

Classification

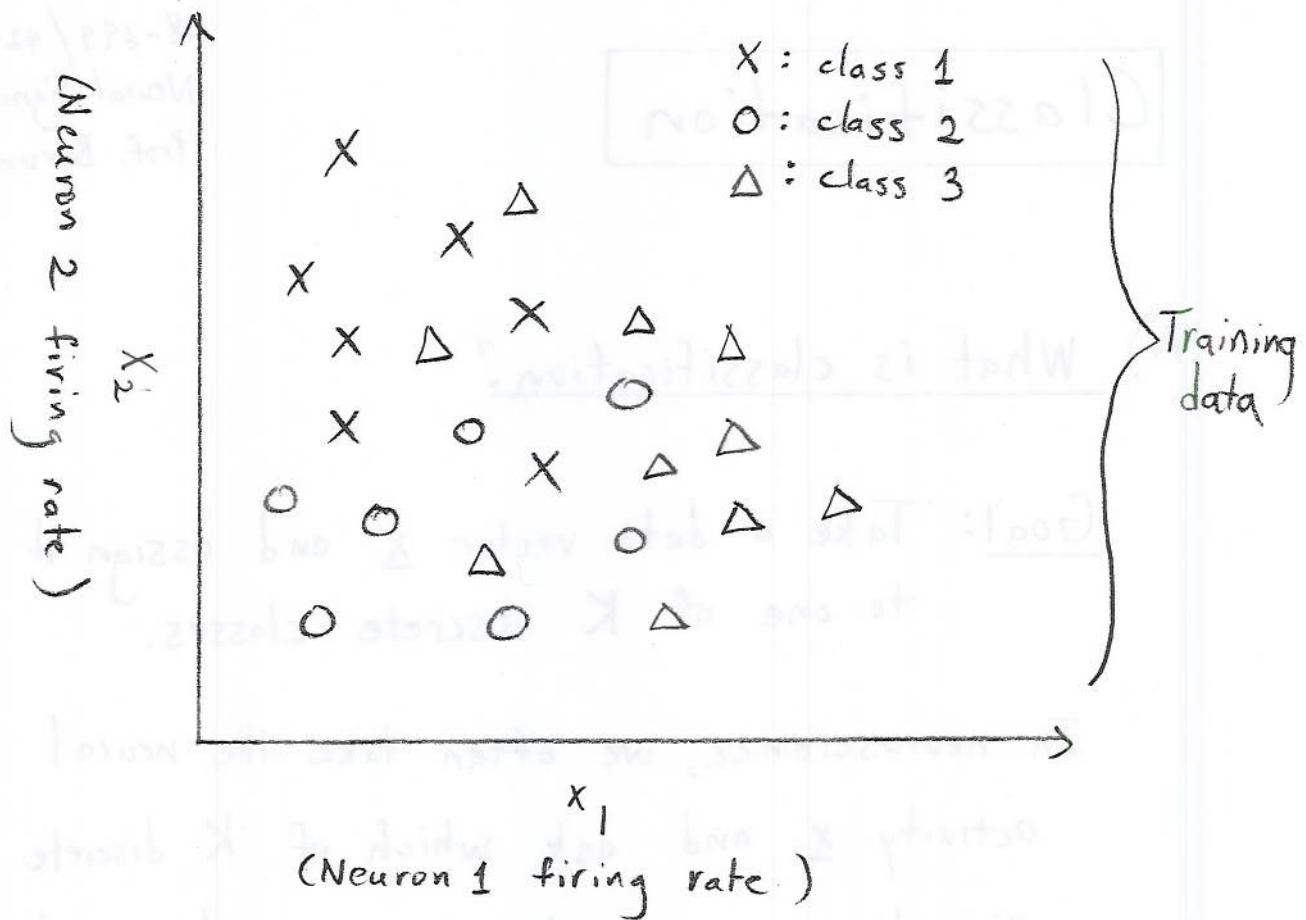
A) What is classification?

Goal: Take a data vector x and assign it to one of K discrete classes.

In neuroscience, we often take the neural activity x and ask which of K discrete stimuli gave rise to the observed neural activity.

To begin, we need labeled training data (i.e., we know the class label of each training data point). We will consider the problem of classification with unlabeled training data later in the course (Ch. 9).

Training data needs labels



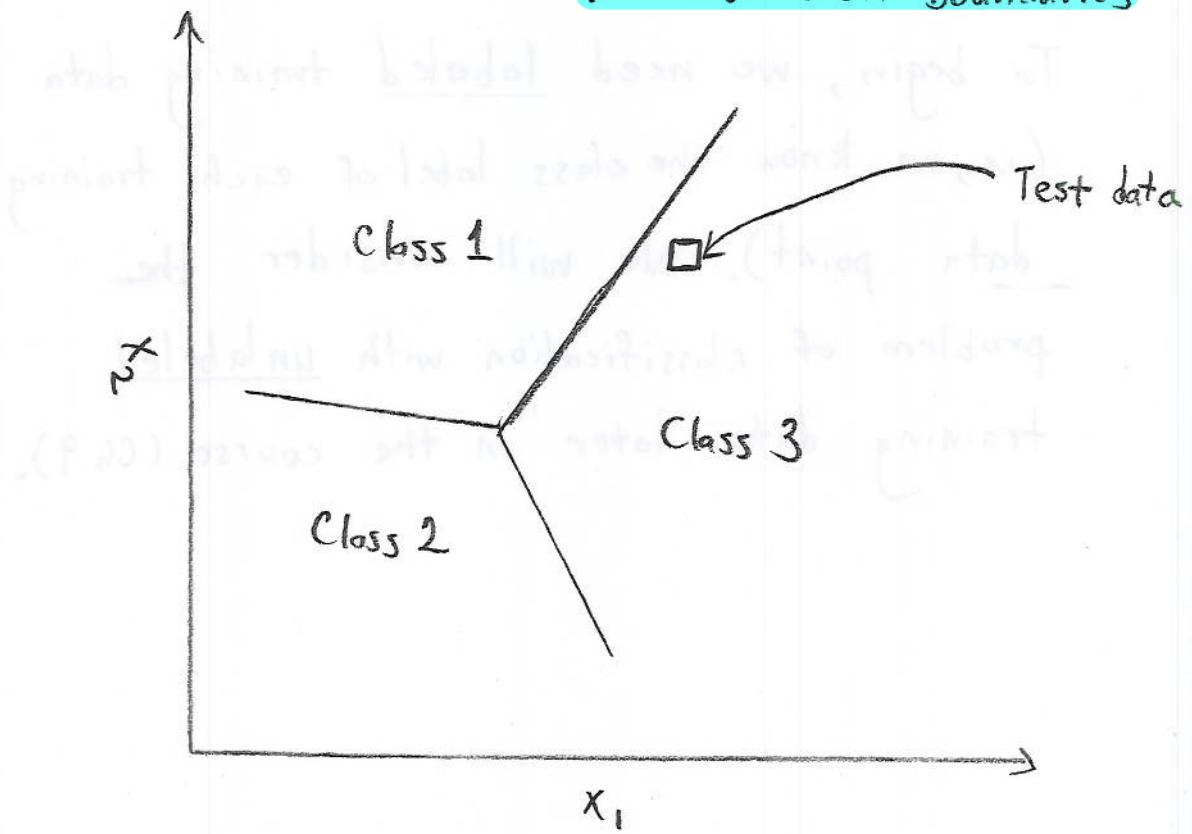
(Neuron 1 firing rate)

(Neuron 2 firing rate)

x : class 1
 o : class 2
 \triangle : class 3

Training data

↓ using methods we will discuss,
 find decision boundaries



For test data \underline{x} , what class does it belong to?

B) Classifying Using Generative Models

Training phase: Finding decision Boundaries

- Fit class-conditional densities $P(\underline{x} | C_k)$ and class priors $P(C_k)$ to training data.
- Probability of train data for given label
($k = 1, \dots, K$)

Test phase:

- Compute $P(C_k | \underline{x})$ using Bayes' rule

$$P(C_k | \underline{x}) = \frac{P(\underline{x} | C_k) P(C_k)}{P(\underline{x})}$$

Posterior

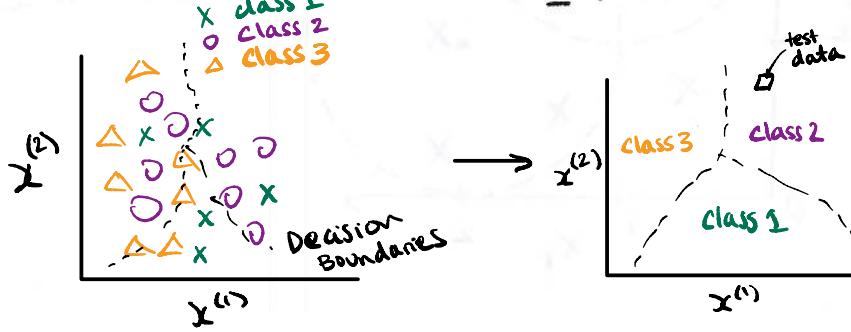
Probability of label given test data

$$= \frac{P(\underline{x} | C_k) P(C_k)}{\sum_{j=1}^K P(\underline{x} | C_j) P(C_j)}$$

All possible combos of training data

- Assign class $\hat{k} = \underset{k}{\operatorname{argmax}} P(C_k | \underline{x})$ most likely label

to test data \underline{x} .



8.1) Generative models

Allows you to
generate fake
data from model

$P(x|C_k)$ and $P(C_k)$ define a "probabilistic generative model". This means that we can generate synthetic data from the model.

For example, say there are two classes and $x \in \mathbb{R}^2$

$$P(C_1) = 0.7$$

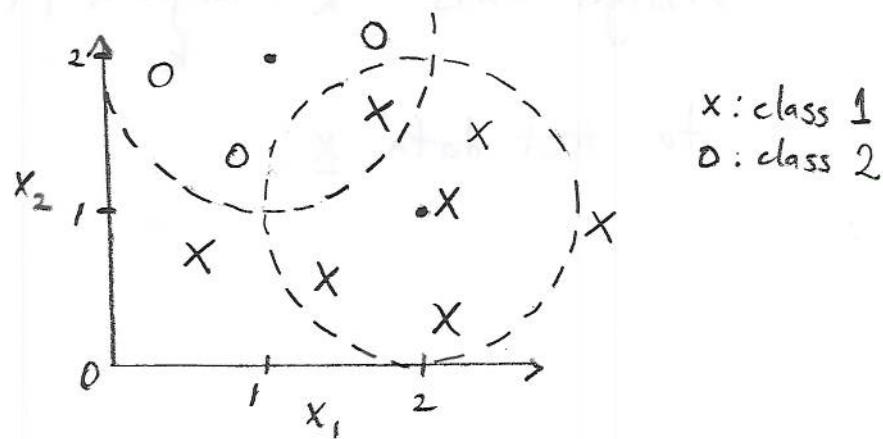
$$P(C_2) = 0.3$$

$$P(x|C_1) = N\left(\begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

$$P(x|C_2) = N\left(\begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

To generate one synthetic data vector x , first flip a biased coin with probability 0.7 of coming up heads.

- If heads, draw from the Gaussian $P(x|C_1)$.
- If tails, draw from the Gaussian $P(x|C_2)$.



Philosophy of generative models:

If we generate synthetic data from the model and it looks a lot like the real data we're trying to model, then we have a good model for our real data.

We can then use the generative model to make optimal inferences, decisions, etc.

Goal: find set of model parameters that best fits your data & best predicts future observations w/ \uparrow likelihood

B.2) Training phase: Maximum likelihood parameter estimation

Maximize the likelihood of the observed data w.r.t. model parameters.

Example: Two classes with Gaussian class-conditional density with shared covariance

Write down probability of training data in terms of model parameters
Training data: $\{x_n, t_n\} \quad n=1, \dots, N$
 $t_n = 1$ denotes class C_1
 $t_n = 0$ denotes class C_2
Now observed data is constant
Tweak parameters

$$P(x_n | C_i) = P(x_n | C_i) P(C_i) \\ = N(x_n | \mu_i, \Sigma_i) \pi_i$$

Let $P(t_n=1) = P(C_1) = \pi$

$$P(t_n=0) = P(C_2) = 1 - \pi$$

For a data point $\underline{x}_n \in \mathbb{R}^D$,

$$P(\underline{x}_n | C_1) = P(\underline{x}_n | C_1) P(C_1) = N(\underline{x}_n | \mu_1, \Sigma) \cdot \pi$$

$$P(\underline{x}_n | C_2) = P(\underline{x}_n | C_2) P(C_2) = N(\underline{x}_n | \mu_2, \Sigma) \cdot (1 - \pi)$$

Data likelihood for N data points together:

$$\begin{aligned} \mathcal{L} &= P(\{\underline{x}_n, t_n\} | \pi, \mu_1, \mu_2, \Sigma) \\ &= \prod_{n=1}^N \underbrace{\left(N(\underline{x}_n | \mu_1, \Sigma) \cdot \pi \right)^{t_n}}_{\substack{\text{train data} \\ \text{train label}}} \underbrace{\left(N(\underline{x}_n | \mu_2, \Sigma) \cdot (1 - \pi) \right)^{1-t_n}}_{\substack{\text{prior} \\ \text{mean class 1} \\ \text{mean class 2} \\ \text{covariance}}} \end{aligned}$$

$$\log \mathcal{L} = \sum_{n=1}^N \left[t_n \log N(\underline{x}_n | \mu_1, \Sigma) + (1 - t_n) \log N(\underline{x}_n | \mu_2, \Sigma) \right]$$

$$+ (1 - t_n) \log (1 - \pi) \quad \boxed{\text{class 1 we want this term}} \quad \boxed{\text{class 2 we want this term}}$$

$$\frac{d}{dx} \log |X| = X^{-T}$$

where

$$\frac{d}{dx} \text{Tr}(X^{-T} A) = -X^{-T} A X^{-T}$$

$$-\frac{1}{2} \log |\Sigma| = -\frac{1}{2} X^{-T}$$

$$\log N(\underline{x}_n | \mu_k, \Sigma) = \frac{1}{2} \text{Tr}(\Sigma^{-1} (\underline{x}_n - \mu_k)(\underline{x}_n - \mu_k)^T)$$

$$\rightarrow 1 \times 1 \\ = \text{Tr}(\text{itself})$$

$$-\frac{1}{2} \log |\Sigma| - \frac{D}{2} \log (2\pi)$$

i) Find $\pi \leftarrow$ probability for class 1 (prior)

$$\frac{\partial \log \mathcal{L}}{\partial \pi} = \sum_{n=1}^N \left[t_n \cdot \frac{1}{\pi} - (1-t_n) \frac{1}{1-\pi} \right] = 0$$

$$(1-\pi) \sum_{n=1}^N t_n - \pi \sum_{n=1}^N (1-t_n) = 0$$

$$(1-\pi) N_1 - \pi (N - N_1) = 0$$

\downarrow let N_1 = number of data points from C_1

$$\boxed{\pi = \frac{N_1}{N}}$$

\leftarrow prior

$$N_1 = \sum_{n=1}^N t_n$$

$$N_2 = \sum_{n=1}^N (1-t_n)$$

ii) Find μ_1

$$\frac{d}{d\mu} (\underline{x}^\top A \underline{x}) = (A + A^\top) \underline{x}$$

$= 2A\underline{x} \leftarrow$ if A is symmetrical

$$\frac{\partial \log \mathcal{L}}{\partial \mu_1} = \sum_{n=1}^N \left(t_n \cdot \frac{1}{2} \cdot 2 \sum^{-1} (\underline{x}_n - \mu_1) \right) = 0$$

\swarrow symmetrical covariance

$$\sum^{-1} \left(\sum_{n=1}^N t_n \underline{x}_n \right) = \sum^{-1} \left(\mu_1 \sum_{n=1}^N t_n \right)$$

$\underline{x}_n \in \mathbb{R}^D$

$$\boxed{\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \underline{x}_n}$$

of label 1 Points

train label

train data

Analogously,

$$\boxed{\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \underline{x}_n}$$

Normal model parameters

$$\text{Tr}(ABC\dots) = \text{Tr}(BC\dots A) = \text{Tr}(C\dots AB)$$

iii) Find $\sum_{D \times D} \frac{\partial \log \mathcal{L}}{\partial \Sigma}$

Focusing only on terms that involve Σ ,

$$\begin{aligned} \log \mathcal{L} = \sum_{n=1}^N & \left[t_n \left(-\frac{1}{2} \text{Tr}(\Sigma^{-1}(\underline{x}_n - \mu_1)(\underline{x}_n - \mu_1)^T) - \frac{1}{2} \log |\Sigma| \right) \right. \\ & \left. + (1-t_n) \left(-\frac{1}{2} \text{Tr}(\Sigma^{-1}(\underline{x}_n - \mu_2)(\underline{x}_n - \mu_2)^T) - \frac{1}{2} \log |\Sigma| \right) \right] \end{aligned}$$

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial \Sigma} = \sum_{n=1}^N & \left[t_n \left(-\frac{1}{2} \cdot -\Sigma'(\underline{x}_n - \mu_1)(\underline{x}_n - \mu_1)^T \Sigma' - \frac{1}{2} \Sigma' \right) \right. \\ & \left. + (1-t_n) \left(-\frac{1}{2} \cdot -\Sigma'(\underline{x}_n - \mu_2)(\underline{x}_n - \mu_2)^T \Sigma' - \frac{1}{2} \Sigma' \right) \right] \\ = & [0] \end{aligned}$$

Rearranging yields

$$\begin{aligned} & \frac{1}{2} \sum_{n \in C_1} (\underline{x}_n - \mu_1)(\underline{x}_n - \mu_1)^T - \frac{1}{2} N_1 \Sigma \\ & + \frac{1}{2} \sum_{n \in C_2} (\underline{x}_n - \mu_2)(\underline{x}_n - \mu_2)^T - \frac{1}{2} N_2 \Sigma = [0] \end{aligned}$$

Σ found in training phase:

$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$, where

↑ weight
fraction of data belonging to this class

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (\underline{x}_n - \mu_1)(\underline{x}_n - \mu_1)^T$$

$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (\underline{x}_n - \mu_2)(\underline{x}_n - \mu_2)^T$$

$$\frac{N_1}{N} \cdot \frac{1}{N_1} \dots$$

Parameters: $\pi, \mu_1, \mu_2, \Sigma$

B.3) Test phase: Assigning a new data point to a class

$$\begin{aligned}
 \hat{k} &= \underset{k}{\operatorname{argmax}} P(C_k | x) \\
 &= \underset{k}{\operatorname{argmax}} \frac{P(x|C_k)P(C_k)}{P(x)} \quad \rightarrow \text{Bayes Rule} \\
 &= \underset{k}{\operatorname{argmax}} P(x|C_k)P(C_k) \quad \rightarrow \text{omit } x \text{ since it doesn't depend on } k \\
 &= \underset{k}{\operatorname{argmax}} \left(\log P(x|C_k) + \log P(C_k) \right) \\
 &= \underset{k}{\operatorname{argmax}} \underbrace{\left(\mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log P(C_k) \right)}_{\text{call this } a_k(x)}
 \end{aligned}$$

Class that maximizes probability of each class
 Classes ($k=1, \dots, K$)
 test data
 likelihood Prior
 same no matter what x

What do the decision boundaries look like in x space?

$$\begin{array}{ccc}
 (k \times D) (D \times D) (D \times k) & (k \times D) (D \times D) (D \times k) \\
 (k \times k) - & (k \times k) + (1 \times k) \\
 (k \times k) + (1 \times k) \\
 \left[\begin{smallmatrix} k_1 & k_2 & \dots \\ k_3 & \dots & \vdots \end{smallmatrix} \right] + \left[\begin{smallmatrix} k_1 \\ k_2 \\ \vdots \end{smallmatrix} \right] = (k \times k)
 \end{array}$$

Why log?

Whatever maximizes x also maximizes $\log(x)$

C) Hyperplanes

A hyperplane is the D-dimensional generalization of a line in 2-dim space and a plane in 3-dim space. A hyperplane is D-1 generalization

A hyperplane is defined as the set of all \underline{x}

such that

\underline{w} determines direction of hyperplane

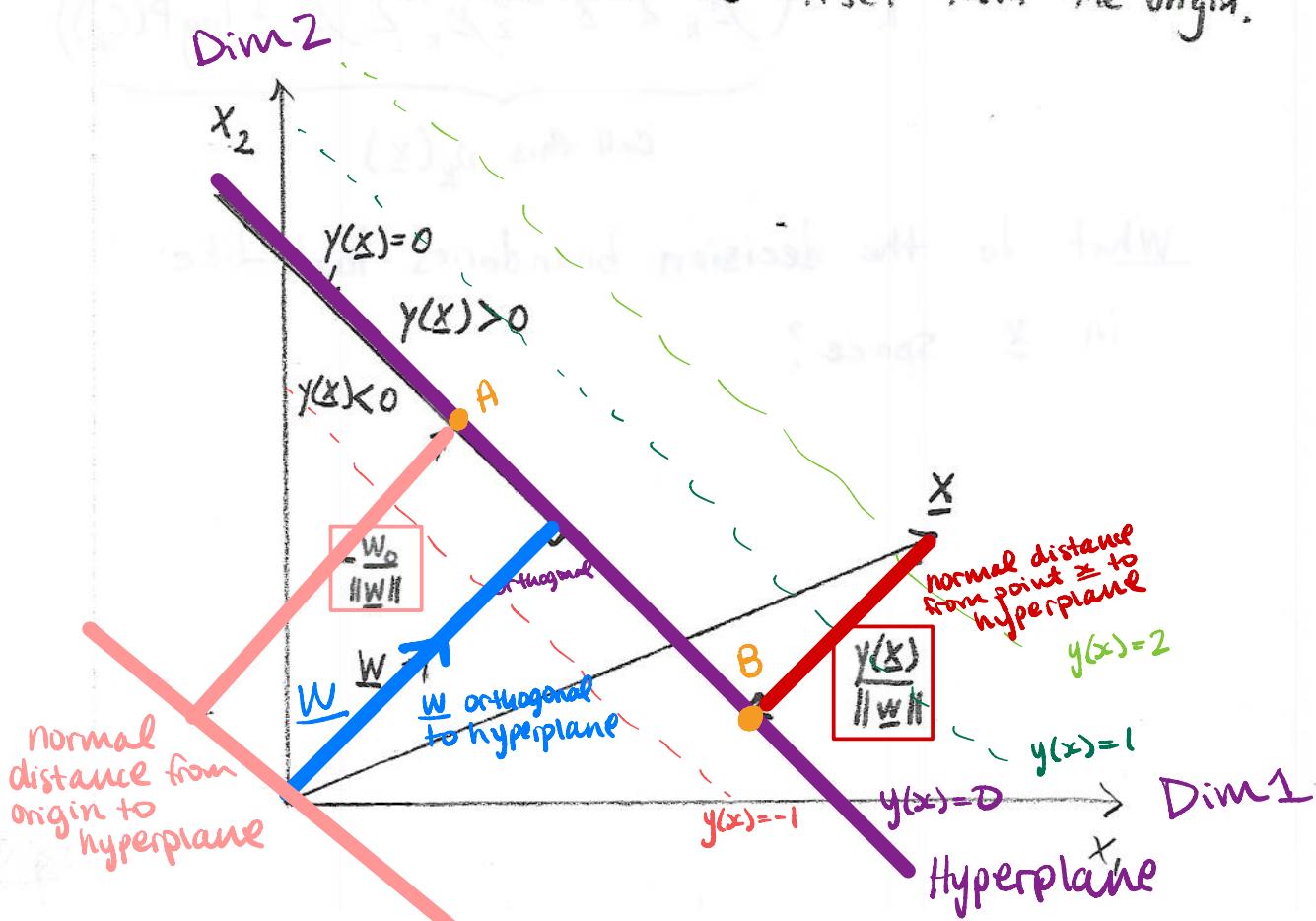
offset of hyperplane w.r.t respect to (0,0)

$$y(\underline{x}) = \underline{w}^T \underline{x} + w_0 = 0 \quad (1)$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
 $(1 \times 1) \quad (D \times 1) \quad (1 \times 1) \quad (1 \times 1)$

\underline{w} determines the direction of the hyperplane

w_0 determines its offset from the origin.



Facts:

i) \underline{w} is orthogonal to hyperplane

Consider two points \underline{x}_A and \underline{x}_B which lie on hyperplane.

$$y(\underline{x}_A) = y(\underline{x}_B) = 0$$

$$\underline{w}^T \underline{x}_A + w_0 = \underline{w}^T \underline{x}_B + w_0$$

$$\underline{w}^T (\underbrace{\underline{x}_A - \underline{x}_B}) = 0$$

vector lying in hyperplane

$\Rightarrow \underline{w}$ is orthogonal to any vector lying in hyperplane.

ii) Normal distance from origin to hyperplane is $-\frac{w_0}{\|\underline{w}\|}$

Let \underline{x} be a point on hyperplane $\Rightarrow \underline{w}^T \underline{x} + w_0 = 0$

Normal distance is projection of \underline{x} onto \underline{w}

$$\left(\frac{\underline{w}}{\|\underline{w}\|} \right)^T \underline{x} = -\frac{w_0}{\|\underline{w}\|}$$

iii) Normal distance from any point \underline{x} to hyperplane

is $\frac{y(\underline{x})}{\|\underline{w}\|}$.

Project \underline{x} onto \underline{w} , then subtract $-\frac{w_0}{\|\underline{w}\|}$

$$\left(\frac{\underline{w}}{\|\underline{w}\|} \right)^T \underline{x} + \frac{w_0}{\|\underline{w}\|} = \frac{y(\underline{x})}{\|\underline{w}\|}$$

$$\underbrace{\left(\boldsymbol{\mu}_k^T \Sigma^{-1} \underline{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log P(C_k) \right)}_{\text{call this } a_k(\underline{x})}$$

D) Linear Decision Boundaries

From p.9, a point \underline{x} is assigned to class C_k

if $a_k(\underline{x}) > a_j(\underline{x})$ for all $j \neq k$.

Thus, the decision boundary between class C_k and class C_j is given by $a_k(\underline{x}) = a_j(\underline{x})$.

Let $a_k(\underline{x}) = \underline{w}_k^T \underline{x} + w_{k0}$, where

↑ in example
w/ 2 classes
this is at 0.5?

$$\underline{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log P(C_k)$$

The decision boundary is thus

$$(\underline{w}_k - \underline{w}_j)^T x + (w_{k0} - w_{j0}) = 0$$

$(D-1)$ dimensional hyperplane in \mathbb{R}^D

Note that this takes the same form as (1), so the decision boundary is a $(D-1)$ dimensional hyperplane in \mathbb{R}^D .

Appendix: Useful matrix properties

$$\frac{d}{d\underline{x}} \underline{x}^T A \underline{x} = (A + A^T) \underline{x} \stackrel{\text{A symmetric.}}{\downarrow} 2A\underline{x}$$

vector $\xrightarrow{\frac{d}{d\underline{x}}}$

$$\frac{d}{dX} \text{Tr}(X^{-1}A) = -X^{-T} A^T X^{-T}$$

matrix $\xrightarrow{\frac{d}{dX}}$

$$\frac{d}{dX} \log|X| = X^{-T}$$

matrix $\xrightarrow{\frac{d}{dX}}$

$$\text{Tr}(ABC\dots) = \text{Tr}(BCD\dots A) = \text{Tr}(CD\dots AB)$$

A good reference is:

<http://www.ee.ic.ac.uk/hp/staff/dmb/matrix/intro.html>

or

Simply google "matrix reference manual".