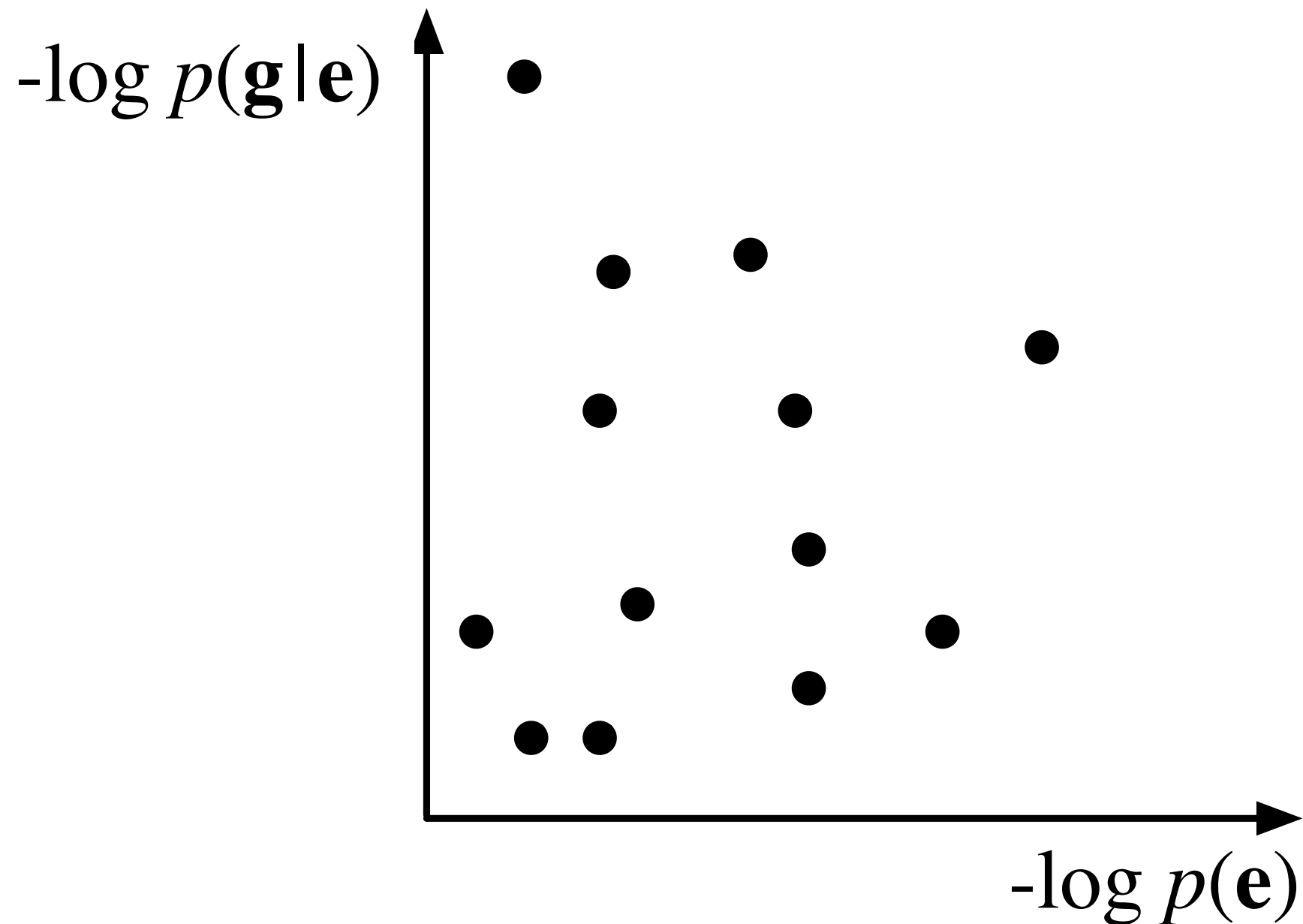


Discriminative Training II: MERT

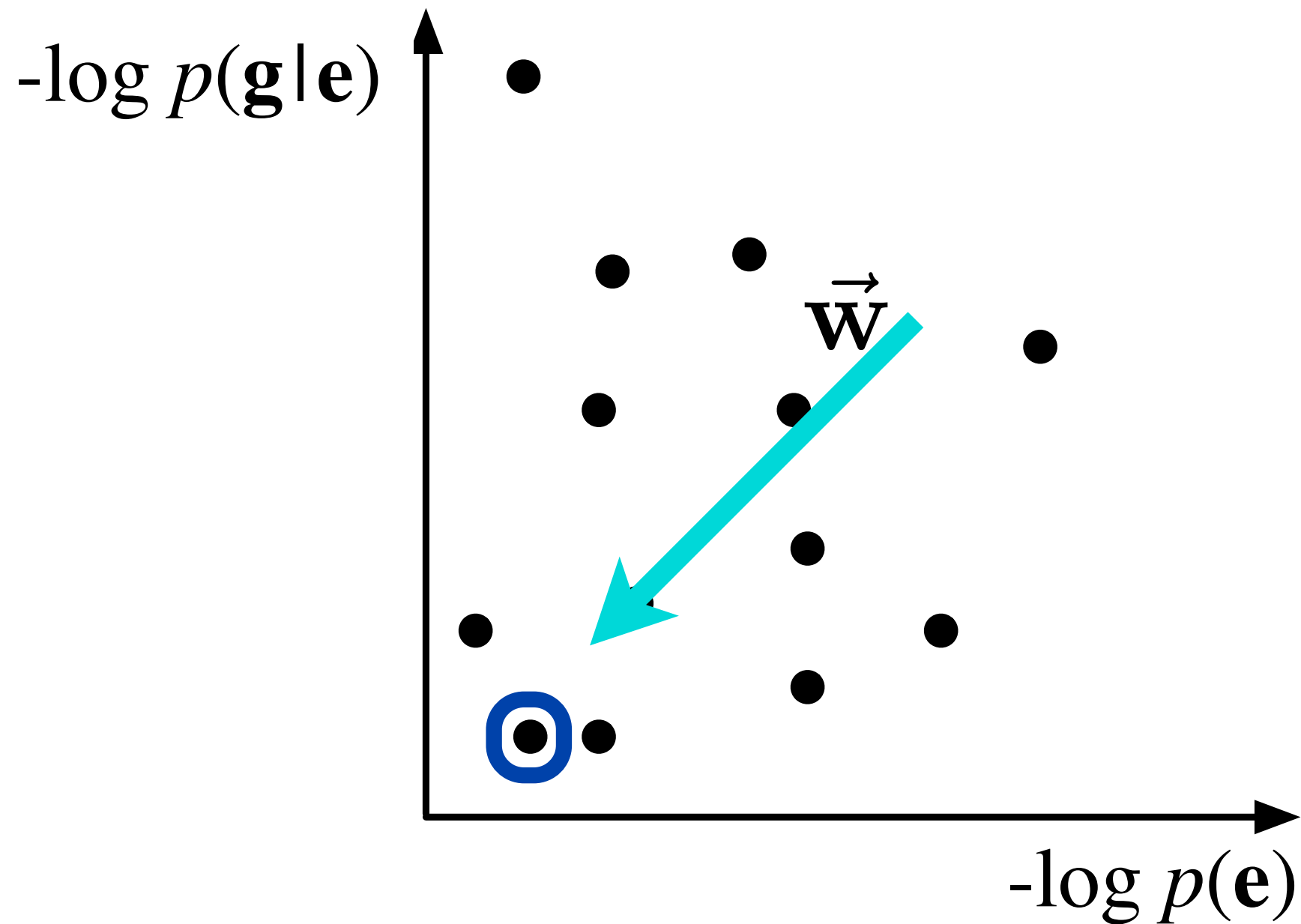
April 3, 2014



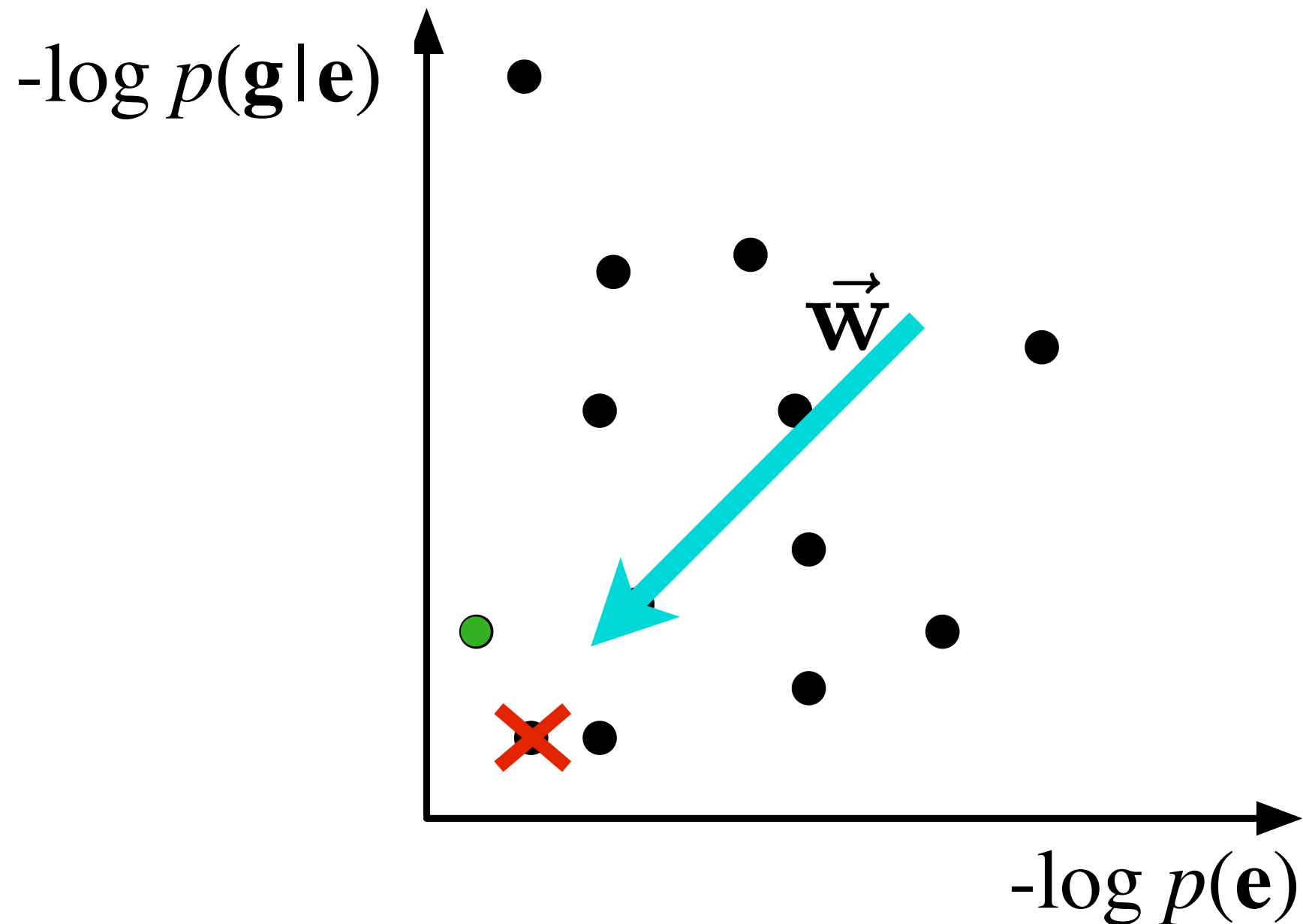
The Noisy Channel



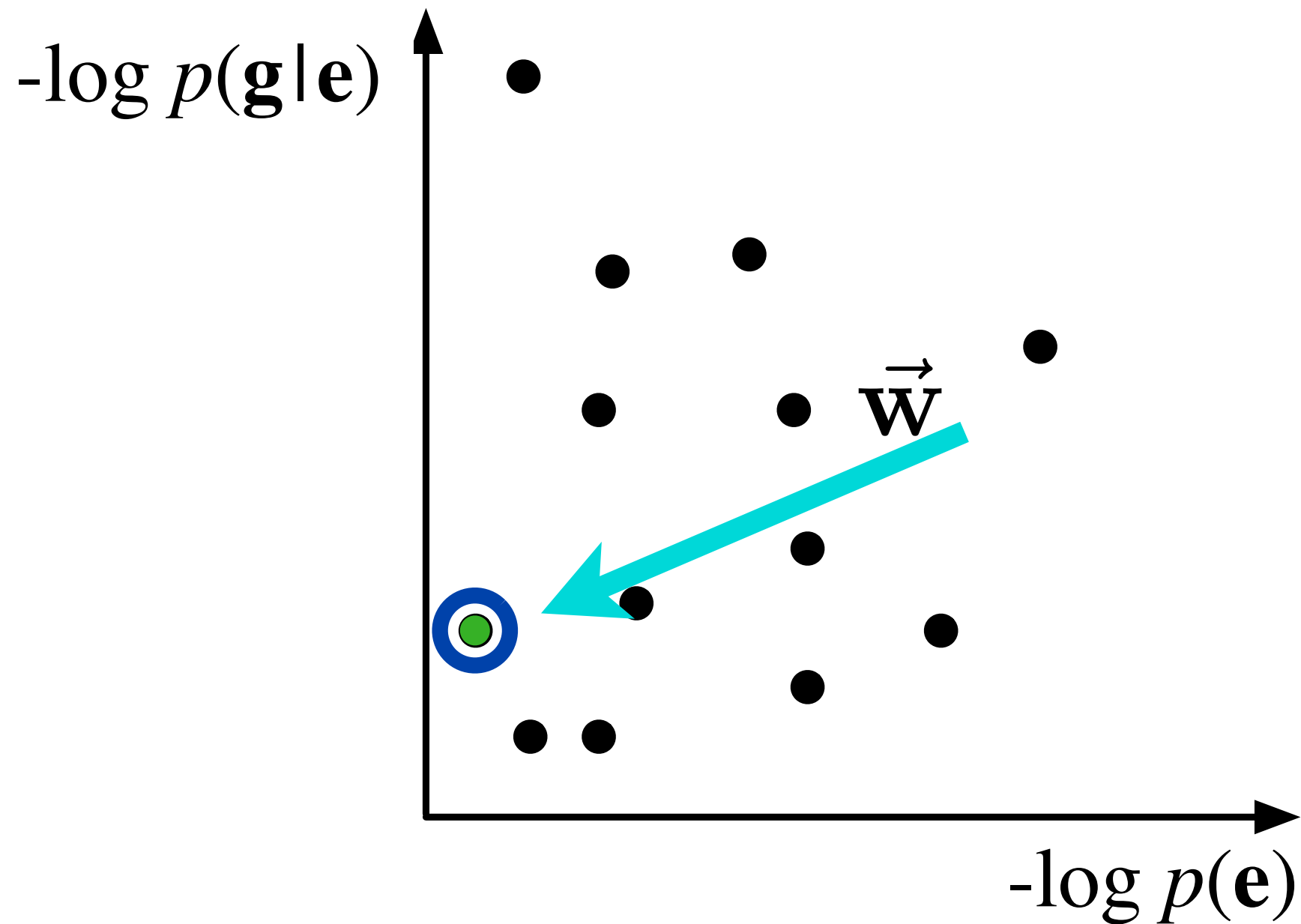
As a Linear Model



As a Linear Model



As a Linear Model

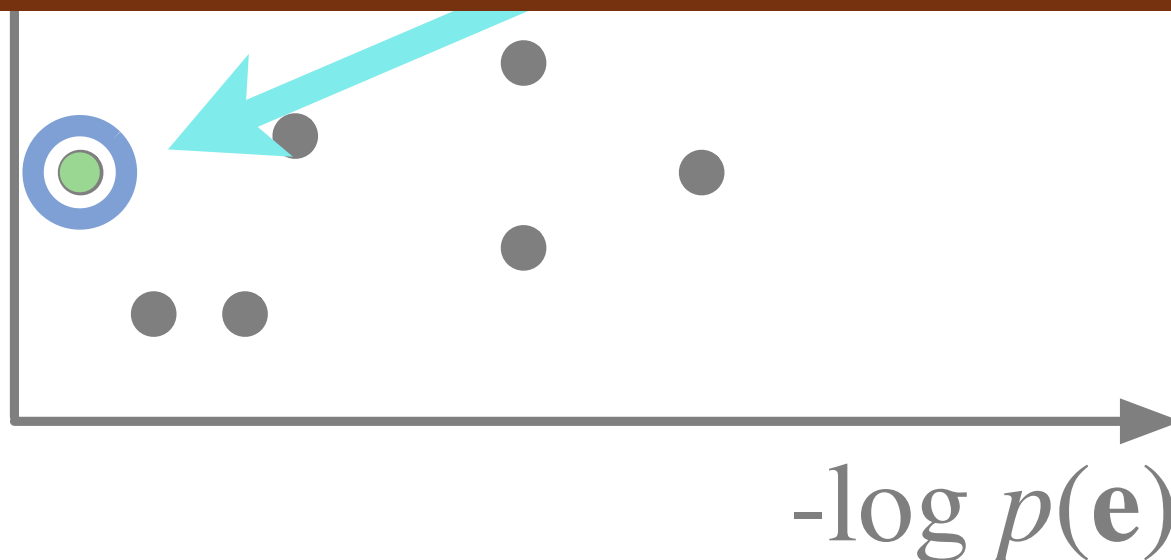


As a Linear Model

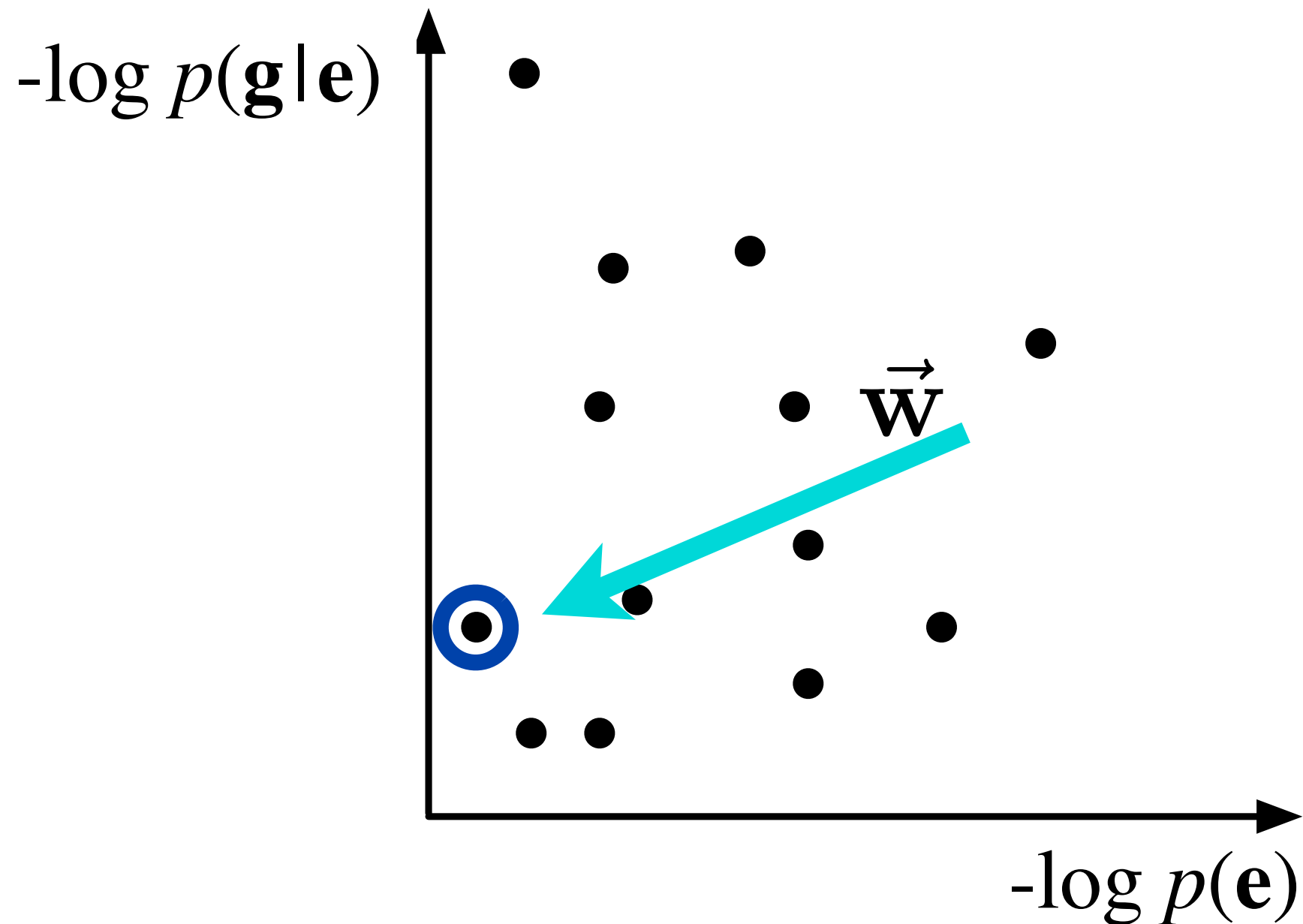
$-\log p(\mathbf{g}|\mathbf{e})$ ↑ •

Improvement 1:

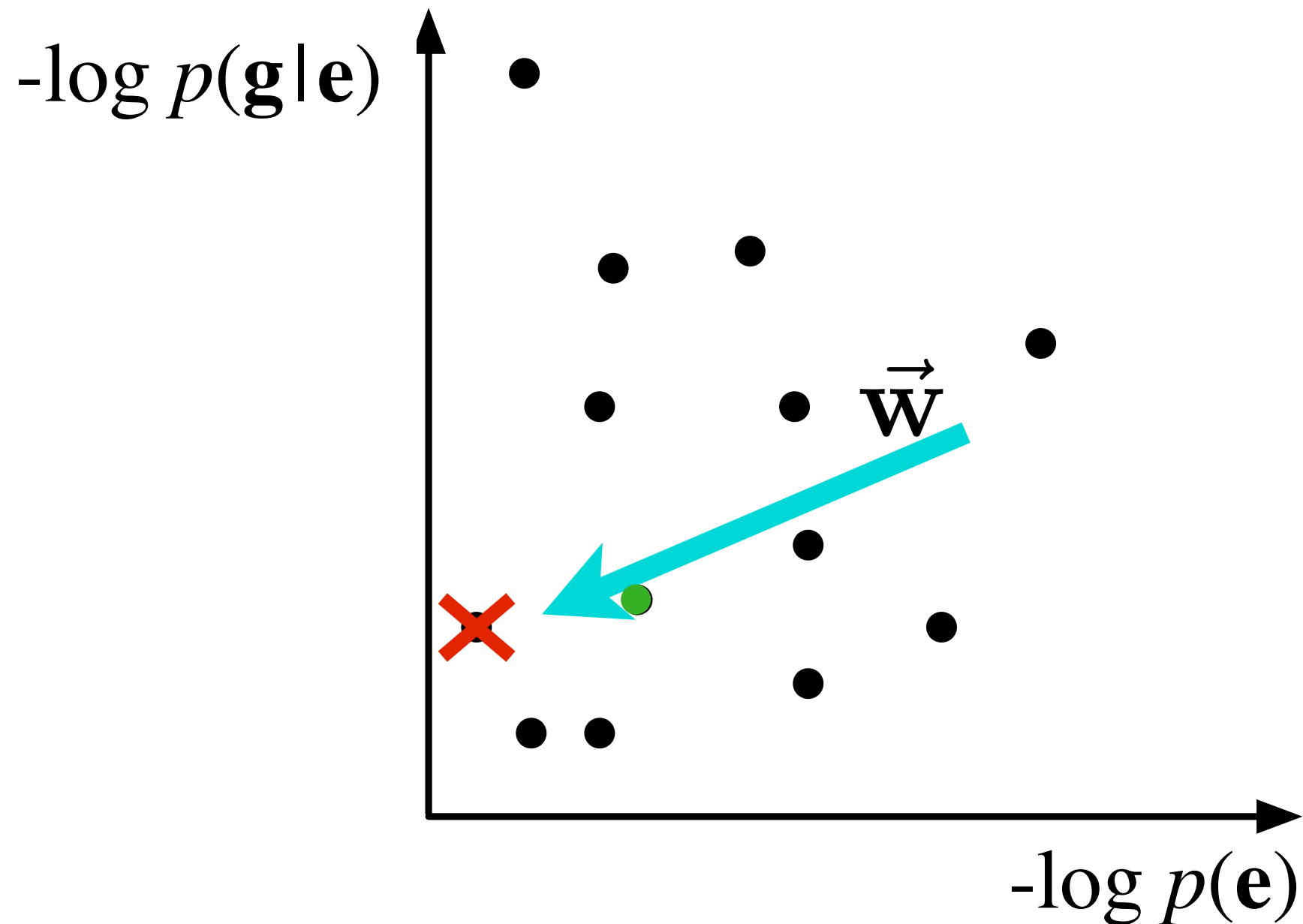
change \vec{w} to find better translations



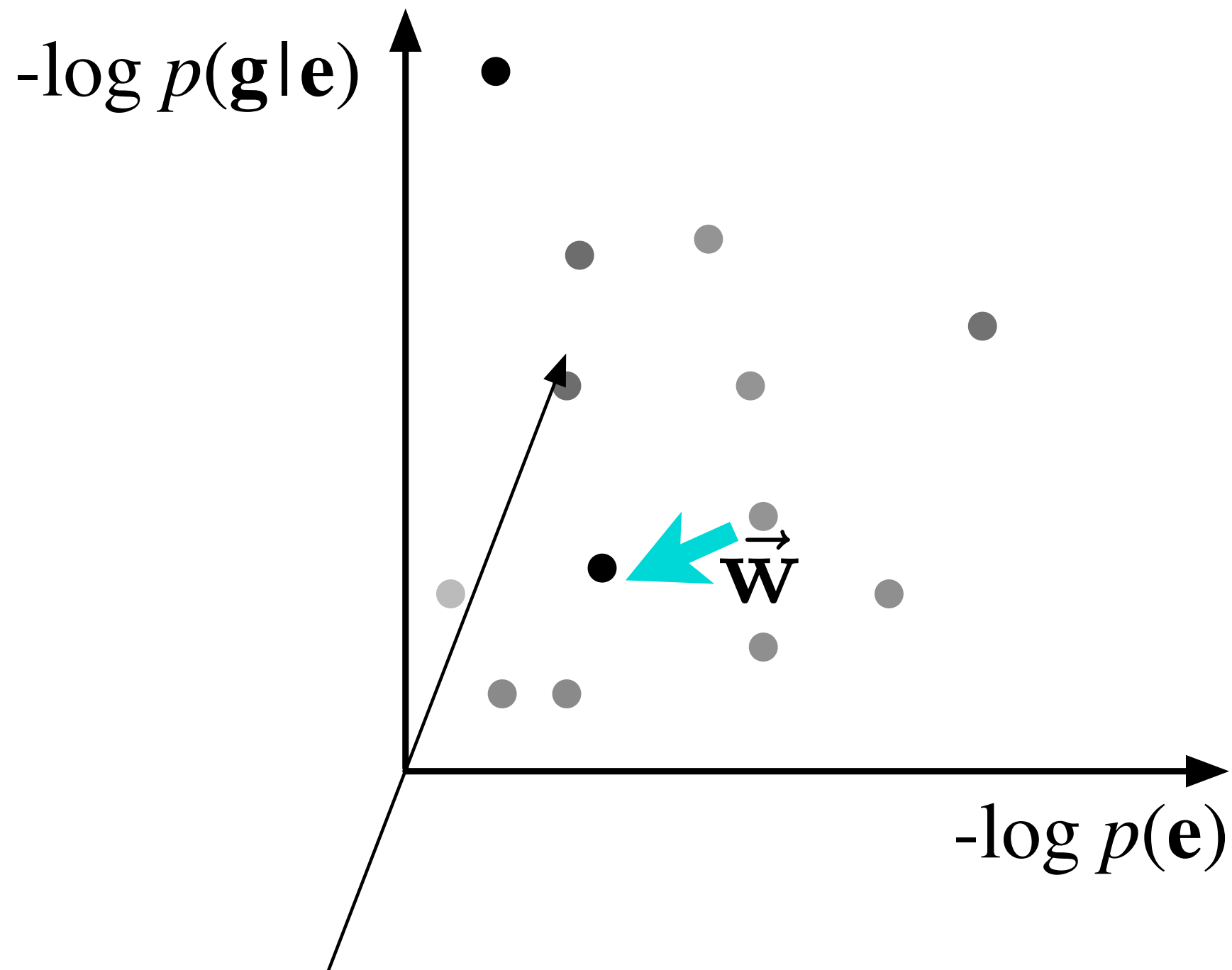
As a Linear Model



As a Linear Model



As a Linear Model

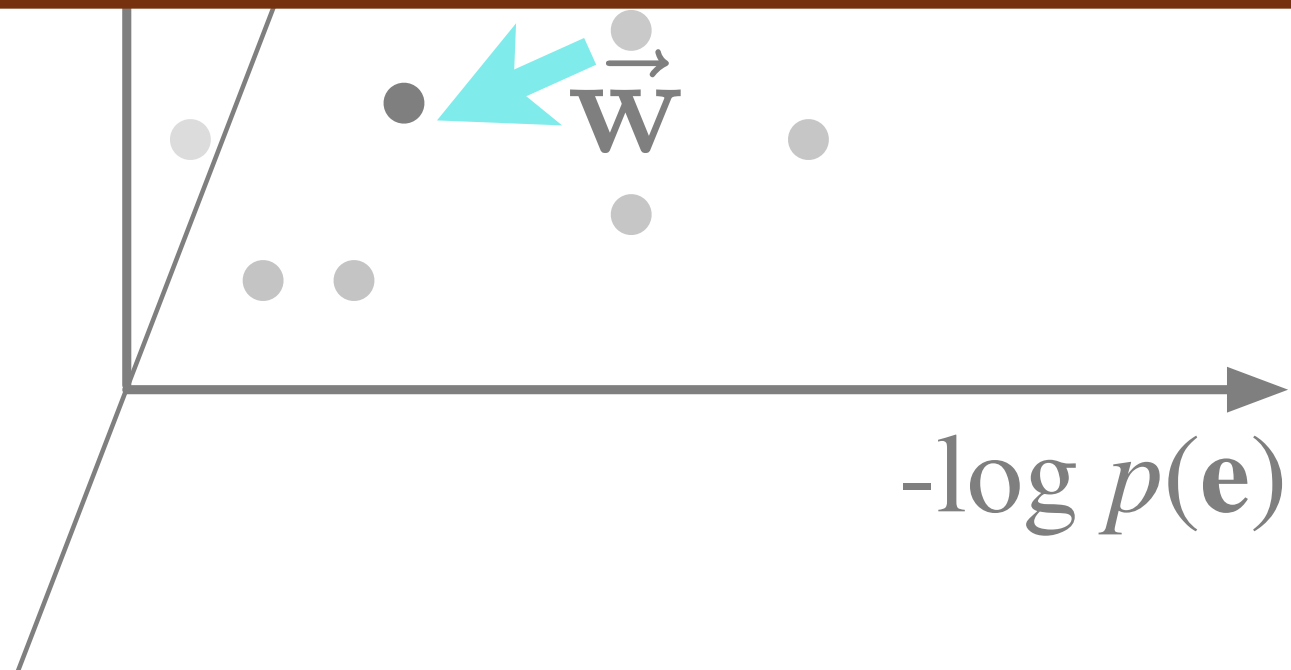


As a Linear Model

$-\log p(\mathbf{g}|\mathbf{e})$ ↑ •

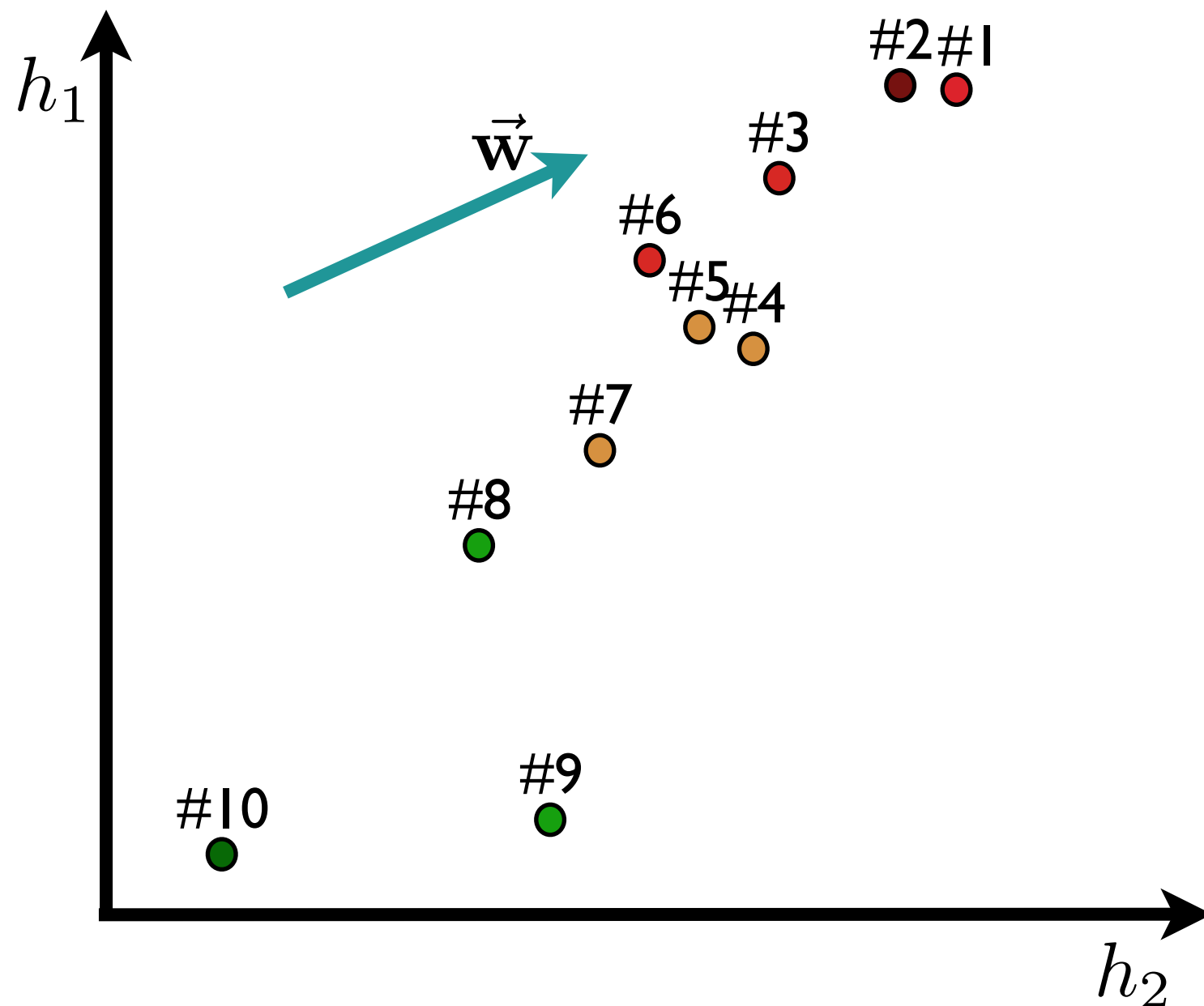
Improvement 2:

Add dimensions to make points separable

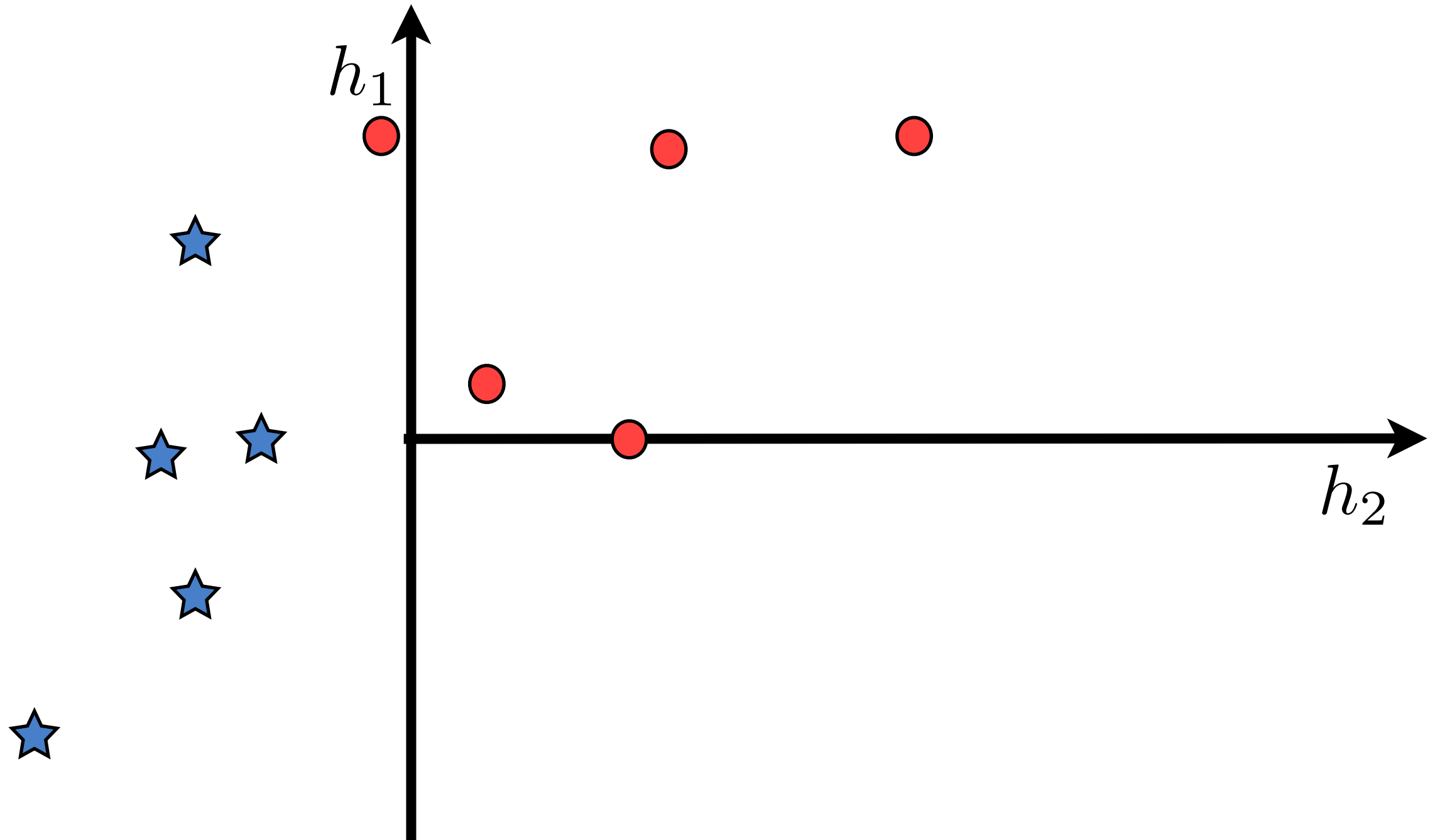


Parameter Learning: Review

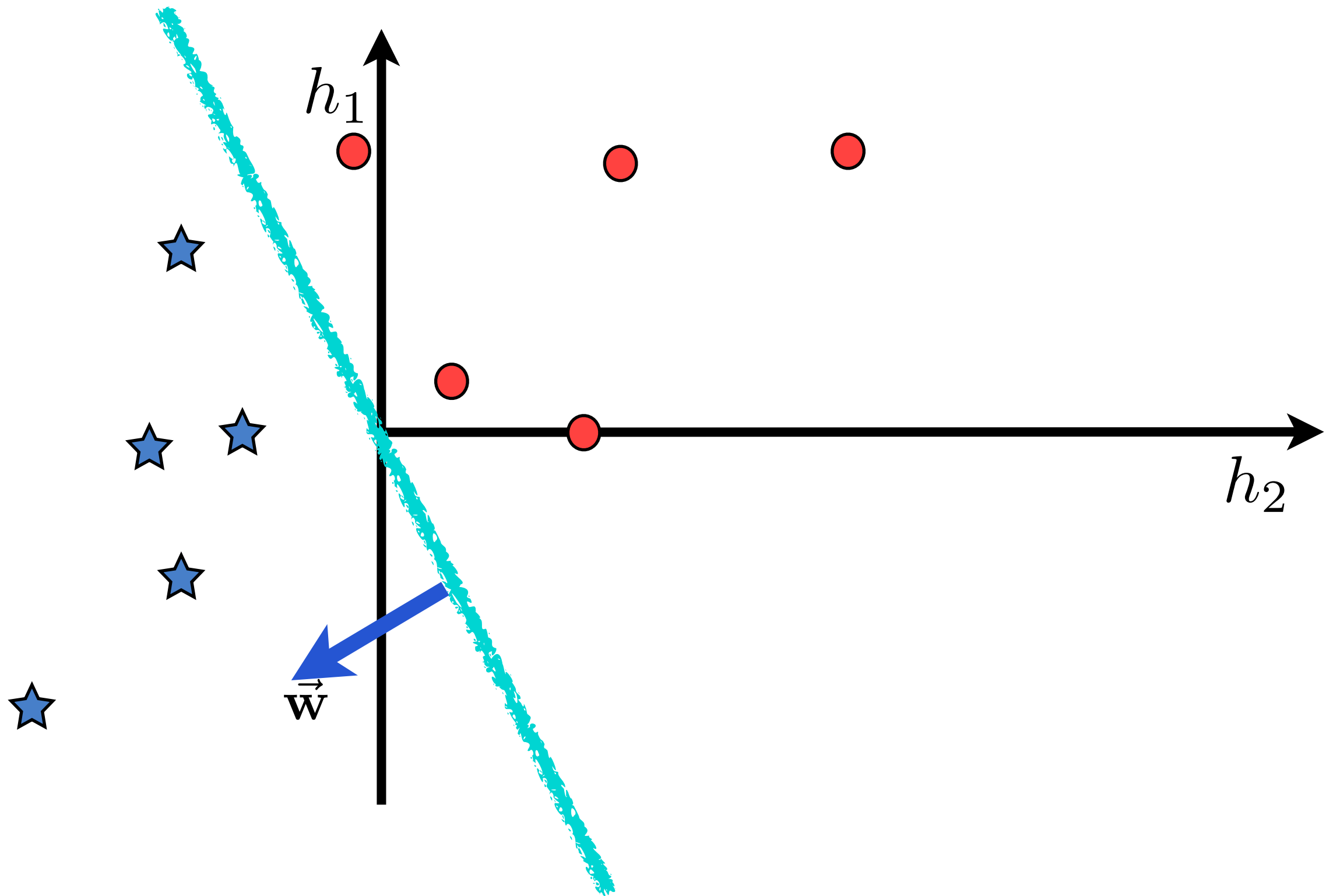
K-Best List Example



- $0.8 \leq \ell < 1.0$
- $0.6 \leq \ell < 0.8$
- $0.4 \leq \ell < 0.6$
- $0.2 \leq \ell < 0.4$
- $0.0 \leq \ell < 0.2$

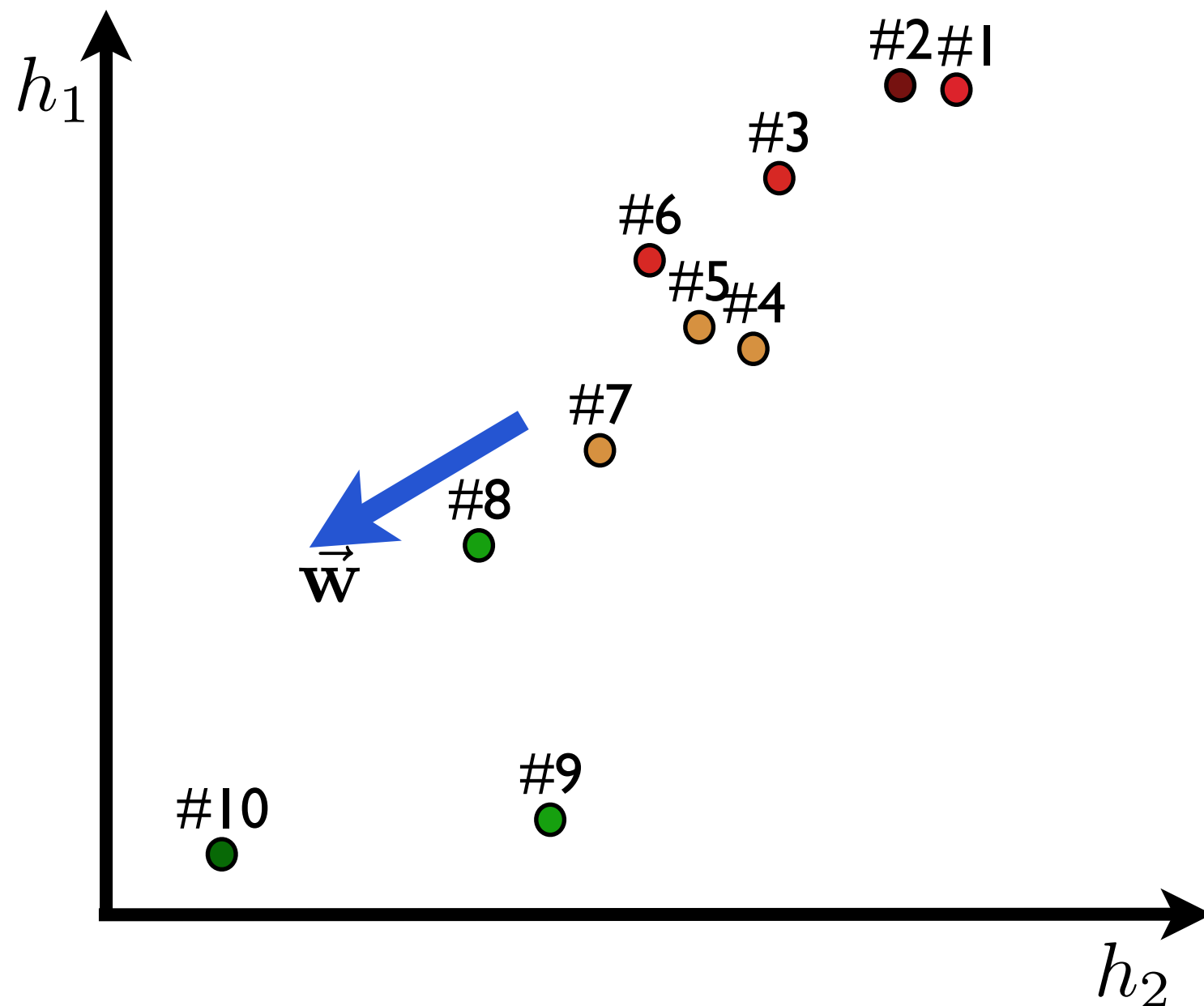


Fit a linear model



Fit a linear model

K-Best List Example



Limitations

- We can't optimize corpus-level metrics, like BLEU, on a test set
 - These don't decompose by sentence!
- We turn now to a kind of “direct cost minimization”

MERT



- **Minimum Error Rate Training**
- Directly target an automatic evaluation metric
 - BLEU is defined at the corpus level
 - MERT optimizes at the corpus level
- **Downsides**
 - Does not deal well with $> \sim 20$ features

MERT

Given weight vector \mathbf{w} , any hypothesis $\langle \mathbf{e}, \mathbf{a} \rangle$ will have a (scalar) score $m = \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$

Now pick a **search vector** \mathbf{v} , and consider how the score of this hypothesis will change:

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \gamma \mathbf{v}$$

MERT

Given weight vector \mathbf{w} , any hypothesis $\langle \mathbf{e}, \mathbf{a} \rangle$ will have a (scalar) score $m = \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$

Now pick a **search vector** \mathbf{v} , and consider how the score of this hypothesis will change:

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \gamma \mathbf{v}$$

$$m = (\mathbf{w} + \gamma \mathbf{v})^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$$

MERT

Given weight vector \mathbf{w} , any hypothesis $\langle \mathbf{e}, \mathbf{a} \rangle$ will have a (scalar) score $m = \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$

Now pick a **search vector** \mathbf{v} , and consider how the score of this hypothesis will change:

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \gamma \mathbf{v}$$

$$\begin{aligned} m &= (\mathbf{w} + \gamma \mathbf{v})^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a}) \\ &= \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a}) + \gamma \mathbf{v}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a}) \end{aligned}$$

MERT

Given weight vector \mathbf{w} , any hypothesis $\langle \mathbf{e}, \mathbf{a} \rangle$ will have a (scalar) score $m = \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$

Now pick a **search vector** \mathbf{v} , and consider how the score of this hypothesis will change:

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \gamma \mathbf{v}$$

$$\begin{aligned} m &= (\mathbf{w} + \gamma \mathbf{v})^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a}) \\ &= \underbrace{\mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})}_b + \gamma \underbrace{\mathbf{v}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})}_a \end{aligned}$$

$$m = a\gamma + b$$

MERT

Given weight vector \mathbf{w} , any hypothesis $\langle \mathbf{e}, \mathbf{a} \rangle$ will have a (scalar) score $m = \mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})$

Now pick a **search vector** \mathbf{v} , and consider how the score of this hypothesis will change:

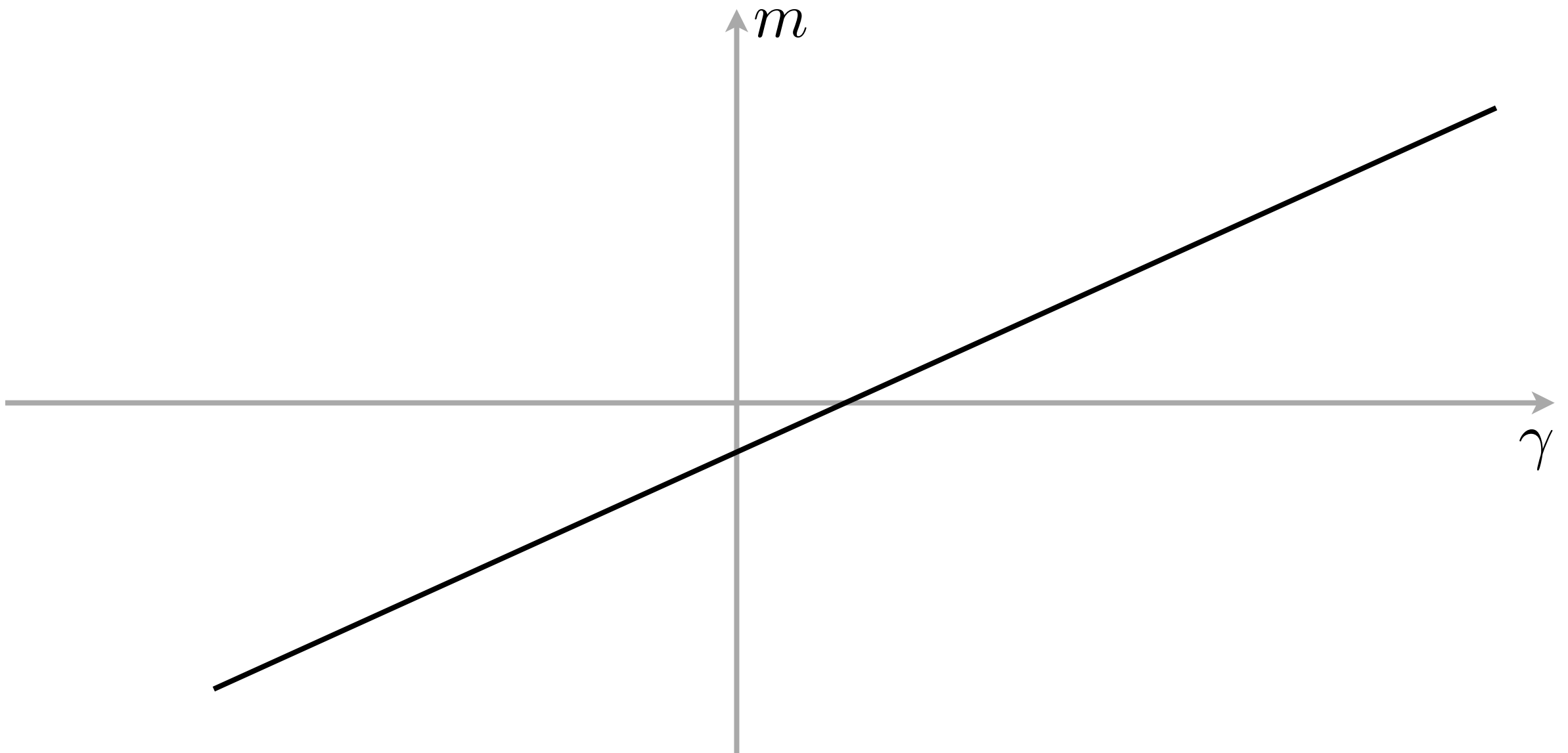
$$\mathbf{w}_{\text{new}} = \mathbf{w} + \gamma \mathbf{v}$$

$$\begin{aligned} m &= (\mathbf{w} + \gamma \mathbf{v})^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a}) \\ &= \underbrace{\mathbf{w}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})}_b + \gamma \underbrace{\mathbf{v}^\top \mathbf{h}(\mathbf{g}, \mathbf{e}, \mathbf{a})}_a \end{aligned}$$

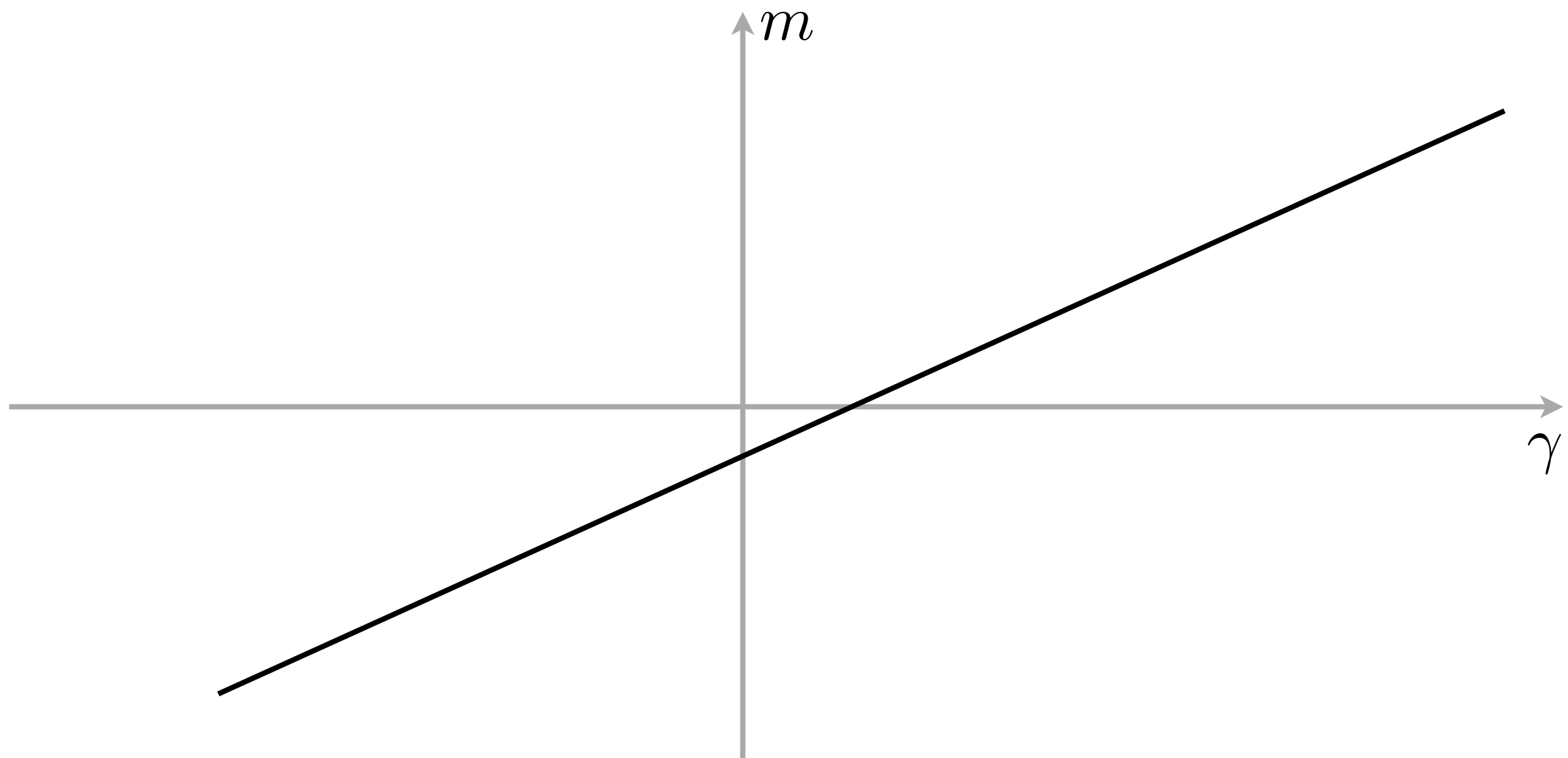
$$m = a\gamma + b$$

Linear function in 2D!

MERT

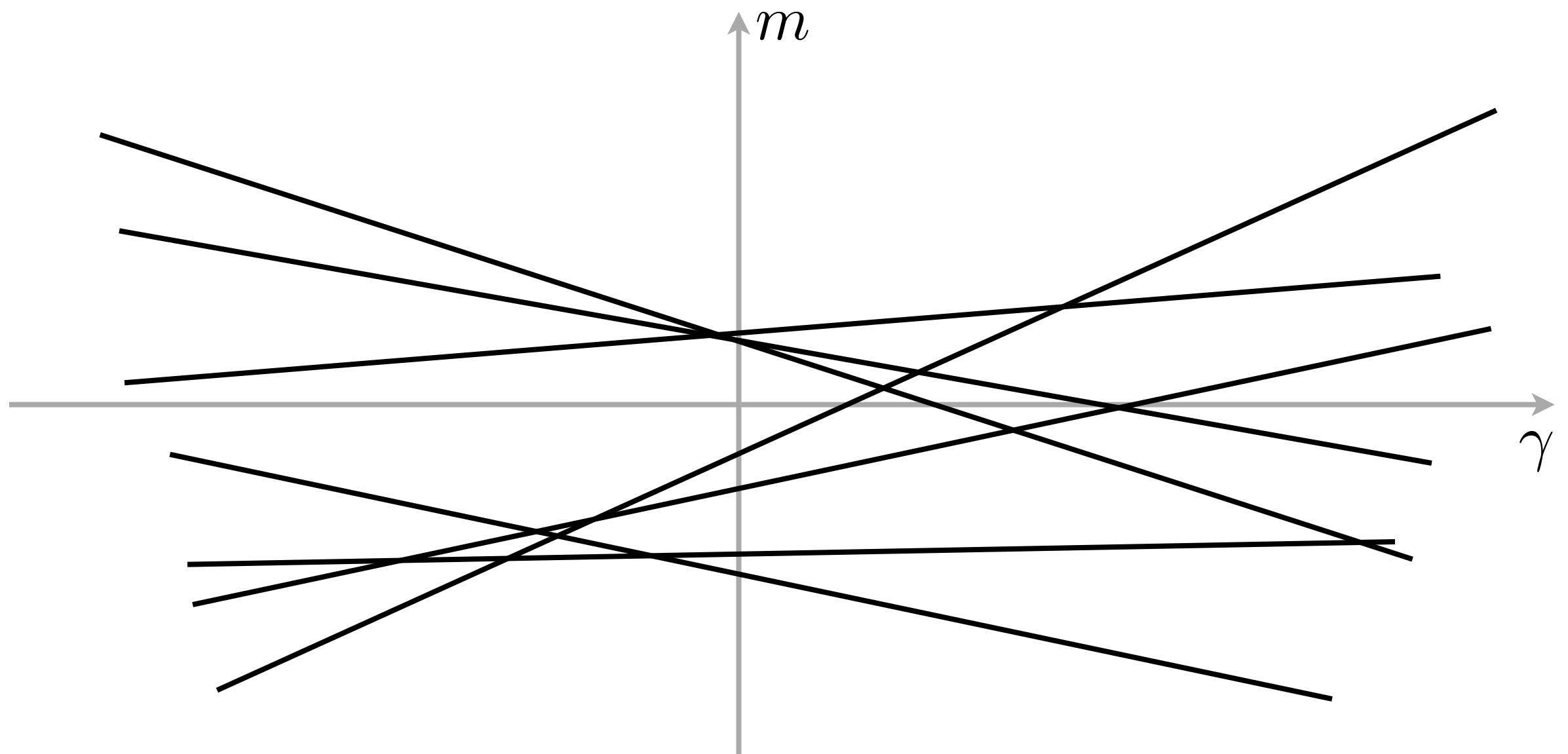


MERT



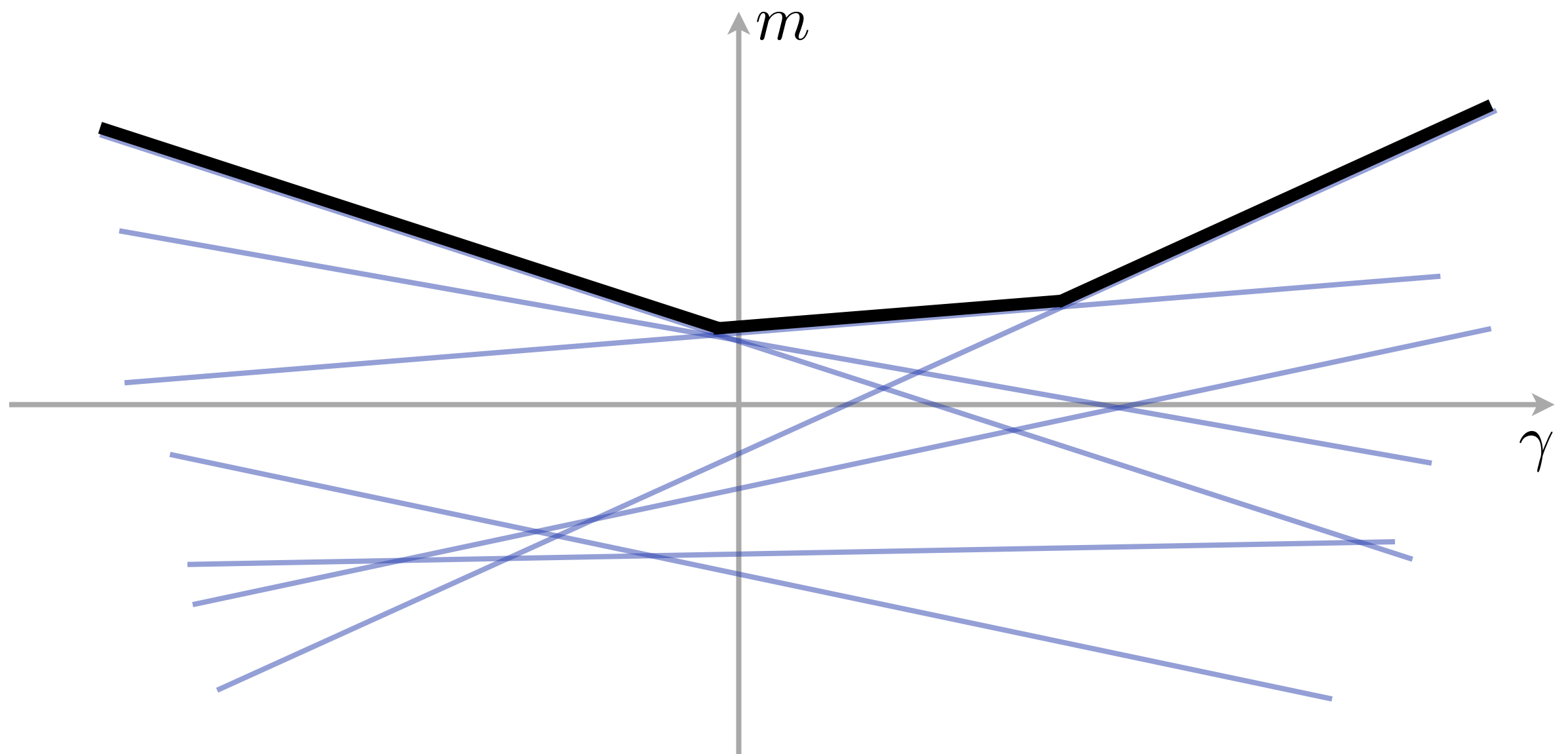
Recall our k-best set $\{\langle \mathbf{e}_i^*, \mathbf{a}_i^* \rangle\}_{i=1}^K$

MERT

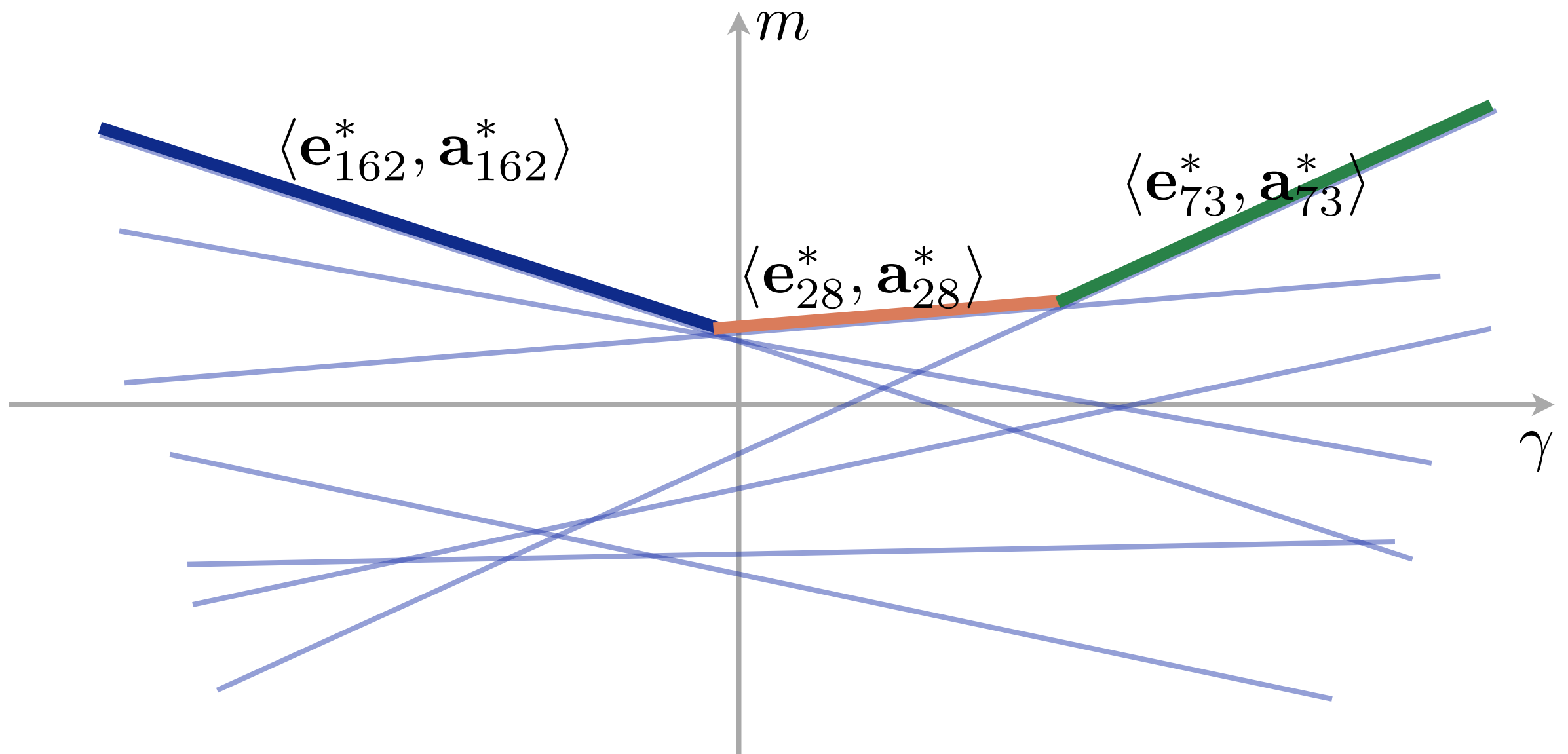


Recall our k-best set $\{\langle \mathbf{e}_i^*, \mathbf{a}_i^* \rangle\}_{i=1}^K$

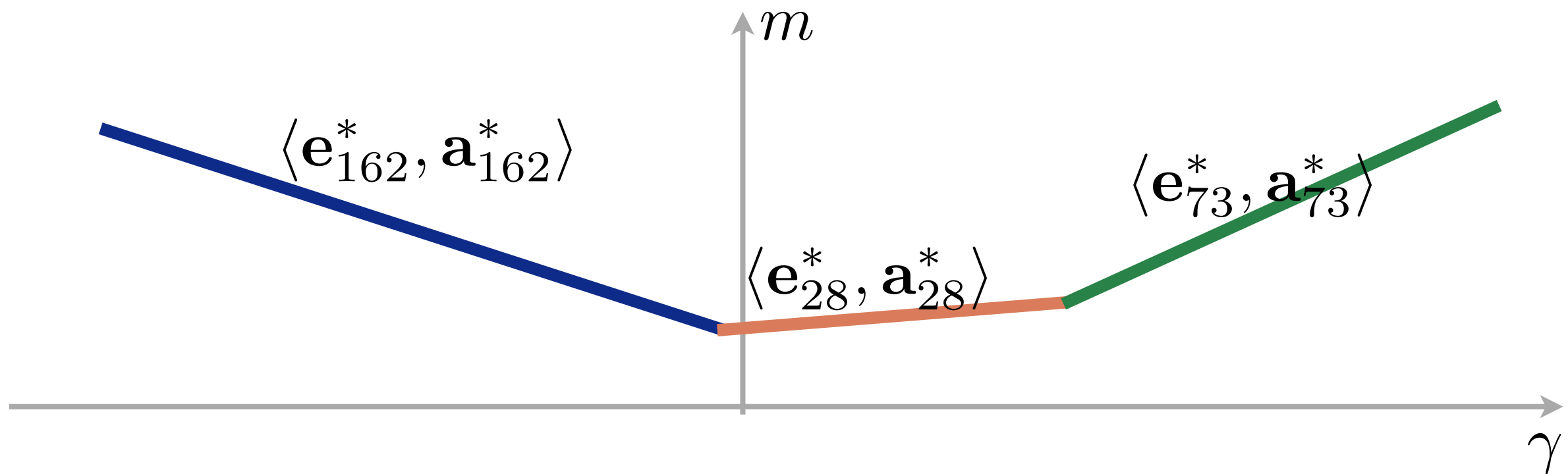
MERT



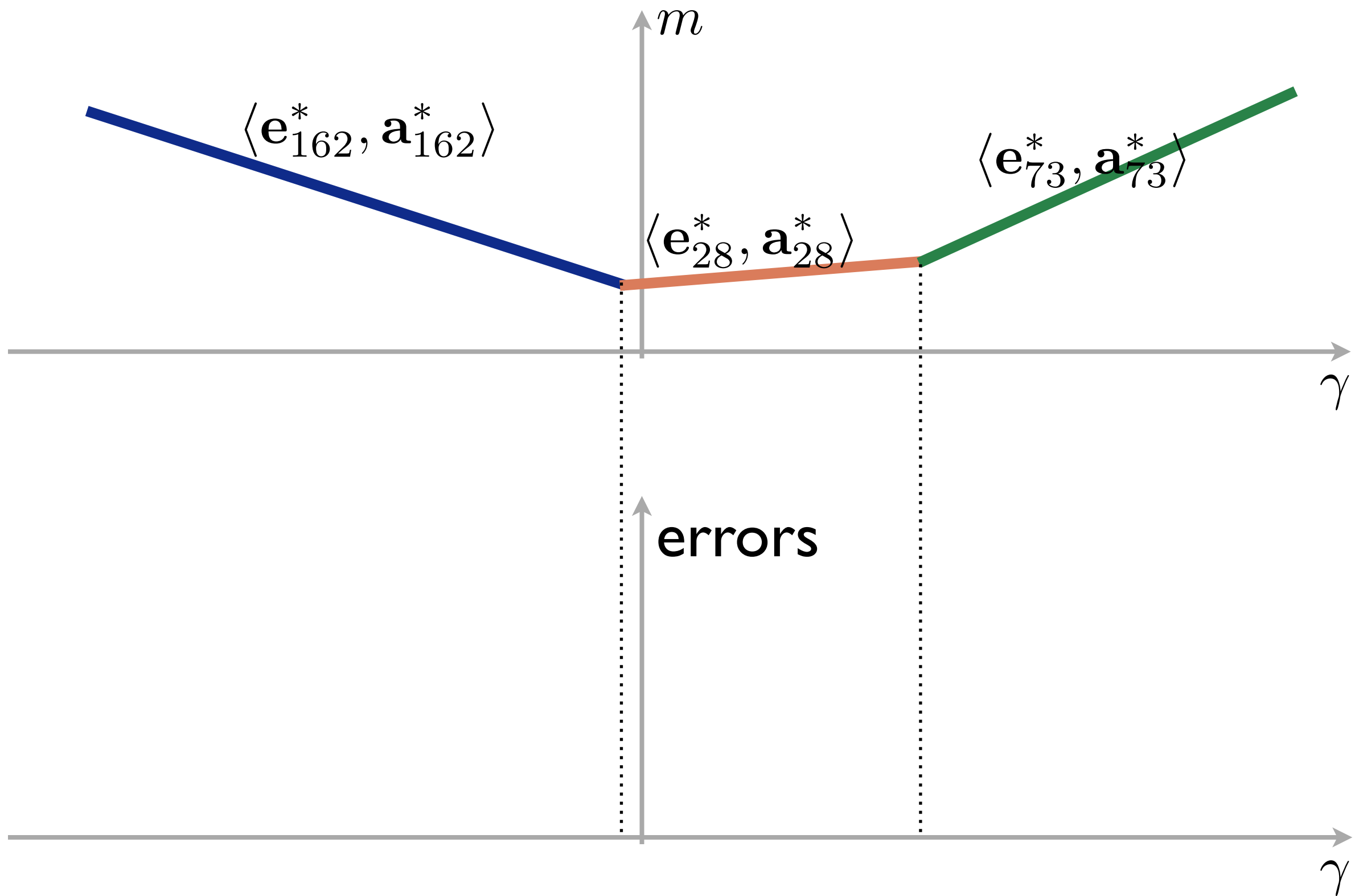
MERT



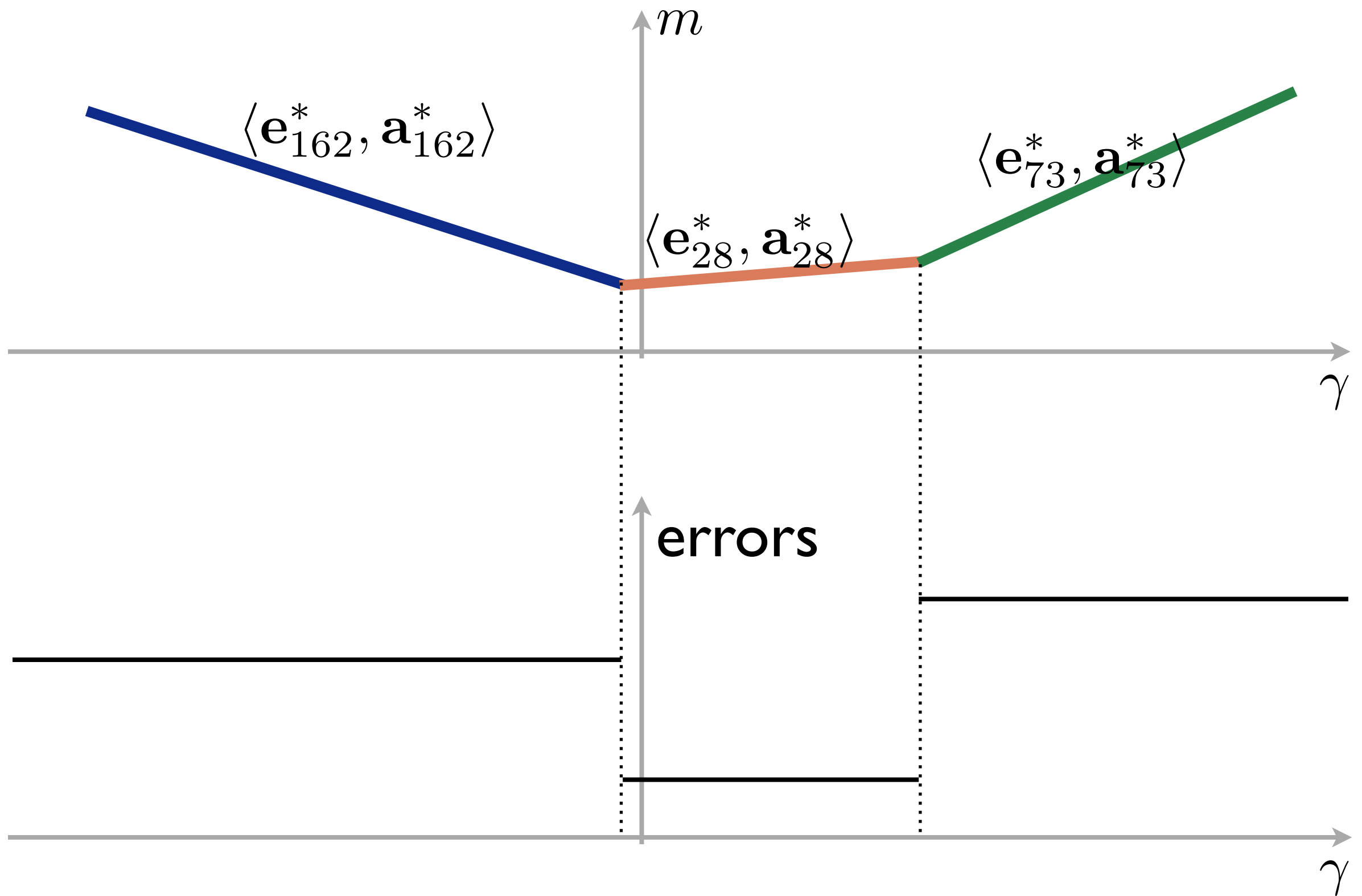
MERT



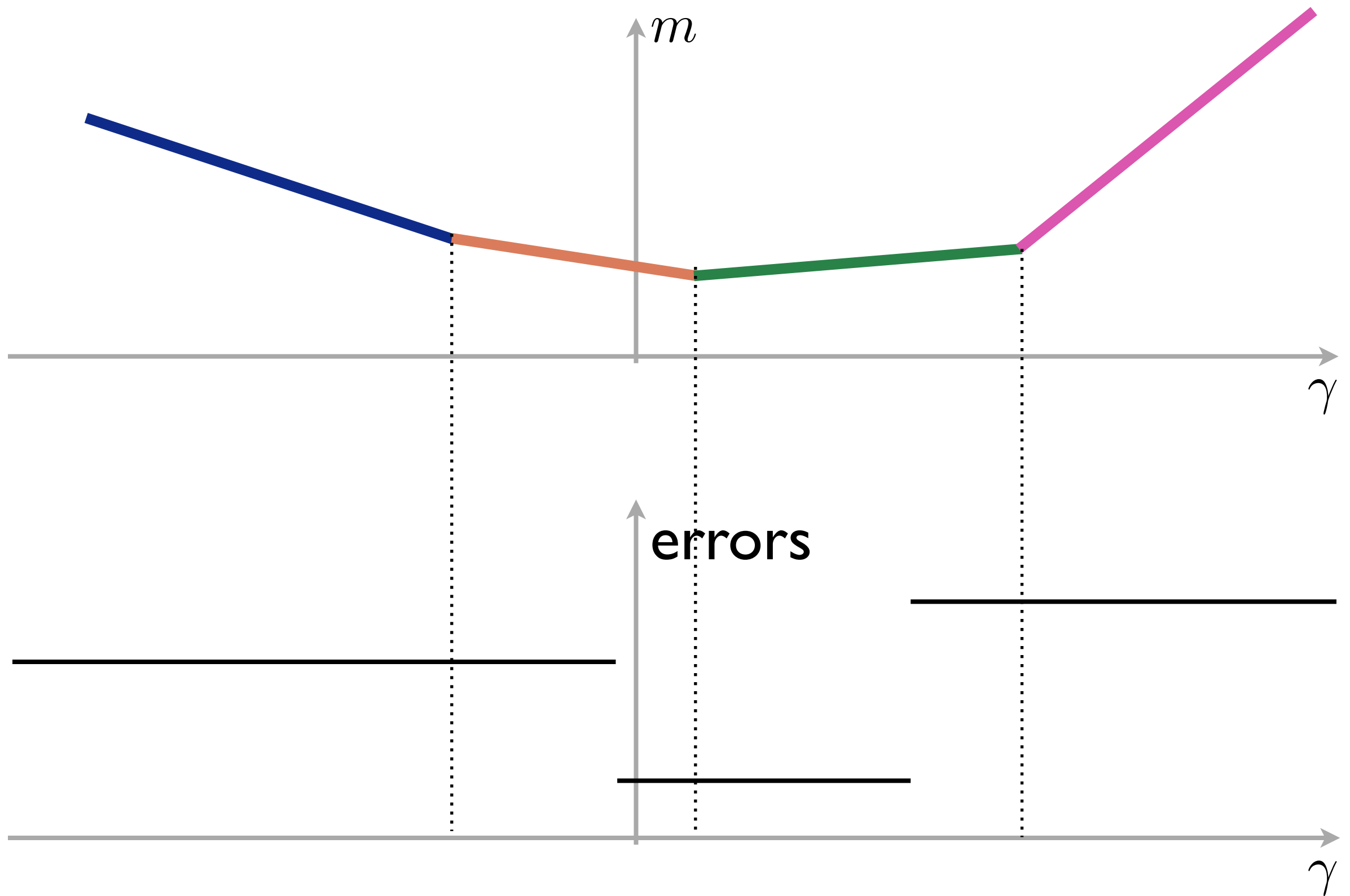
MERT



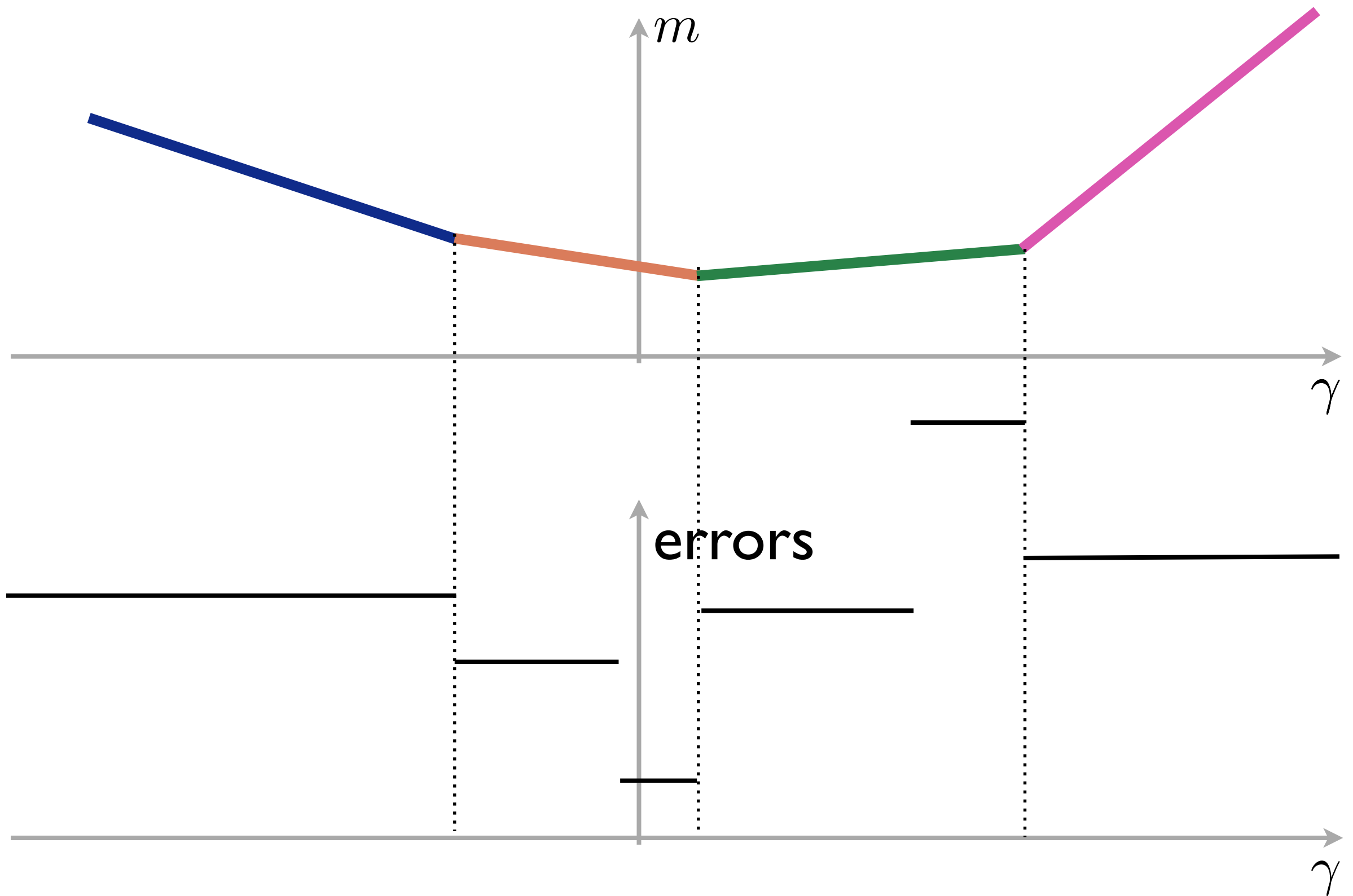
MERT

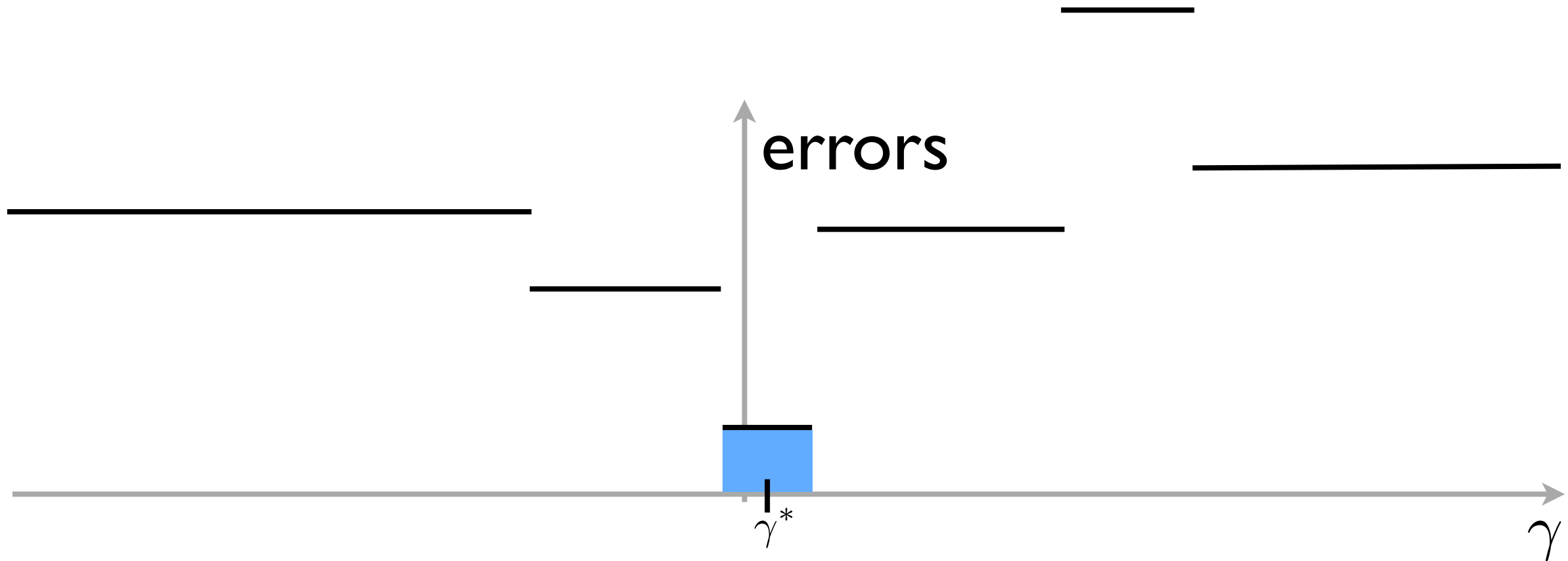


MERT



MERT





Let $\mathbf{w}_{\text{new}} = \gamma^* \mathbf{v} + \mathbf{w}$

MERT

- In practice “errors” are sufficient statistics for evaluation metrics (e.g., BLEU)
- Can maximize or minimize
- How do you pick the search direction?

Dynamic Programming

MERT

Other Algorithms

- Given a hypergraph translation space
- In the Viterbi (Inside) algorithm, there are two operations
 - **Multiplication** (extend path)
 - **Maximization** (choose between paths)
- **Semirings** generalize these to compute other quantities

Semirings

semiring	\mathbb{K}	\oplus	\otimes	$\bar{0}$	$\bar{1}$	notes
Boolean	$\{0,1\}$	\vee	\wedge	0	1	idempotent
count	$\mathbb{N}_0 \cup \{\infty\}$	$+$	\times	0	1	
probability	$\mathbb{R}_+ \cup \{\infty\}$	$+$	\times	0	1	
tropical	$\mathbb{R} \cup \{-\infty, \infty\}$	\max	$+$	$-\infty$	0	idempotent
log	$\mathbb{R} \cup \{-\infty, \infty\}$	\oplus_{\log}	$+$	$-\infty$	0	

Inside Algorithm

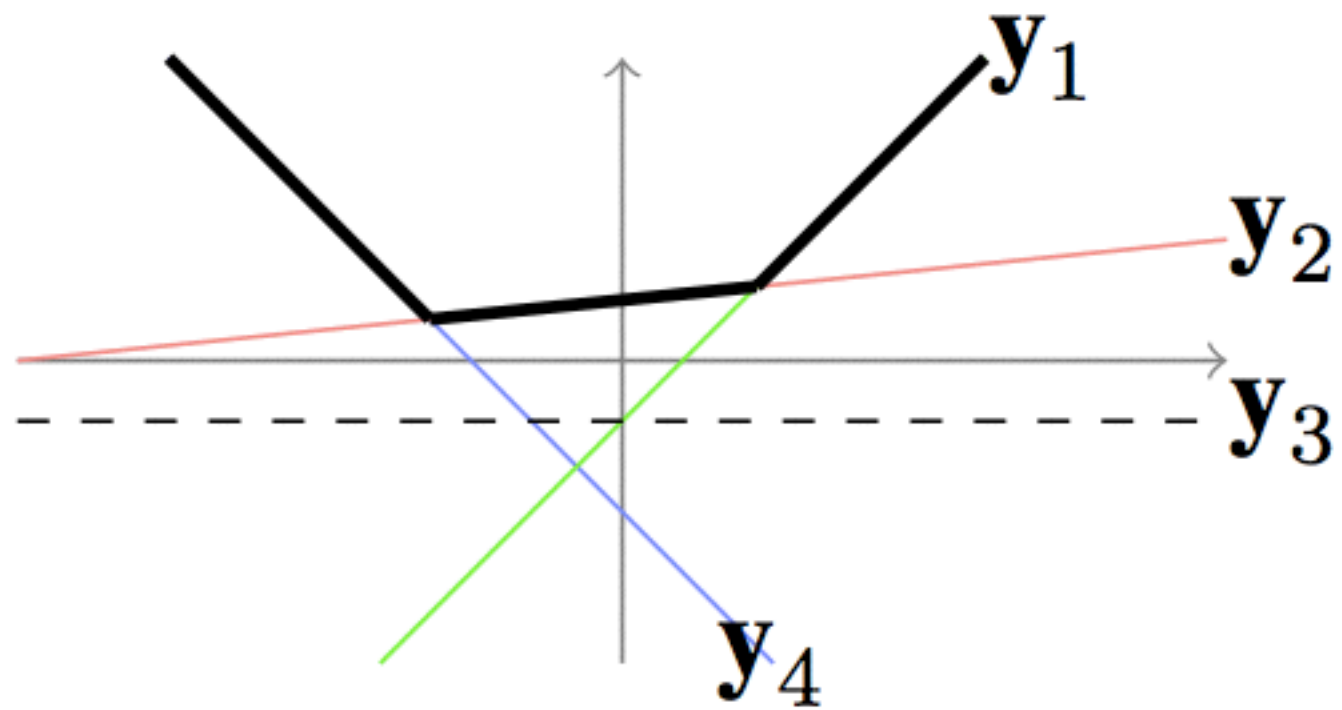
$$\alpha(q_{goal}) = \bigoplus_{\mathbf{d} \in \mathcal{G}} \bigotimes_{e \in \mathbf{d}} w(e)$$

```
1: function INSIDE( $\mathcal{G}, K$ )                                ▷  $\mathcal{G}$  is an acyclic hypergraph and  $K$  is a semiring
2:   for  $q$  in topological order in  $\mathcal{G}$  do
3:     if  $B(q) = \emptyset$  then
4:        $\alpha(q) \leftarrow \bar{1}$                                 ▷ assume states with no in-edges are axioms
5:     else
6:        $\alpha(q) \leftarrow \bar{0}$ 
7:       for all  $e \in B(q)$  do                                ▷ all in-coming edges to node  $q$ 
8:          $k \leftarrow w(e)$ 
9:         for all  $r \in \mathbf{t}(e)$  do                                ▷ all tail (previous) nodes of edge  $e$ 
10:           $k \leftarrow k \otimes \alpha(r)$ 
11:           $\alpha(q) \leftarrow \alpha(q) \oplus k$ 
12:   return  $\alpha$ 
```

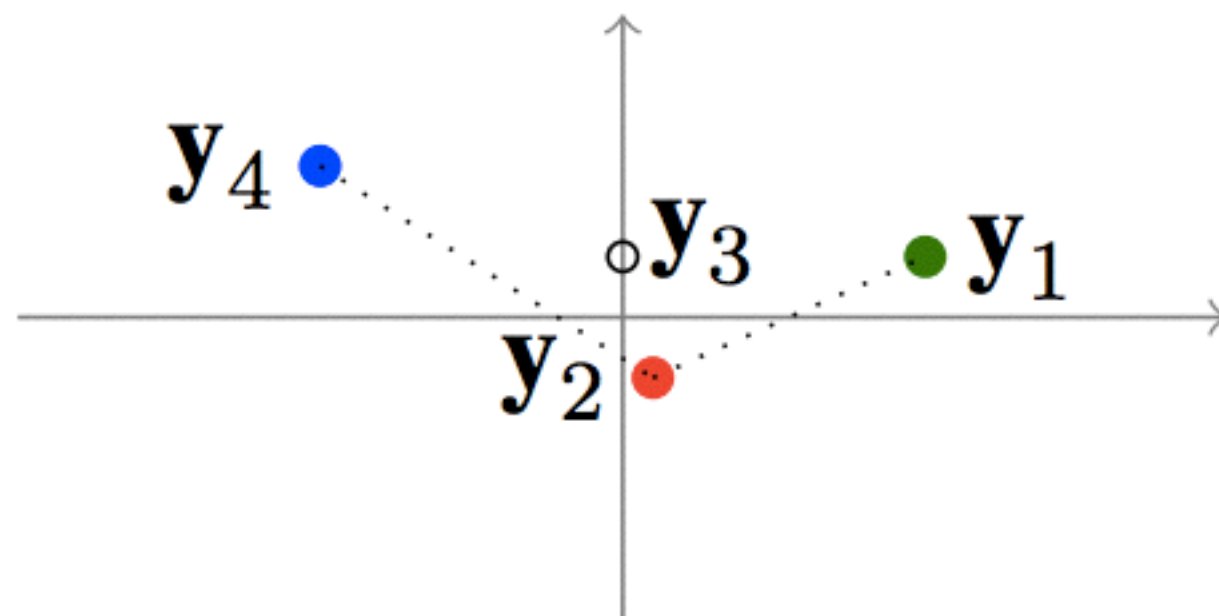
Point-Line Duality

- Represent a set of lines as a set of points (and vice-versa)
- $y = mx + b \Rightarrow (m, -b)$
- The slope between dual points is the intersection x-axis of the pair of lines
- An upper envelope is dual to a **lower convex hull**

Primal



Dual



Convex Hull Semiring

Definition 2. The Convex Hull Semiring.

Let $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ be defined as follows:

\mathbb{K}	<i>A set of points in the plane that are the extreme points of a convex hull.</i>
$A \oplus B$	$\text{conv}[A \cup B]$
$A \otimes B$	<i>convex hull of the Minkowski sum, i.e.,</i> $\text{conv}\{(a_1 + b_1, a_2 + b_2) \mid (a_1, a_2) \in A \wedge (b_1, b_2) \in B\}$
$\bar{0}$	\emptyset
$\bar{1}$	$\{(0, 0)\}$

Theorem 1. *The Convex Hull Semiring fulfills the semiring axioms and is commutative and idempotent.*

Theorem 2

- The Inside algorithm with the computes the convex hull dual to the MERT upper envelope generated from the ∞ -best list of derivations

Summary

- Evaluation metrics
 - Figure out how well we're doing
 - Figure out if a feature helps
 - **Train your system**
- What's a great way to improve translation?
 - **Improve evaluation!**