# Machine Learning I:  Supervised Methods

### B. Keith Jenkins

## Announcements

- Homework 1 is due Friday

- Slido poll questions in lecture will begin soon - initially just a trial run

## Today's lecture

- Discriminant functions and approaches to multiclass problems $(C > 2)$  (part 2)

  - One vs. Rest (OvR)  (part 2)

  - One vs. One (OvO)

  - Maximal Value Method (MVM)

  - Summary

- Computational complexity (part 1)

  *deferred to Lecture 5.*

## 2. Multiclass $(C > 2)$ problems

→ Can we pose a $C$-class problem $(C > 2)$ as a set of 2-class problems? Yes.

(i) Use $C$ discriminant fcns: $g_k(\underline{x})$, $k = 1, 2, \cdots, C$.

__One vs. rest (OvR)__   (also called One vs-all)
to define $C$ 2-class problems.
Each 2-class problem:

$$S_k' \quad \text{vs.} \quad \overline{S_k'} \qquad \begin{array}{l}(\text{e.g., hs. price decr.}\\ \text{vs. hs. price not decr.})\\ (\text{cat vs. not cat})\end{array}$$

for $k = 1, 2, \cdots, C.$

[see plot below]

Combine results:

OvR __Decision rule__ : $\underline{x} \in \Gamma_k$ IFF $\underline{x} \in \Gamma_k'$ AND $\underline{x} \in \overline{\Gamma_j'} \ \forall \ j \neq k.$
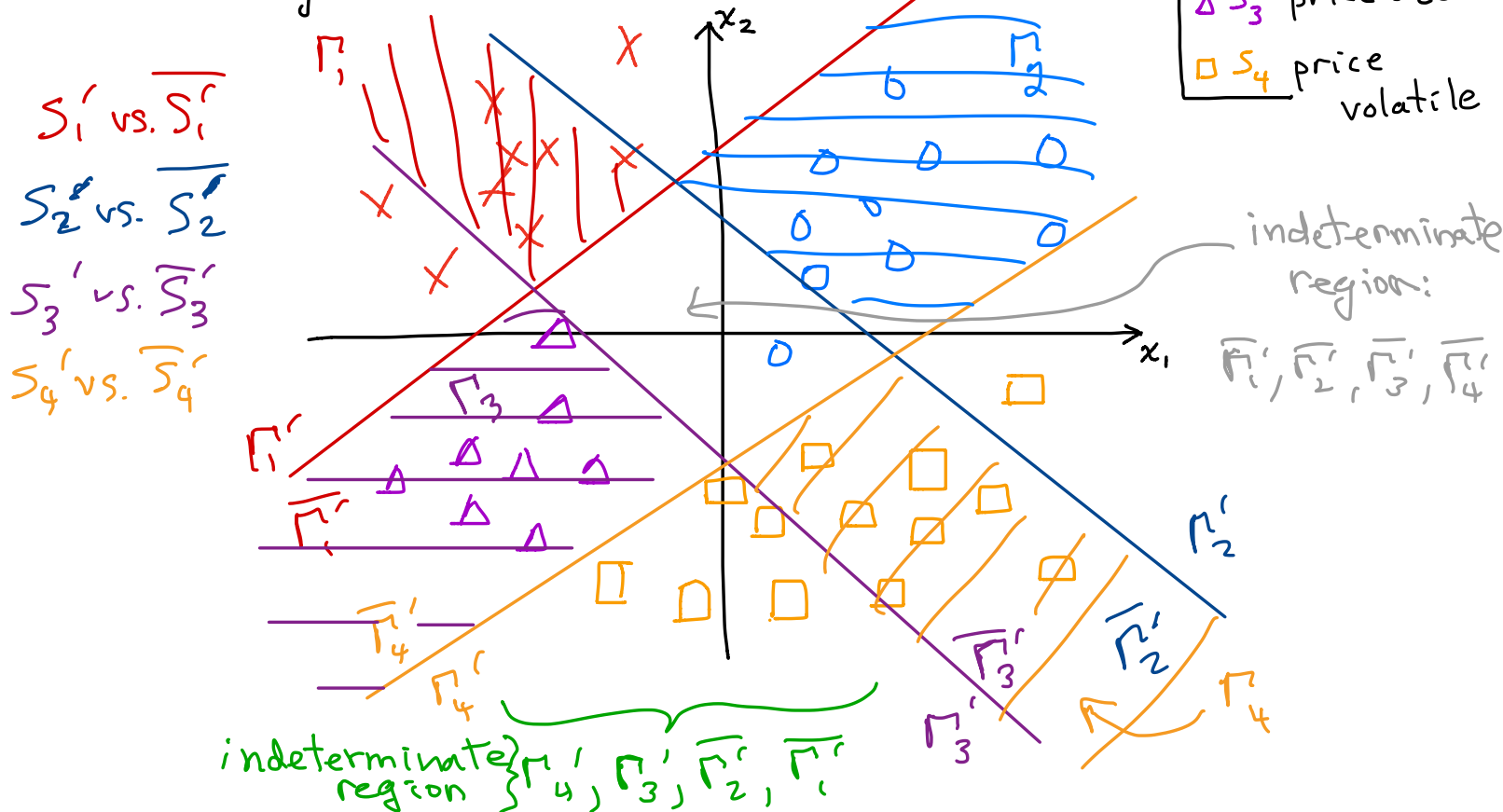
our default OvR decision rule. Other OvR decision rules are possible.

# OvR



$\underline{x} \longrightarrow$ | $S_1' \text{ vs. } \bar{S}_1'$ | $\longrightarrow \Gamma_1' \text{ or } \bar{\Gamma}_1'$

$\underline{x} \longrightarrow$ | $S_2' \text{ vs } \bar{S}_2'$ |

$\vdots$

$\underline{x} \longrightarrow$ | $S_c' \text{ vs. } \bar{S}_c'$ | $\Gamma_c' \text{ or } \bar{\Gamma}_c'$

OvR default decision rule $\longrightarrow \underline{x} \in \Gamma_k'$

Example: Consider a $C = 4$-class problem with $D = 2$ features

assume: each 2-class classifier is linear.

OvR method

Training dataset:

$S_1'$ vs. $\overline{S_1'}$

$S_2'$ vs. $\overline{S_2'}$

$S_3'$ vs. $\overline{S_3'}$

$S_4'$ vs. $\overline{S_4'}$



| | |
|---|---|
| $\times$ $S_1$ | price incr. |
| $\circ$ $S_2$ | price const. |
| $\triangle$ $S_3$ | price decr. |
| $\square$ $S_4$ | price volatile |

indeterminate region:

$\overline{\Gamma_1'}, \overline{\Gamma_2'}, \overline{\Gamma_3'}, \overline{\Gamma_4'}$

indeterminate region $\{\Gamma_4', \Gamma_3', \overline{\Gamma_2'}, \overline{\Gamma_1'}\}$

Apply 4 2-class problems: each $S_k'$ vs $\overline{S_k'}$

$\Gamma_1$ is defined by: all $\underline{x} \in \{\Gamma_1'$ and $\overline{\Gamma_2'}$ and $\overline{\Gamma_3'}$ and $\overline{\Gamma_4'}\}$
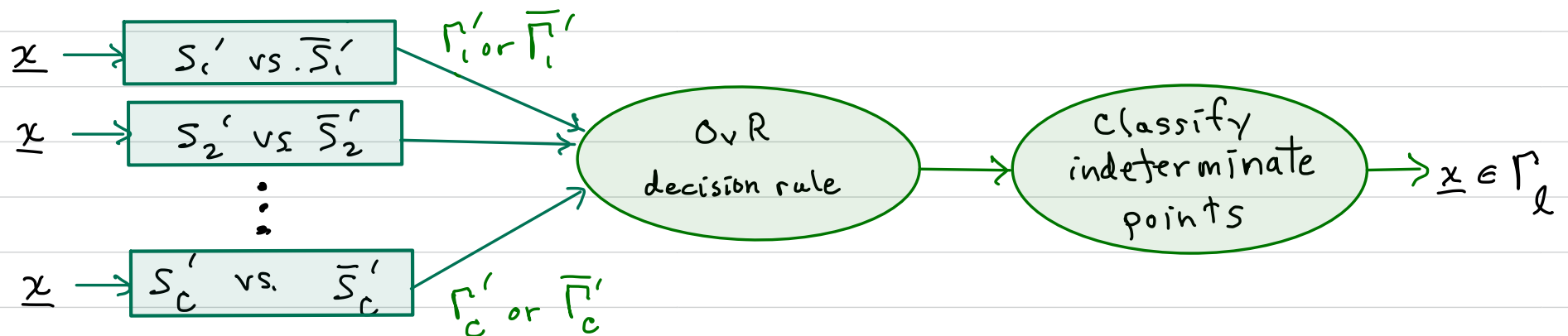
(for our default OvR decision rule.)

# Comments on OvR method

1. | **Def :** Given a dataset $\mathcal{D}$. If all data points $\underline{x}_i^{(k)}$ of class $S_k$ be separated from all points of all other classes by a linear boundary, and this holds for all $k = 1, 2, \cdots, C$, then $\mathcal{D}$ is <u>totally linearly separable</u>.

2. In practice, for methods that can result in indeterminate regions, often an additional ad-hoc rule is used to classify the indeterminate points; for example:

For $\underline{x}$ in indeterminate region,
$$\underline{x} \in \Gamma_k \quad \text{iff} \quad k = \underset{k}{\operatorname{argmax}} \left[ g_k(\underline{x}) \right] \quad \longrightarrow \quad \underline{x} \in \Gamma_k$$

3. Alternate OvR decision rules can also be used. For example, if there is a confidence measure for each binary classification result for a given point $\underline{x}$ :

Confidence ( $S_k'$ given $\underline{x}$ )     (from $S_k'$ vs $\overline{S}_k'$ classification of $\underline{x}$)

then final decision rule can be:

$$\underline{x} \in \Gamma_{k^*} \text{ iff } \quad k^* = \underset{k}{\text{argmax}} \left[ \text{Confidence} (\Gamma_k' \text{ given } \underline{x}) \right].$$

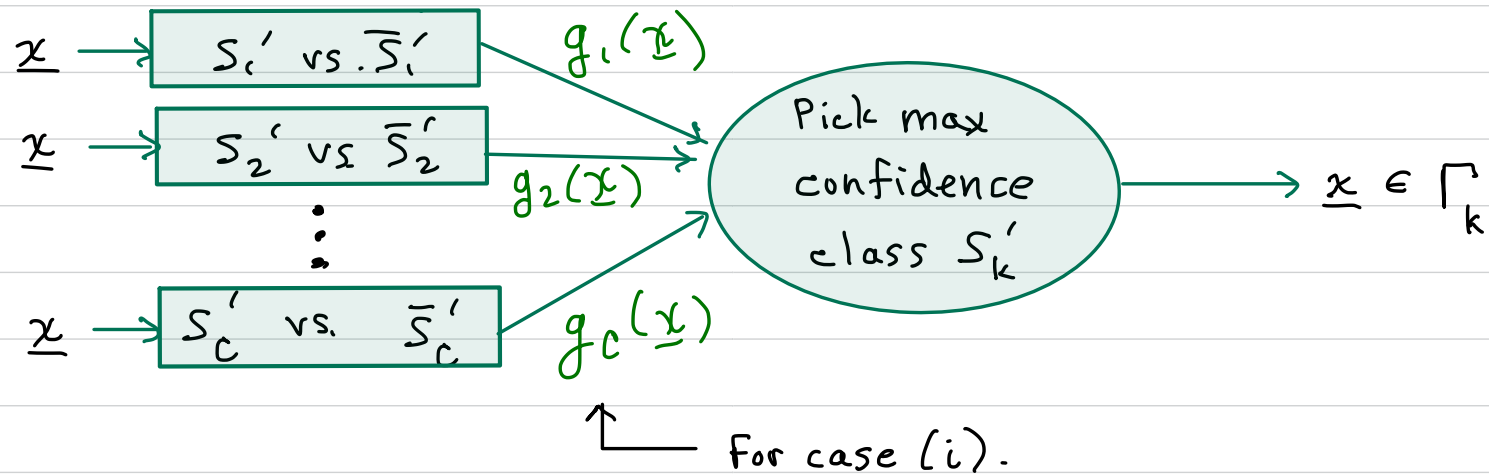→ Alternate OvR decision rule.

Examples of confidence measures:

(i) Confidence ( $\Gamma_k'$ given $\underline{x}$ ) $= g_k (\underline{x})^*$     ⎫ Confidence in
(ii) Confidence ( $\Gamma_k'$ given $\underline{x}$ ) $= P(S_k' | \underline{x})$     ⎬ $S_k'$ vs. $\overline{S}_k'$ result, for input point $\underline{x}$.

$^*$ Assumes scale of $g_j(\underline{x})$ is comparable to $g_k(\underline{x})$.

In this case we have:



$$x \longrightarrow \boxed{S_1' \text{ vs. } \bar{S}_1'} \xrightarrow{g_1(x)}$$

$$x \longrightarrow \boxed{S_2' \text{ vs } \bar{S}_2'} \xrightarrow{g_2(x)}$$

$$\vdots$$

$$x \longrightarrow \boxed{S_c' \text{ vs. } \bar{S}_c'} \xrightarrow{g_c(x)}$$

Pick max confidence class $S_k'$ $\longrightarrow$ $x \in \Gamma_k$

For case (i).

(ii) <u>One vs. one</u> (OvO)    (or all vs. all)   (or all pairs)

→ set of $S_k$ vs. $S_j$ decisions.  (all possible pairs)

Use discr. fcns $g_{kj}(\underline{x})$.   Each 2-class problem:

$$g_{kj} \underset{\Gamma_j}{\overset{\Gamma_k}{\gtrless}} 0 \quad ; \qquad g_{jk}(\underline{x}) = -g_{kj}(\underline{x})$$

Combine results:

> <u>Decision rule</u> (OvO, default):
>
> $\underline{x} \in \Gamma_k$ iff $g_{kj}(\underline{x}) > 0 \quad \forall\, j \neq k.$

How many 2-class problems?    $\binom{c}{2} = \dfrac{c(c-1)}{2}$

Example: Consider a $C = 4$-class problem with $D = 2$ features
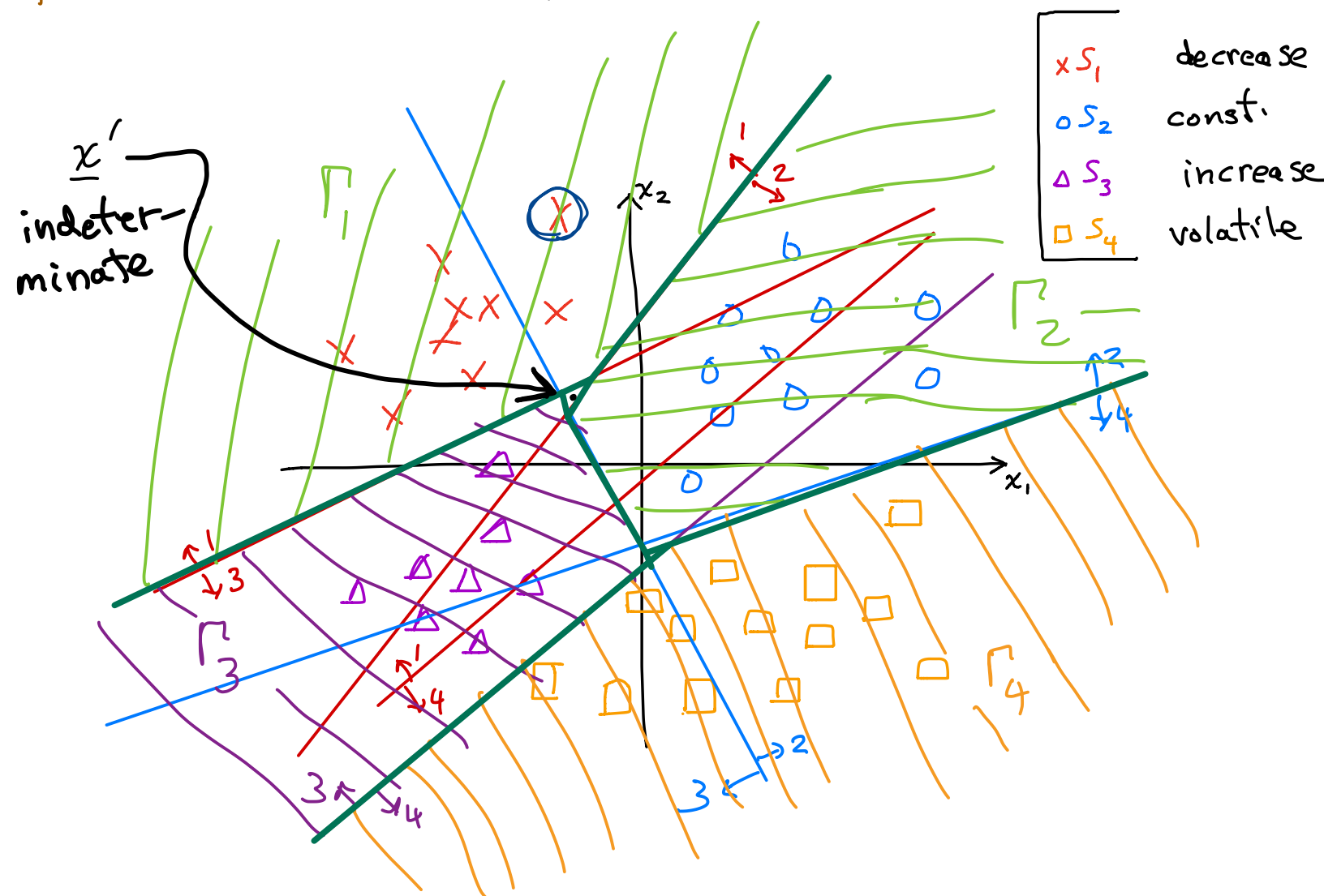
Same data as previous example.

Number of 2-class problems $= \dfrac{C(C-1)}{2} = \dfrac{4 \cdot 3}{2} = 6$

Training dataset:



× $S_1$
∘ $S_2$
△ $S_3$
□ $S_4$

## One vs. one

At $\underline{x}'$: can't be: $\Gamma_1, \Gamma_2, \Gamma_3, \Gamma_4$ $\Rightarrow$ indeterminate



Comments (this example):

1. has indeterm. regions (smaller than OvR).
2. No errors on training data

# Comments on OvO method

1. 
> **Def**: If $\exists$ $\frac{c(c-1)}{2}$ linear separating boundaries $H_{kj}$, such that $H_{kj}$ separates all data points of $S_k$ from all data points of $S_j$, $j \neq k$, then the data is _pairwise linearly separable._

2. For any indeterminate points, an additional ad hoc rule can be used to classify them.

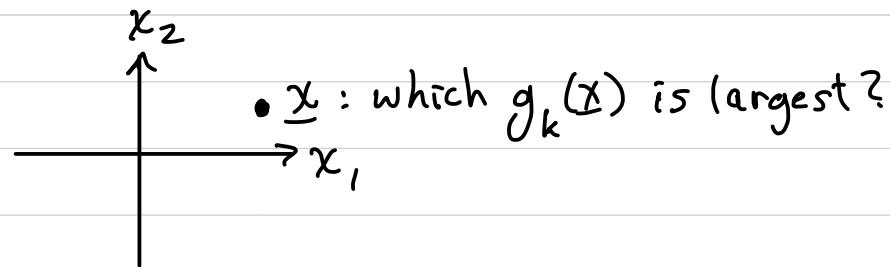3. An alternate decision rule for OvO:

   Take vote: how many 2-class classifiers decided $S_k'$ over $S_j'$?
   $\underline{x} \in \Gamma_k$ if $S_k'$ has the largest number of votes.

(iii) <u>Maximal Value Method</u> (MVM)

- for the linear case, called a <u>linear machine</u>.
- is <u>not</u> based on a set of 2-class classifiers.

Let each class $S_k$ have 1 discriminant function: $g_k(\underline{x})$

Decision rule:     $g_k(\underline{x}) > g_j(\underline{x}) \quad \forall j \neq k \quad \Rightarrow \quad \underline{x} \in \Gamma_k$



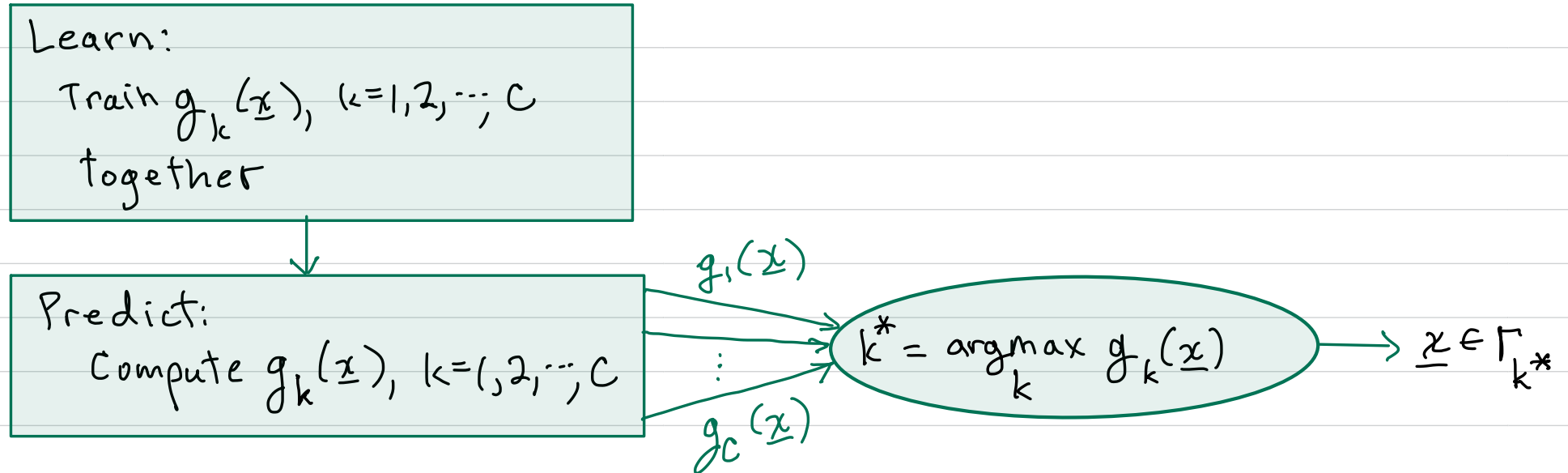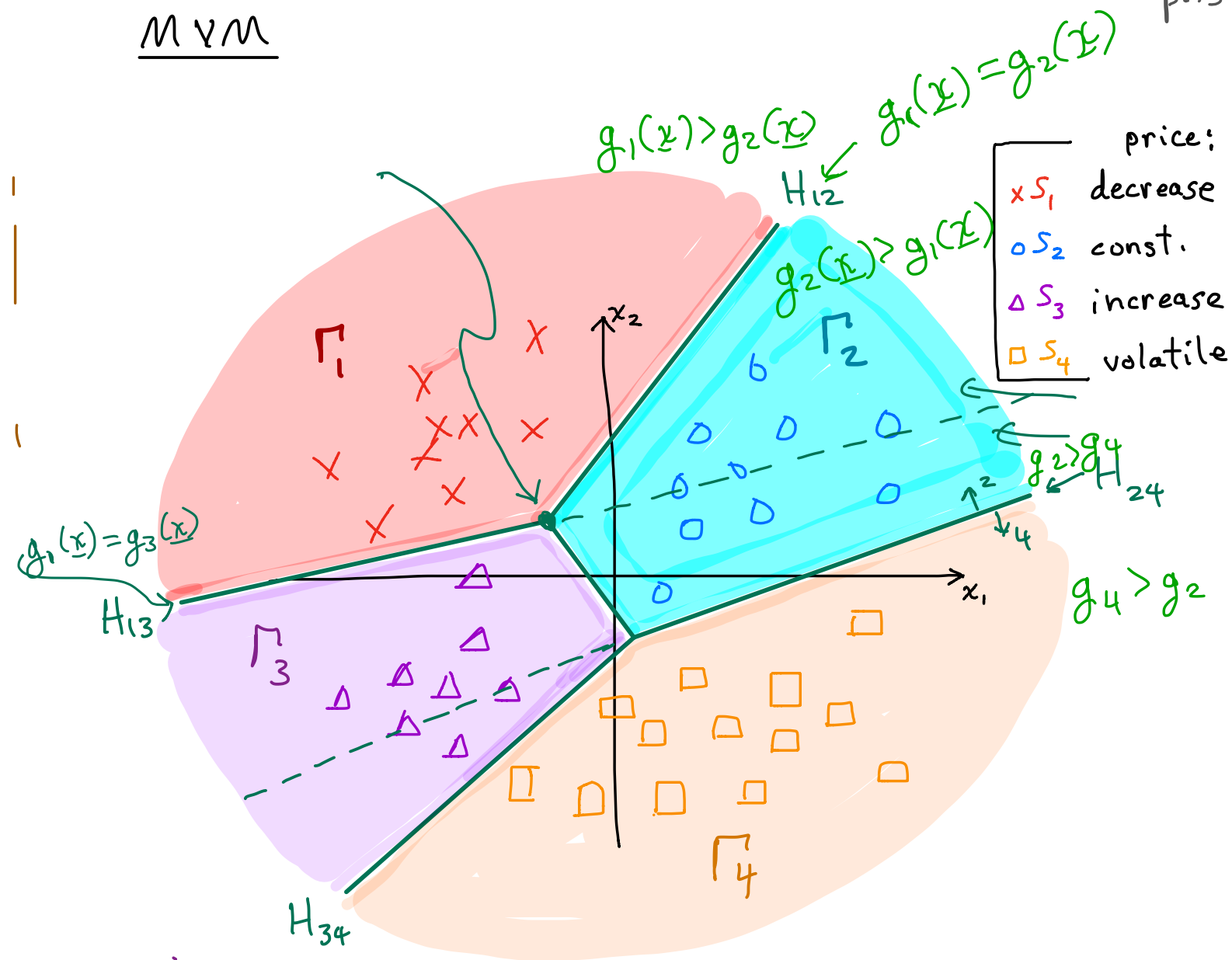$\bullet \underline{x}$ : which $g_k(\underline{x})$ is largest?

Decision rule can also be expressed as:

$$\underline{x} \in \Gamma_k \text{ iff } k = \arg\max_m \{ g_m(\underline{x}) \}$$   MVM decision rule

$\Rightarrow$ Decision boundary between $\Gamma_k$ and $\Gamma_j$ is: $g_k(\underline{x}) = g_j(\underline{x})$

(some boundaries may be redundant)

## MVM

Learn:

Train $g_k(\underline{x})$, $k=1,2,\cdots,C$
together

Predict:

Compute $g_k(\underline{x})$, $k=1,2,\cdots,C$

$g_1(\underline{x})$

$\vdots$

$g_c(\underline{x})$

$k^* = \underset{k}{\arg\max} \; g_k(\underline{x})$

$\underline{x} \in \Gamma_{k^*}$

# MVM

$g_1(\underline{x}) > g_2(\underline{x})$   $g_1(\underline{x}) = g_2(\underline{x})$

$H_{12}$

$g_2(\underline{x}) > g_1(\underline{x})$

price:

x $S_1$ decrease

o $S_2$ const.

△ $S_3$ increase

□ $S_4$ volatile

$\Gamma_1$

$x_2$

$\Gamma_2$

$g_2 > g_4$   $H_{24}$

$g_1(\underline{x}) = g_3(\underline{x})$

$g_4 > g_2$

$H_{13}$

$\Gamma_3$

$x_1$

$\Gamma_4$

$H_{34}$

*

$H_{jk}: \quad g_j(\underline{x}) = g_k(\underline{x})$

# Comments on MVM

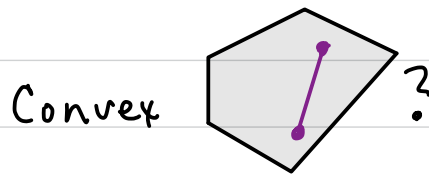1.   <u>Def:</u> If there exists $g_i(\underline{x})$, $i = 1, 2, \cdots, C$, such that

$$g_k \left[ \underline{x}_m^{(k)} \right] > g_j \left[ \underline{x}_m^{(k)} \right]$$

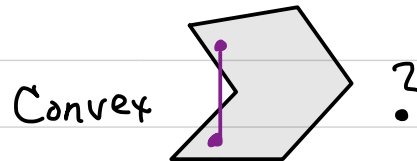$$\forall m = 1, 2, \cdots, N_k, \quad \forall j \neq k, \quad \forall k$$

with all $g_i(\underline{x})$ expressible as linear functions, then the dataset is <u>linearly separable</u>.

2.  No indeterminate regions (if unlikely special cases are avoided, like $g_k(\underline{x}) = g_j(\underline{x})$ over some region).

3.  If $g_k(\underline{x})$ are linear $\forall k$, then decision regions $\Gamma_k$ are convex.

Convex     ?       Convex     ?

Yes                  No.

# Summary of multiclass classification approaches

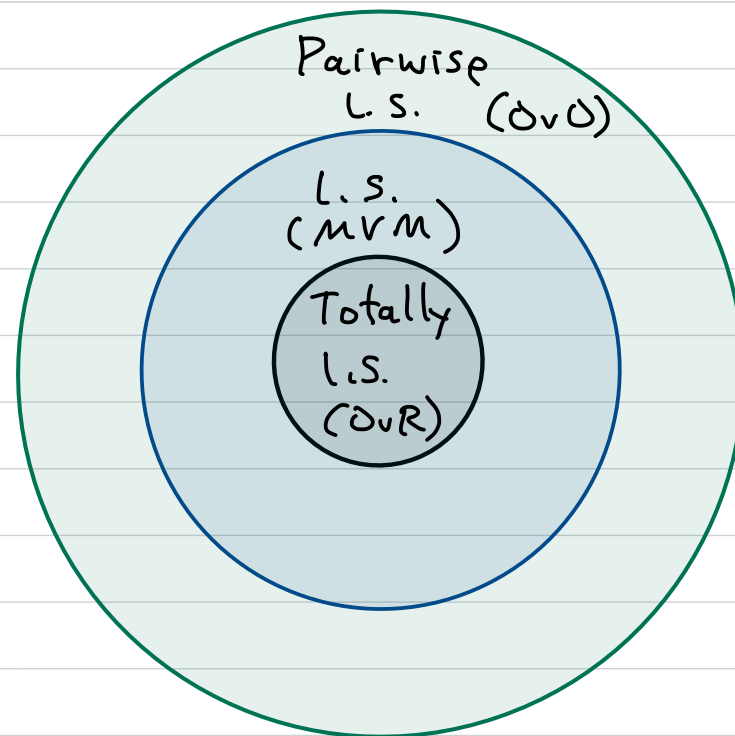| Method | Linear $g(\underline{x})$ defines type of linear separability |
|--------|-----------------------------------------------|
| OvR | Totally linearly separable |
| OvO | Pairwise linearly separable |
| MVM | Linearly separable |

How powerful is each method?
→ How large is the class of datasets that each method can fully separate using linear $g(\underline{x})$?

→ Venn diagram:

## Nonlinear case

Can use any of the above 3 methods ($OvR$, $OvO$, $MVM$), by letting each $g_k(\underline{x})$ or each $g_{kj}(\underline{x})$ be a nonlinear function of $\underline{x}$.

(More later)