &= \pi_0 + \int_{y_0} \pi_1 \cdot f(y|H_1) - \pi_0 \cdot f(y|H_0) dy
 \end{aligned}$$

Q: How should we set the decision rule for y_0 ?

for y , if $\pi_1 \cdot f(y|H_1) - \pi_0 \cdot f(y|H_0) < 0$ then detect as H_0 R_0
 else if $\pi_1 \cdot f(y|H_1) - \pi_0 \cdot f(y|H_0) \geq 0$ then detect as H_1 R_1

$$\checkmark \frac{p(y|H_1)}{P(y|H_0)} \geq \frac{\pi_0}{\pi_1} \quad \text{"MAP rule"}$$

Gaussian Discriminant Analysis (GDA)

Q: Decision Boundaries for the Two Possible Cases

- $\Sigma_0 \neq \Sigma_1$
- $\Sigma_0 = \Sigma_1$

$$\begin{array}{ccc}
 & H_0 & \\
 \frac{1}{\sqrt{2\pi|\Sigma_0|}} \exp -\frac{1}{2}(x - \mu_0)^t \Sigma_0^{-1} (x - \mu_0) \cdot P[\mathcal{H}_0 = 0] & \geq & \frac{1}{\sqrt{2\pi|\Sigma_1|}} \exp -\frac{1}{2}(x - \mu_1)^t \Sigma_1^{-1} (x - \mu_1) \cdot P[\mathcal{H}_0 = 1] \\
 \underbrace{\hspace{10em}}_{PC_1} & & \underbrace{\hspace{10em}}_{PC_2} \\
 & H_1 & \\
 \ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi|\Sigma_0|}} - \frac{1}{2}(x - \mu_0)^t \Sigma_0^{-1} (x - \mu_0) & \geq & \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi|\Sigma_1|}} - \frac{1}{2}(x - \mu_1)^t \Sigma_1^{-1} (x - \mu_1) \\
 f_0(x) & & f_1(x)
 \end{array}$$

- Quadratic discriminant functions!
- Quadratic decision boundary!

Q: Decision Boundaries for the Two Possible Cases

- $\Sigma_0 = \Sigma_1$

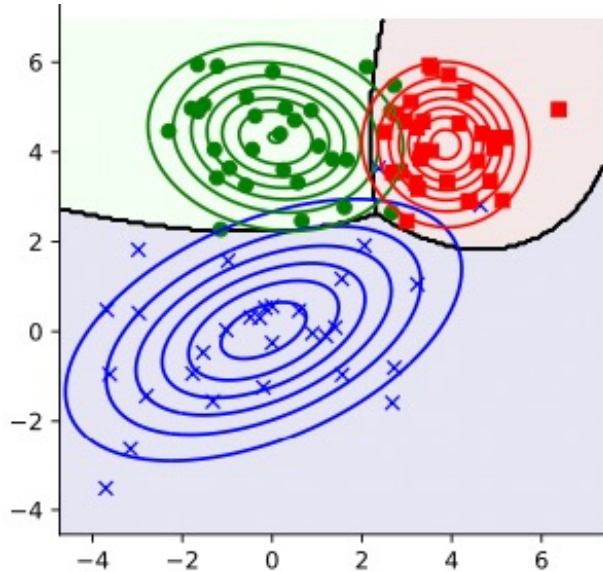
$$\begin{array}{ccc}
 & H_0 & \\
 \frac{1}{\sqrt{2\pi|\Sigma_0|}} \exp -\frac{1}{2}(x - \mu_0)^t \Sigma_0^{-1} (x - \mu_0) \cdot P[\mathcal{H}_0 = 0] & \gtrless & \frac{1}{\sqrt{2\pi|\Sigma_1|}} \exp -\frac{1}{2}(x - \mu_1)^t \Sigma_1^{-1} (x - \mu_1) \cdot P[\mathcal{H}_0 = 1] \\
 & H_1 & \\
 \ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi|\Sigma_0|}} - \frac{1}{2}(x - \mu_0)^t \Sigma_0^{-1} (x - \mu_0) & \gtrless & \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi|\Sigma_1|}} - \frac{1}{2}(x - \mu_1)^t \Sigma_1^{-1} (x - \mu_1) \\
 & \updownarrow & \\
 \ln P[\mathcal{H}_0] + \mu_0^t \Sigma_0^{-1} x - \frac{1}{2} \mu_0^t \Sigma_0^{-1} \mu_0 & \gtrless & \ln P[\mathcal{H}_1] + \mu_1^t \Sigma_1^{-1} x - \frac{1}{2} \mu_1^t \Sigma_1^{-1} \mu_1 \\
 & f_0(x) & f_1(x)
 \end{array}$$

Handwritten red annotations: $x^T \Sigma_0^{-1} x$ (crossed out), $x^T \Sigma_1^{-1} x$ (crossed out), and red lines striking through the Gaussian terms in the middle row.

- Linear discriminant functions!
- Linear decision boundary!

From Murphy Figure 9.2

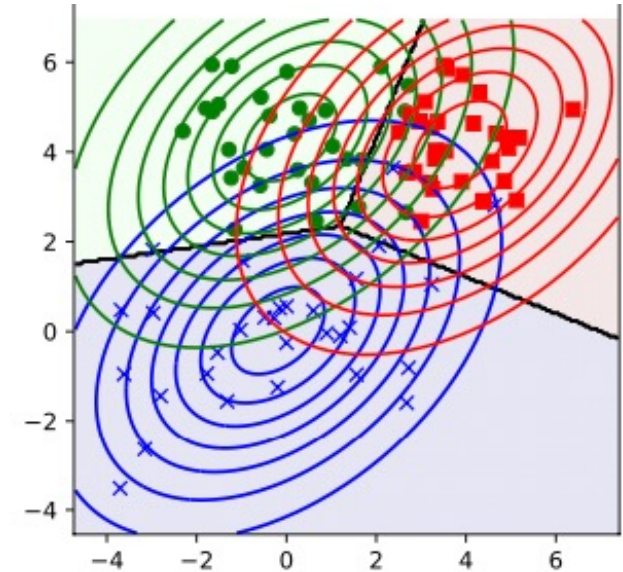
[Quadratic Discriminant Analysis]



$$\Sigma_0 \neq \Sigma_1$$

(Unconstrained Covariance)

[Linear Discriminant Analysis]



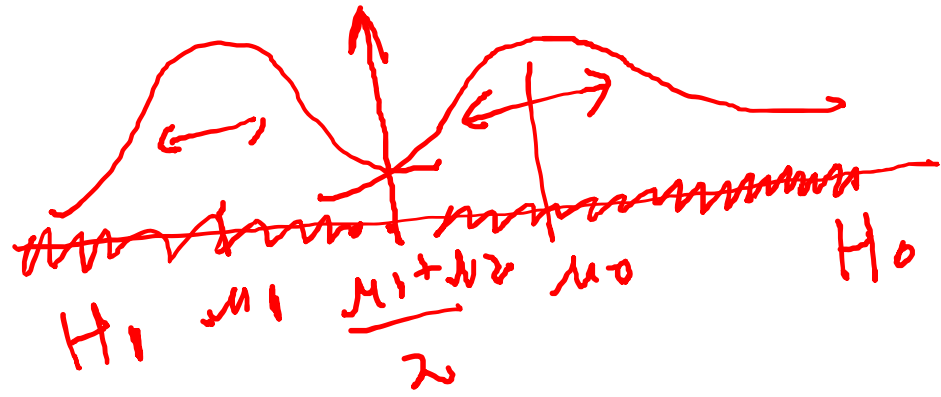
$$\Sigma_0 = \Sigma_1$$

(Tied Covariance)

QDA becomes LDA as assuming tied Covariance.

case 1]

- scalar feature
- $\sigma_0 = \sigma_1 = \sigma$
- $P[\mathcal{H}_0] = P[\mathcal{H}_1]$



$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_0)^2 \geq \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_1)^2$$

$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_0)^2 \geq \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_1)^2$$

$$\frac{1}{\sigma^2}x\mu_0 - \frac{1}{2\sigma^2}\mu_0^2 \geq \frac{1}{\sigma^2}x\mu_1 - \frac{1}{2\sigma^2}\mu_1^2$$

$$x(\mu_0 - \mu_1) \geq \frac{1}{2}(\mu_0^2 - \mu_1^2)$$

$$x \geq \frac{1}{2}(\mu_0 + \mu_1)$$

WLOG if $(\mu_0 > \mu_1)$
Binary classification
decision rule!

case 2]

- feature vector
- $\Sigma_0 = \Sigma_1 = \sigma I$, isotropic
- $P[\mathcal{H}_0] = P[\mathcal{H}_1]$

H_0

$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_0)^2 \gtrless \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_1)^2$$

H_1



$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi\sigma^N}} + -\frac{1}{2\sigma^2}(x - \mu_0)^t(x - \mu_0) \gtrless \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi\sigma^N}} + -\frac{1}{2\sigma^2}(x - \mu_1)^t(x - \mu_1)$$

$$\frac{1}{\sigma^2}\mu_0^t x - \frac{1}{2\sigma^2}\mu_0^t \mu_0 \gtrless \frac{1}{\sigma^2}\mu_1^t x - \frac{1}{2\sigma^2}\mu_1^t \mu_1$$

$$(\mu_0 - \mu_1)^t x \gtrless \frac{1}{2}(\mu_0 - \mu_1)^t(\mu_0 + \mu_1)$$



projection to $(\mu_0 - \mu_1)$

then the decision rule same as the scalar²² case.

case 3]

- feature vector
- $\Sigma_0 = \Sigma_1 = \Sigma$, anisotropic
- $P[\mathcal{H}_0] = P[\mathcal{H}_1]$

Q: projection to $(\mu'_0 - \mu'_1)$
then the decision rule same as the scalar
case?

$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_0)^2 \gtrless \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(x - \mu_1)^2$$

$$\ln P[\mathcal{H}_0] + \ln \frac{1}{\sqrt{2\pi|\Sigma|}} - \frac{1}{2}(x - \mu_0)^t \Sigma^{-1} (x - \mu_0) \gtrless \ln P[\mathcal{H}_1] + \ln \frac{1}{\sqrt{2\pi|\Sigma|}} - \frac{1}{2}(x - \mu_1)^t \Sigma^{-1} (x - \mu_1)$$

$$\mu_0^t \Sigma^{-1} x - \frac{1}{2} \mu_0^t \Sigma^{-1} \mu_0 \gtrless \mu_1^t \Sigma^{-1} x - \frac{1}{2} \mu_1^t \Sigma^{-1} \mu_1$$

$$(\mu_0^t - \mu_1^t) \Sigma^{-1} x \gtrless \frac{1}{2} \mu_0^t \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^t \Sigma^{-1} \mu_1$$

$$(\mu_0 - \mu_1)^t E \Lambda^{-1} E^t x \gtrless \frac{1}{2} \mu_0^t E \Lambda^{-1} E^t \mu_0 - \frac{1}{2} \mu_1^t E \Lambda^{-1} E^t \mu_1$$