# COS234:
# Introduction To Machine Learning

## Prof. Yoram Singer

PRINCETON UNIVERSITY

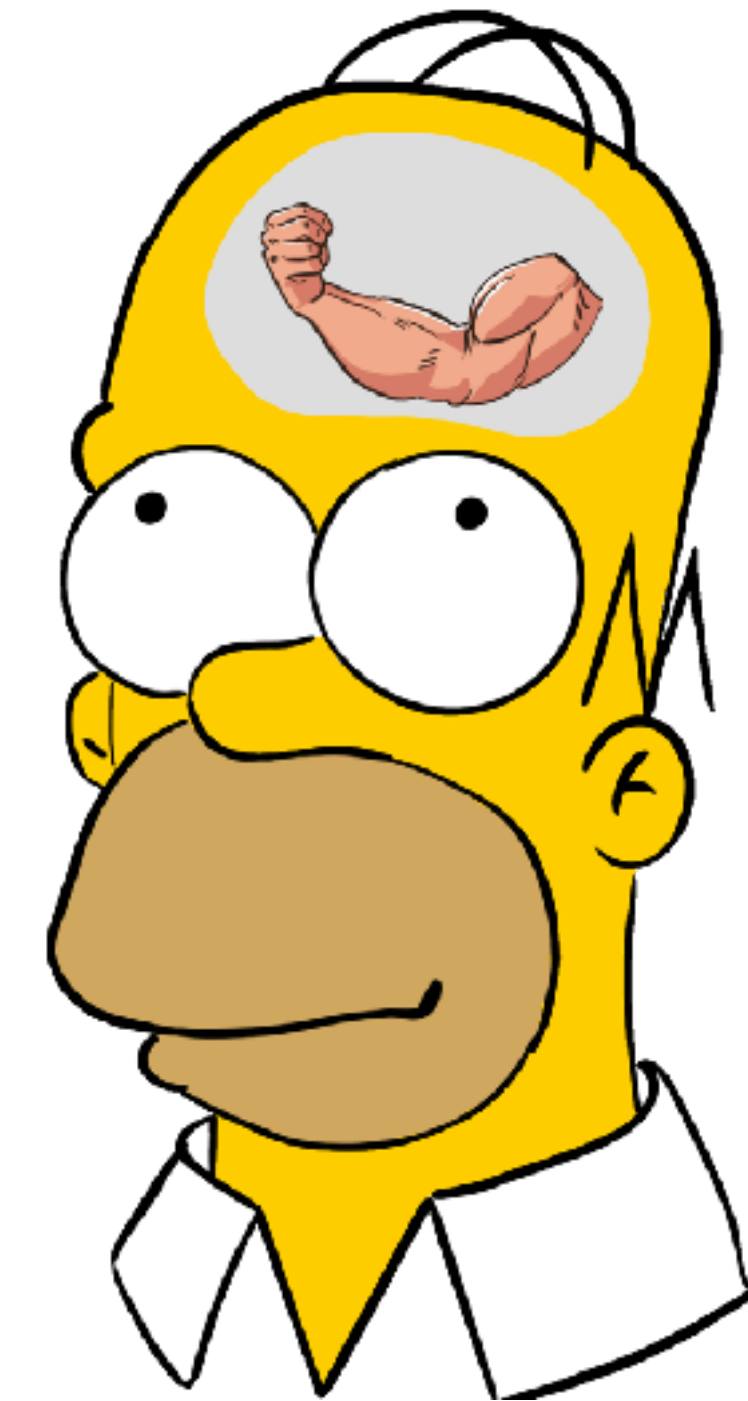## Topic: Multiclass Learning

# Phoneme Classification

2

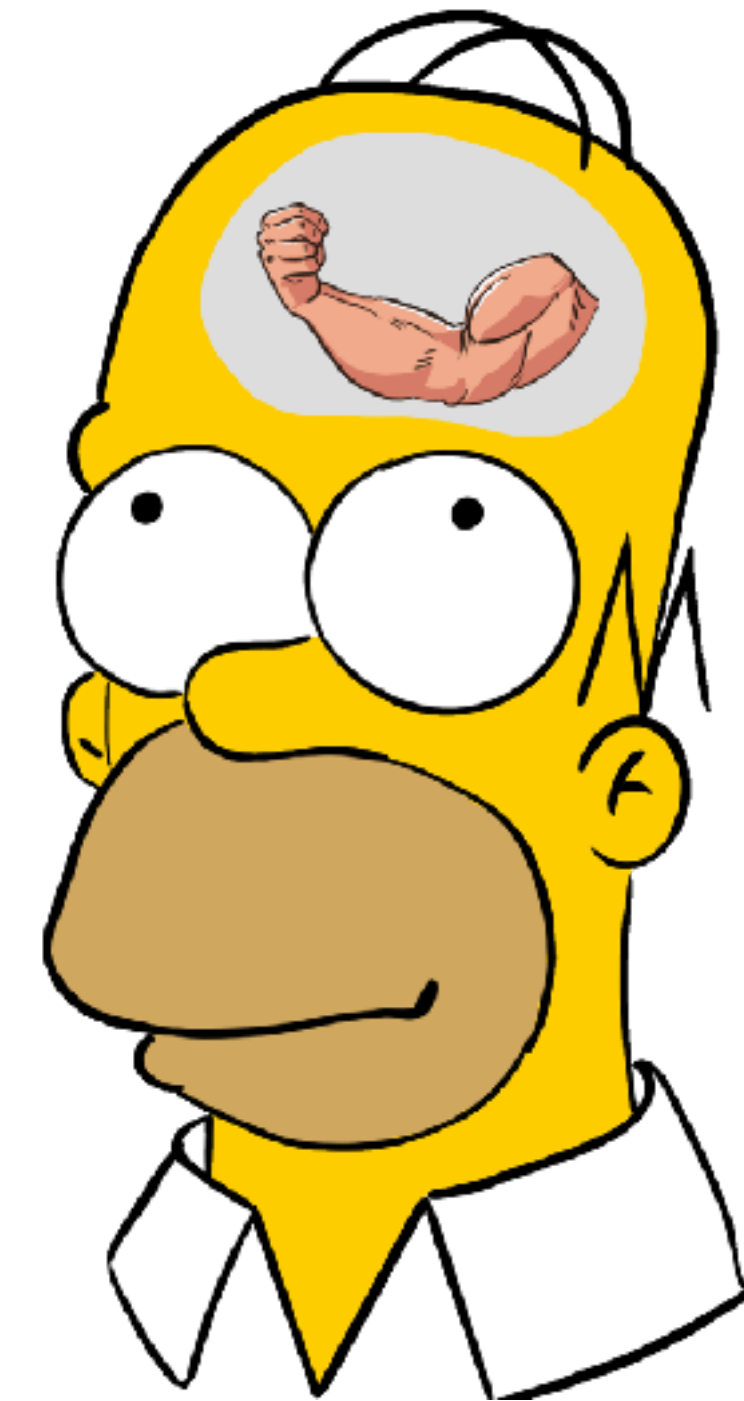# Optical Character Recognition

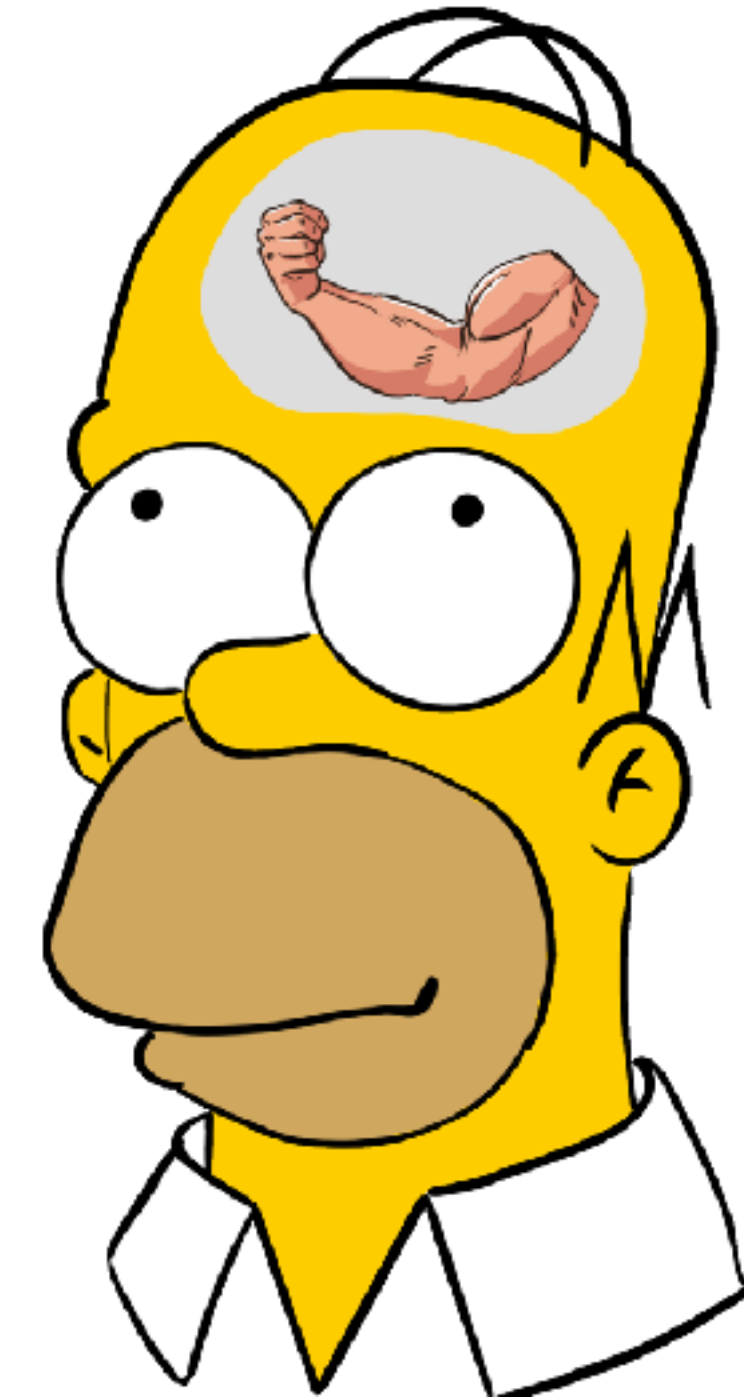# Multiclass: Problem Setting

# Multiclass: Problem Setting
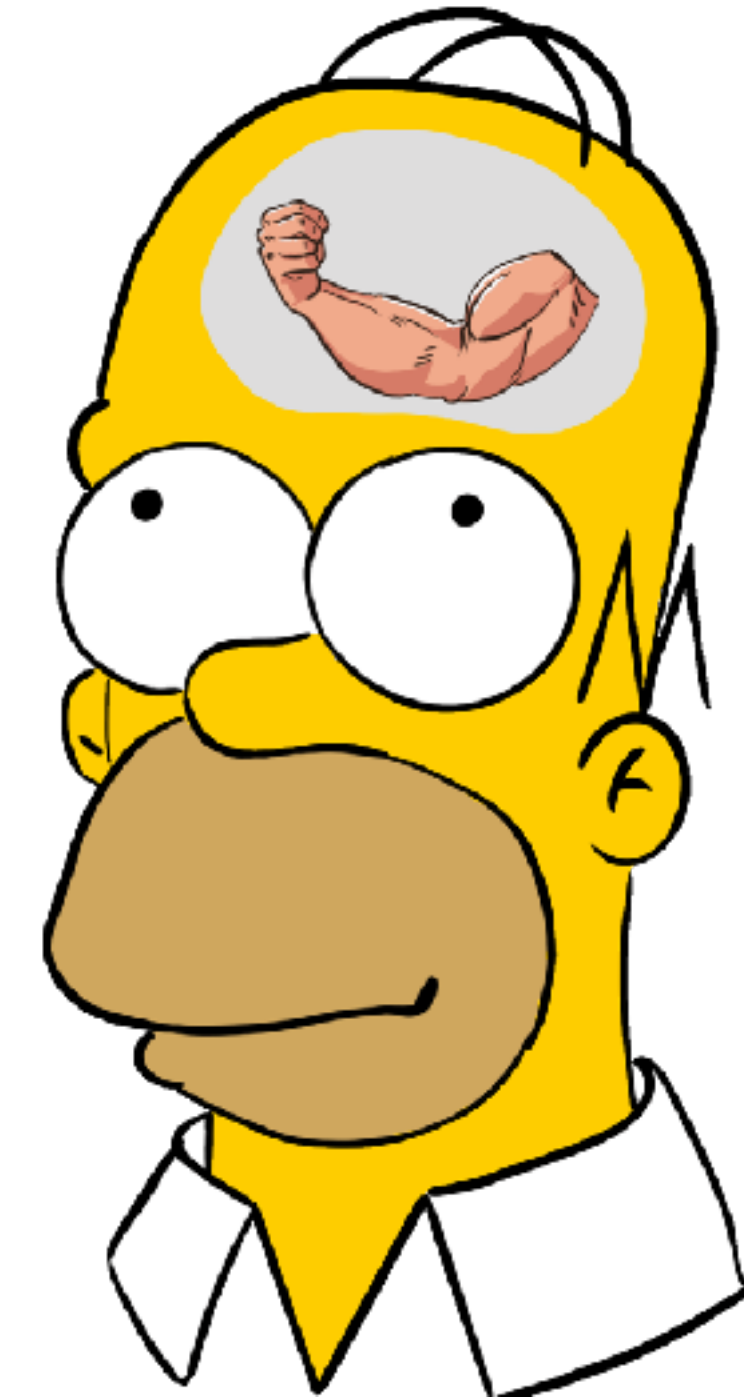
- Instances: $\mathbf{x} \in \mathbf{R}^d$

# Multiclass: Problem Setting

- Instances: $\mathbf{x} \in \mathbf{R}^d$

- Labels: $y \in [k] = \{1, 2, \ldots, k\}$

# Multiclass: Problem Setting

- Instances: $\mathbf{x} \in \mathbf{R}^d$

- Labels: $y \in [k] = \{1, 2, \ldots, k\}$

- Multiclass predictor $h : \mathbf{R}^d \to [k]$

# Multiclass: Problem Setting

- Instances: $\mathbf{x} \in \mathbf{R}^d$

- Labels: $y \in [k] = \{1, 2, \ldots, k\}$

- Multiclass predictor $h : \mathbf{R}^d \to [k]$

- Classification error / mistake: $\mathbf{1}\big[h(\mathbf{x}) \neq y\big]$

# Multiclass: Problem Setting

- Instances: $\mathbf{x} \in \mathbf{R}^d$

- Labels: $y \in [k] = \{1, 2, \ldots, k\}$

- Multiclass predictor $h : \mathbf{R}^d \to [k]$

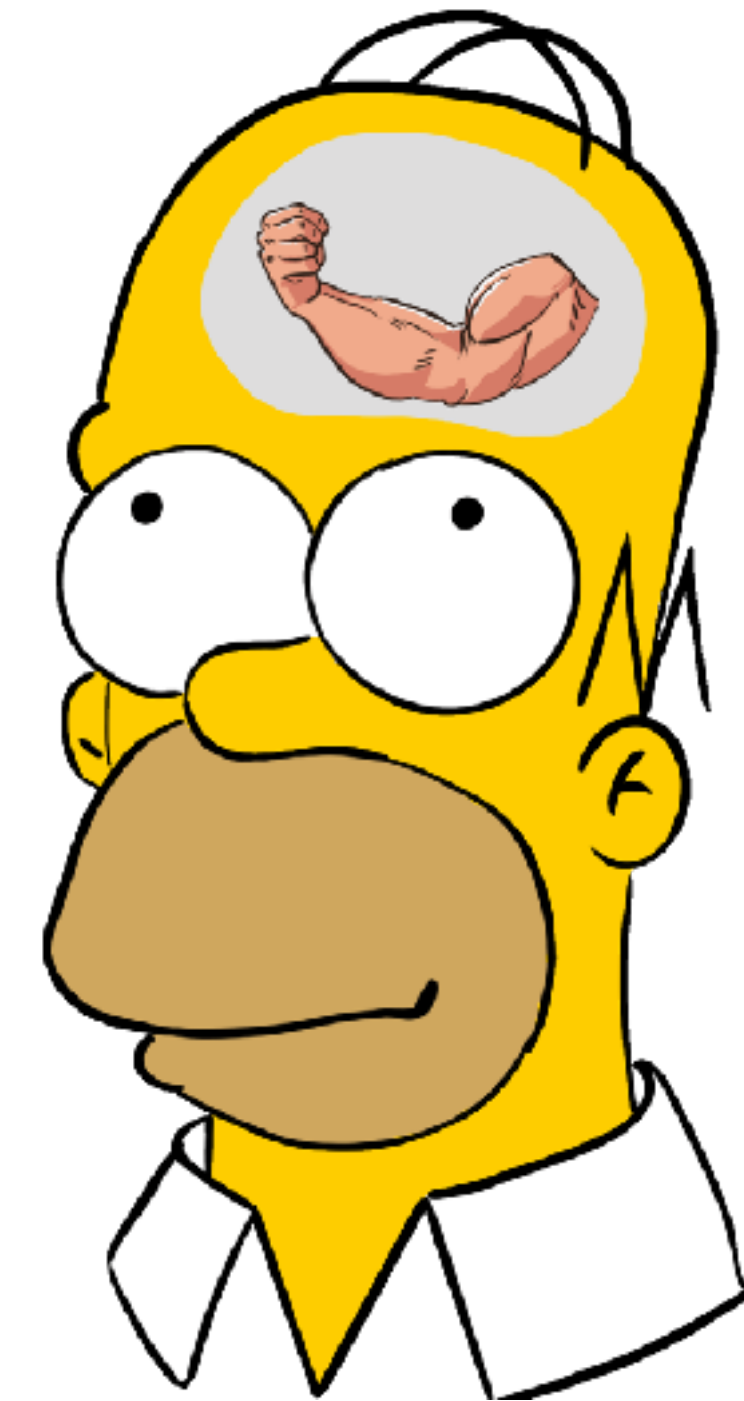- Classification error / mistake: $\mathbf{1}\big[h(\mathbf{x}) \neq y\big]$
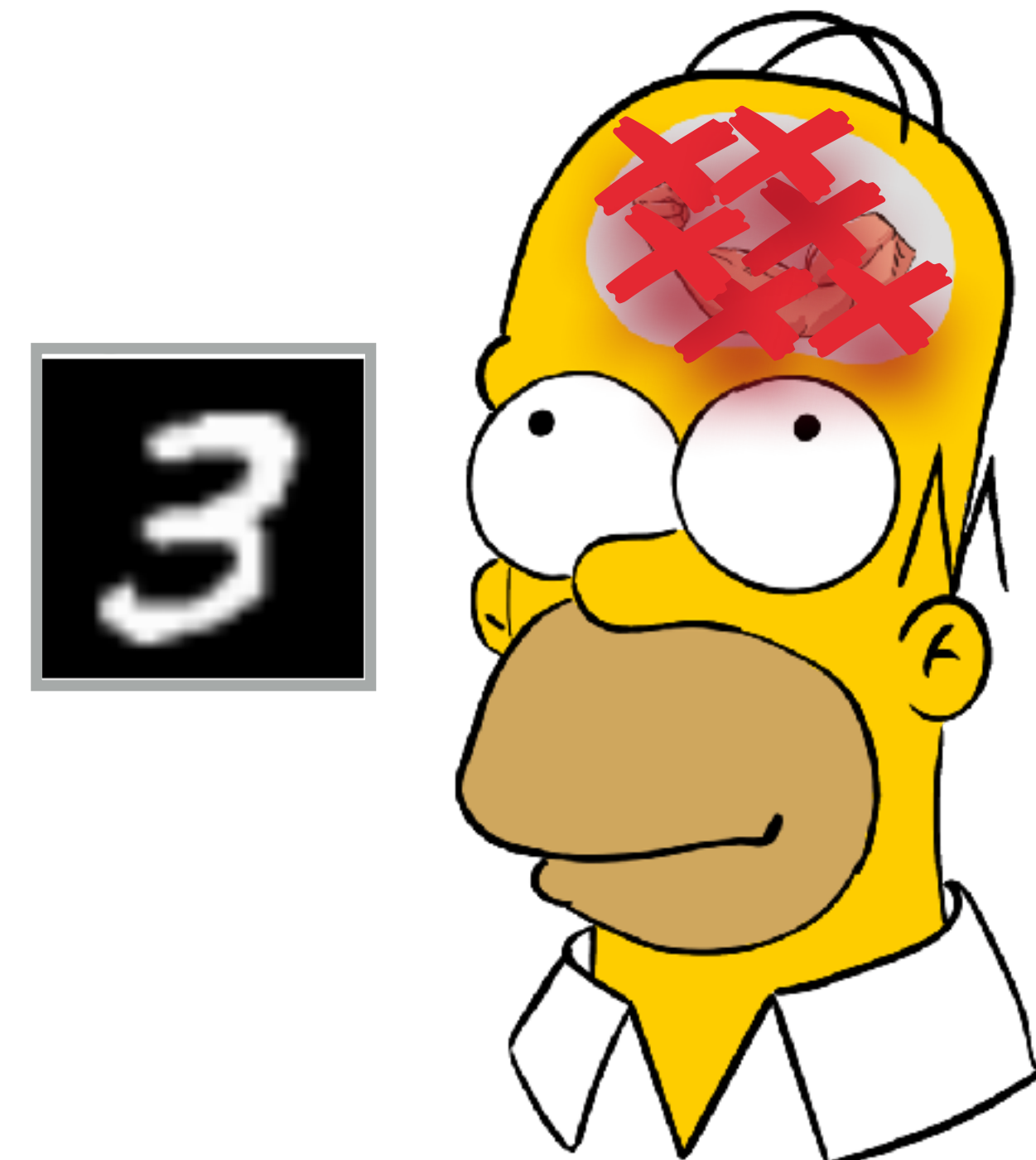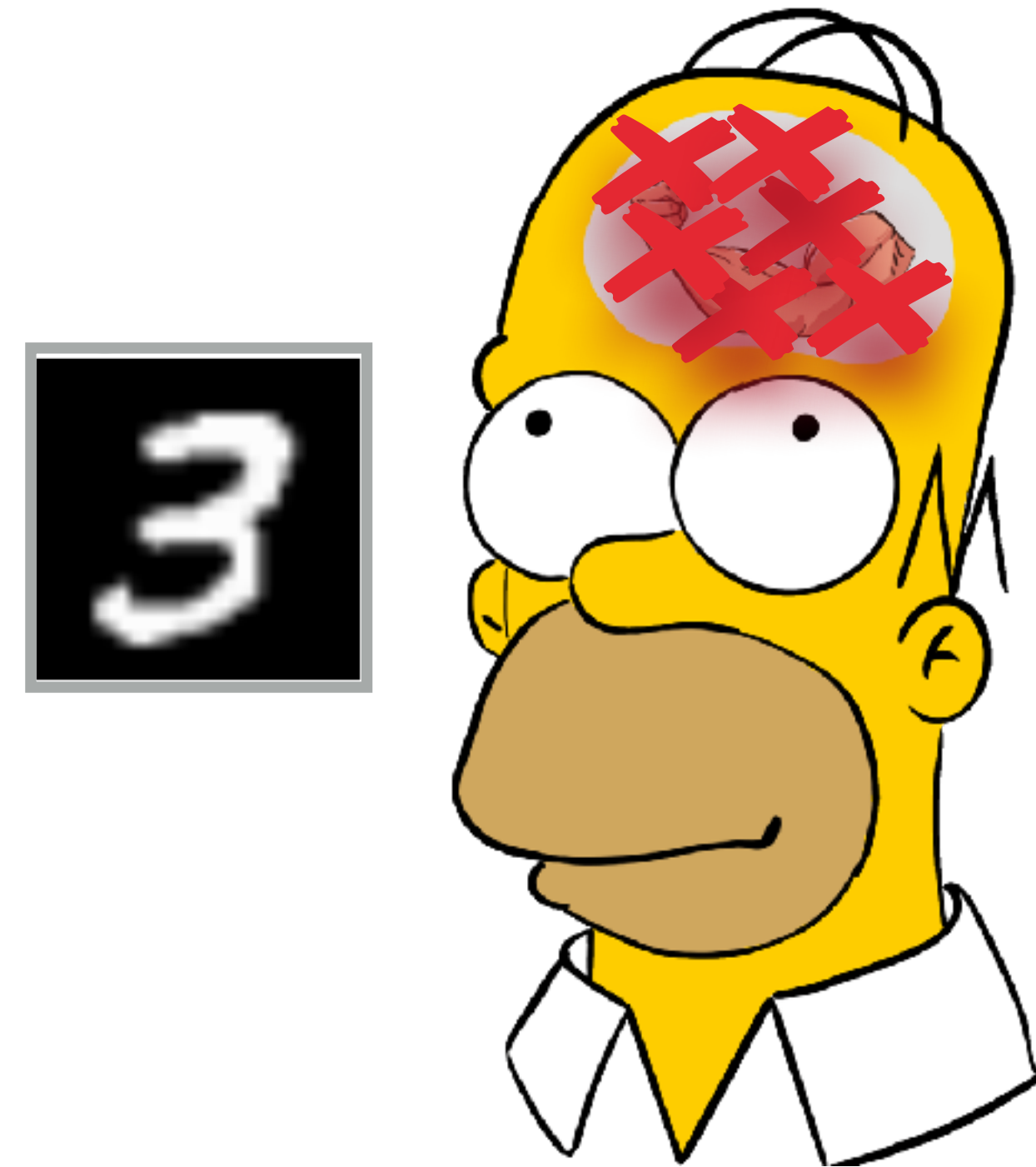
# Multiclass: Problem Setting

- Instances: $\mathbf{x} \in \mathbf{R}^d$

- Labels: $y \in [k] = \{1, 2, \ldots, k\}$

- Multiclass predictor $h : \mathbf{R}^d \to [k]$

- Classification error / mistake: $\mathbf{1}\big[h(\mathbf{x}) \neq y\big]$

- As in binary case: minimizing prediction mistakes is NP-hard

# Prediction

I. Predictor h($\mathbf{x}$) can be a general function

II. Need to express confidence in predicted class

Instead of $\mathbf{h} : \mathbf{R}^d \rightarrow [k]$ use $\mathbf{h} : \mathbf{R}^d \times [k] \rightarrow \mathbf{R}$:

where h($\mathbf{x}$,$c$) confidence that label of $\mathbf{x}$ is $c$

# Winner Takes All

**X** | 0.1 | 7 | -1 | 6.2 |

$h_1$

2.4

**X** | 0.1 | 7 | -1 | 6.2 |

$h_2$

7.1

• • •

**X** | 0.1 | 7 | -1 | 6.2 |

$h_k$

-3

Predicted class:  $\hat{y} = \arg\max_{j} h_j(\mathbf{x})$

# Winner Takes All



Predicted class: $\hat{y} = \arg\max_j h_j(\mathbf{x})$

# Linear Multiclass Predictors

Linear predictor for class $j$ : $h_j(\mathbf{x})$ is $\mathbf{w}_j \cdot \mathbf{x}$

# Linear Multiclass Predictors

Linear predictor for class j : $h_j(\mathbf{x})$ is $\mathbf{w}_j \cdot \mathbf{x}$

Construct matrix W of size k × d whose j'th row is $\mathbf{w}_j$

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ \cdots \\ \cdots \\ \cdots \\ -\mathbf{w}_k- \end{bmatrix}$$

# Linear Multiclass Predictors

Linear predictor for class $j$ : $h_j(\mathbf{x})$ is $\mathbf{w}_j \cdot \mathbf{x}$

Construct matrix $W$ of size $k \times d$ whose $j$'th row is $\mathbf{w}_j$

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ \cdots \\ \cdots \\ -\mathbf{w}_k- \end{bmatrix}$$

Predicted scores: $\mathbf{z} = W\mathbf{x}$

$$\mathbf{z} = \begin{bmatrix} W_{11} & \cdots & W_{1d} \\ \vdots & \cdots & \vdots \\ W_{k1} & \cdots & W_{kd} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

# Linear Multiclass Predictors

Linear predictor for class $j$ : $h_j(\mathbf{x})$ is $\mathbf{w}_j \cdot \mathbf{x}$

Construct matrix $W$ of size $k \times d$ whose $j$'th row is $\mathbf{w}_j$

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ \cdots \\ \cdots \\ \cdots \\ -\mathbf{w}_k- \end{bmatrix}$$

Predicted scores: $\mathbf{z} = W\mathbf{x}$

$$\mathbf{z} = \begin{bmatrix} W_{11} & \cdots & W_{1d} \\ \vdots & \cdots & \vdots \\ W_{k1} & \cdots & W_{kd} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

Predicted label: $\hat{y} = \arg \max_{j=1}^{k} z_j$

# One vs. Rest (One vs. All)

- Learn **k** binary linear predictors

- **j**'th predictor distinguishes **j**'th class from the rest

- Learning scheme:

  I. Transform $S \mapsto S_1, S_2, \ldots, S_k$ where $S_j = \left\{ \left( \mathbf{x}_i, (-1)^{\mathbf{1}[y_i \neq j]} \right) \right\}_{i=1}^{m}$

  II. For j = 1,...,k learn a linear classifier $\mathbf{w}_j$ from $S_j$

- Inference: $\hat{\mathbf{y}} = \arg \max_{j=1}^{k} \mathbf{z}_j = \arg \max_{j=1}^{k} \mathbf{w_j} \cdot \mathbf{x}$

# Example

Original training set:  $S = \{(\mathbf{x}_1, 2), (\mathbf{x}_2, 4), (\mathbf{x}_3, 2), (\mathbf{x}_4, 3), (\mathbf{x}_5, 1)\}$

Results in four binary-labeled datasets:

| $S_1$ | $S_2$ | $S_3$ | $S_4$ |
|---|---|---|---|
| $(\mathbf{x}_1, -)$ | $(\mathbf{x}_1, +)$ | $(\mathbf{x}_1, -)$ | $(\mathbf{x}_1, -)$ |
| $(\mathbf{x}_2, -)$ | $(\mathbf{x}_2, -)$ | $(\mathbf{x}_2, -)$ | $(\mathbf{x}_2, +)$ |
| $(\mathbf{x}_3, -)$ | $(\mathbf{x}_3, +)$ | $(\mathbf{x}_3, -)$ | $(\mathbf{x}_3, -)$ |
| $(\mathbf{x}_4, -)$ | $(\mathbf{x}_4, -)$ | $(\mathbf{x}_4, +)$ | $(\mathbf{x}_4, -)$ |
| $(\mathbf{x}_5, +)$ | $(\mathbf{x}_5, -)$ | $(\mathbf{x}_5, -)$ | $(\mathbf{x}_5, -)$ |

# Deficiencies of OvA

# Deficiencies of OvA

Predictors are trained independently and then used jointly

# Deficiencies of OvA

Predictors are trained independently and then used jointly

Albeit trained independently "compete" during inference

# Deficiencies of OvA

Predictors are trained independently and then used jointly

Albeit trained independently "compete" during inference

Resulting binary problems might be overly difficult
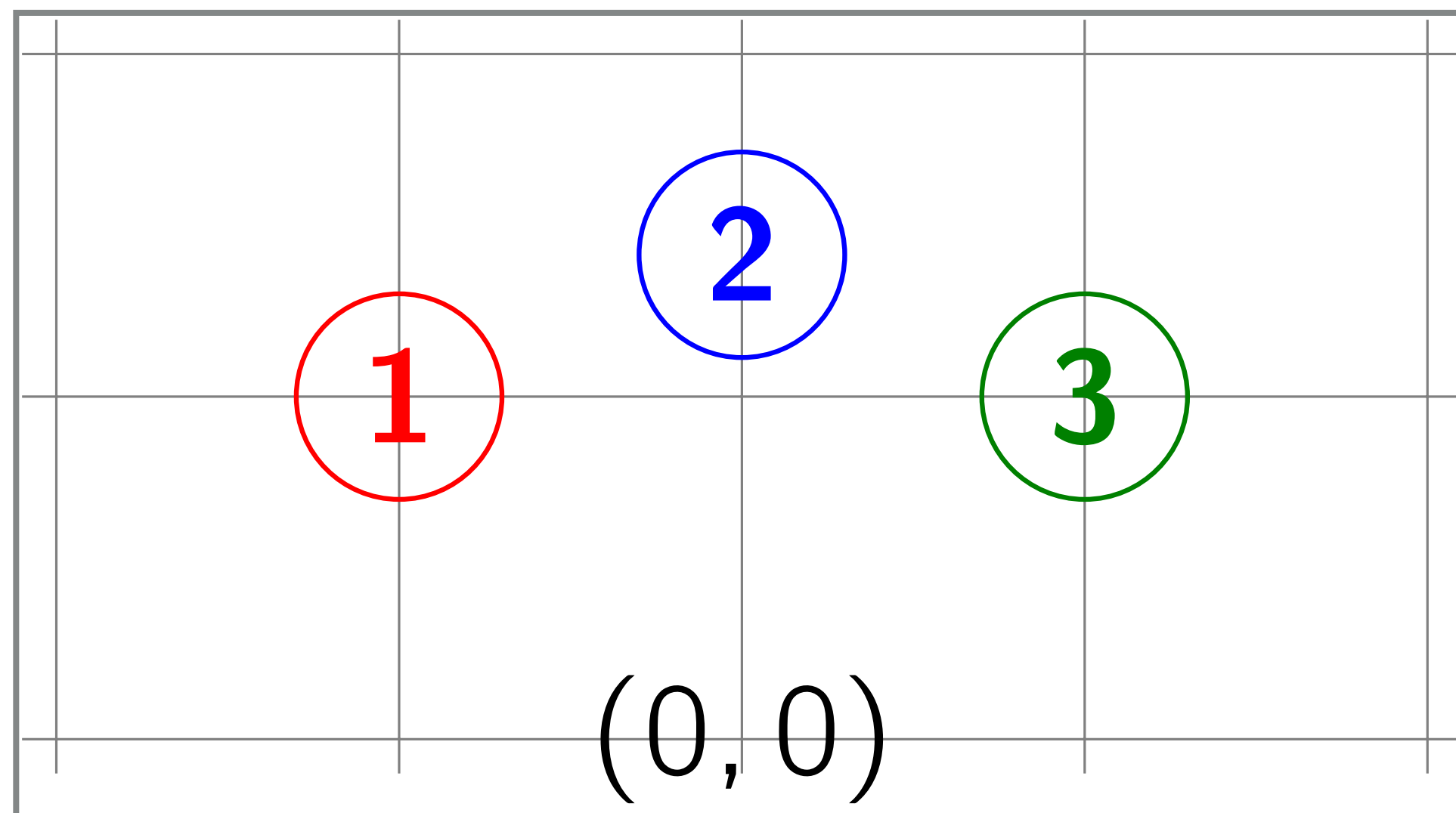
# Deficiencies of OvA

Predictors are trained independently and then used jointly

Albeit trained independently "compete" during inference

Resulting binary problems might be overly difficult
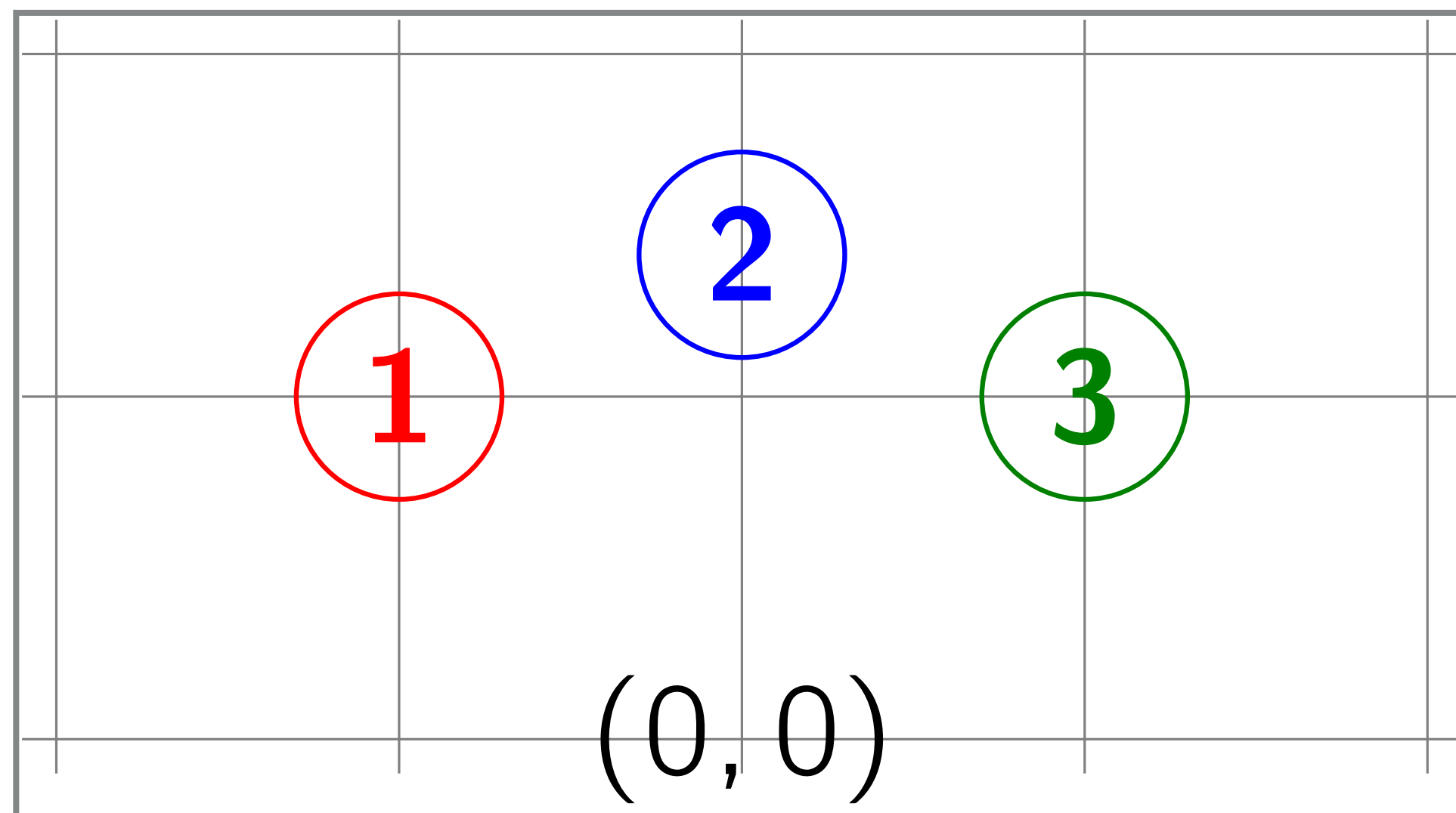
OvA would fail for setting:

# Deficiencies of OvA

Predictors are trained independently and then used jointly

Albeit trained independently "compete" during inference

Resulting binary problems might be overly difficult

OvA would fail for setting:



(0, 0)

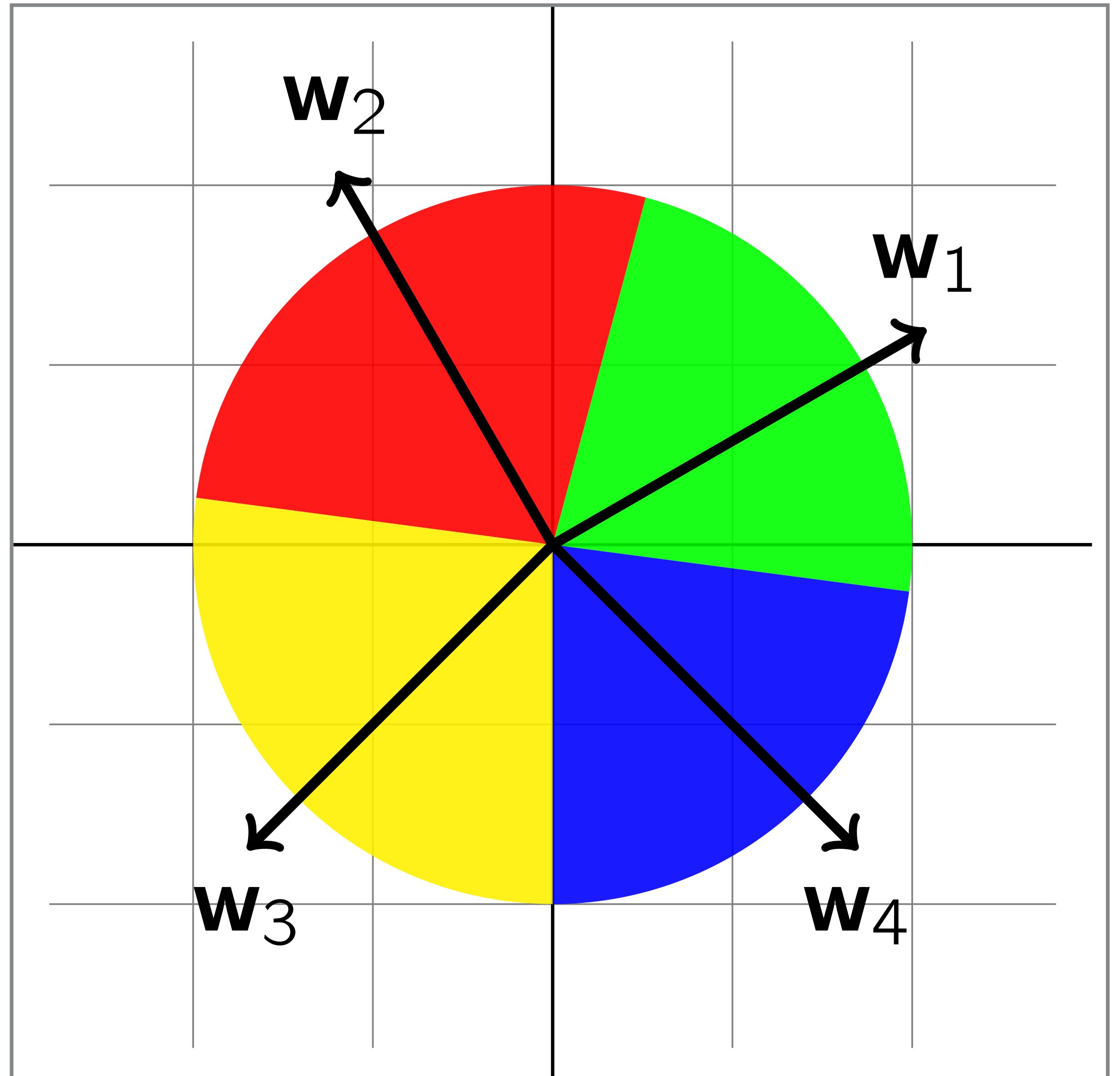While data is linearly separable using:

$$W = \begin{bmatrix} -1 & 1 \\ 0 & \sqrt{2} \\ 1 & 1 \end{bmatrix}$$

# Multiclass Margin

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ -\mathbf{w}_3- \\ -\mathbf{w}_4- \end{bmatrix} \in \mathbf{R}^{4x2}$$

Assume : $\|\mathbf{w}_j\| = 1$   $\|\mathbf{x}\| = 1$

$$\sphericalangle(\mathbf{w}, \mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

# Multiclass Margin

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ -\mathbf{w}_3- \\ -\mathbf{w}_4- \end{bmatrix} \in \mathbf{R}^{4x2}$$

Assume : $\|\mathbf{w}_j\| = 1$   $\|\mathbf{x}\| = 1$

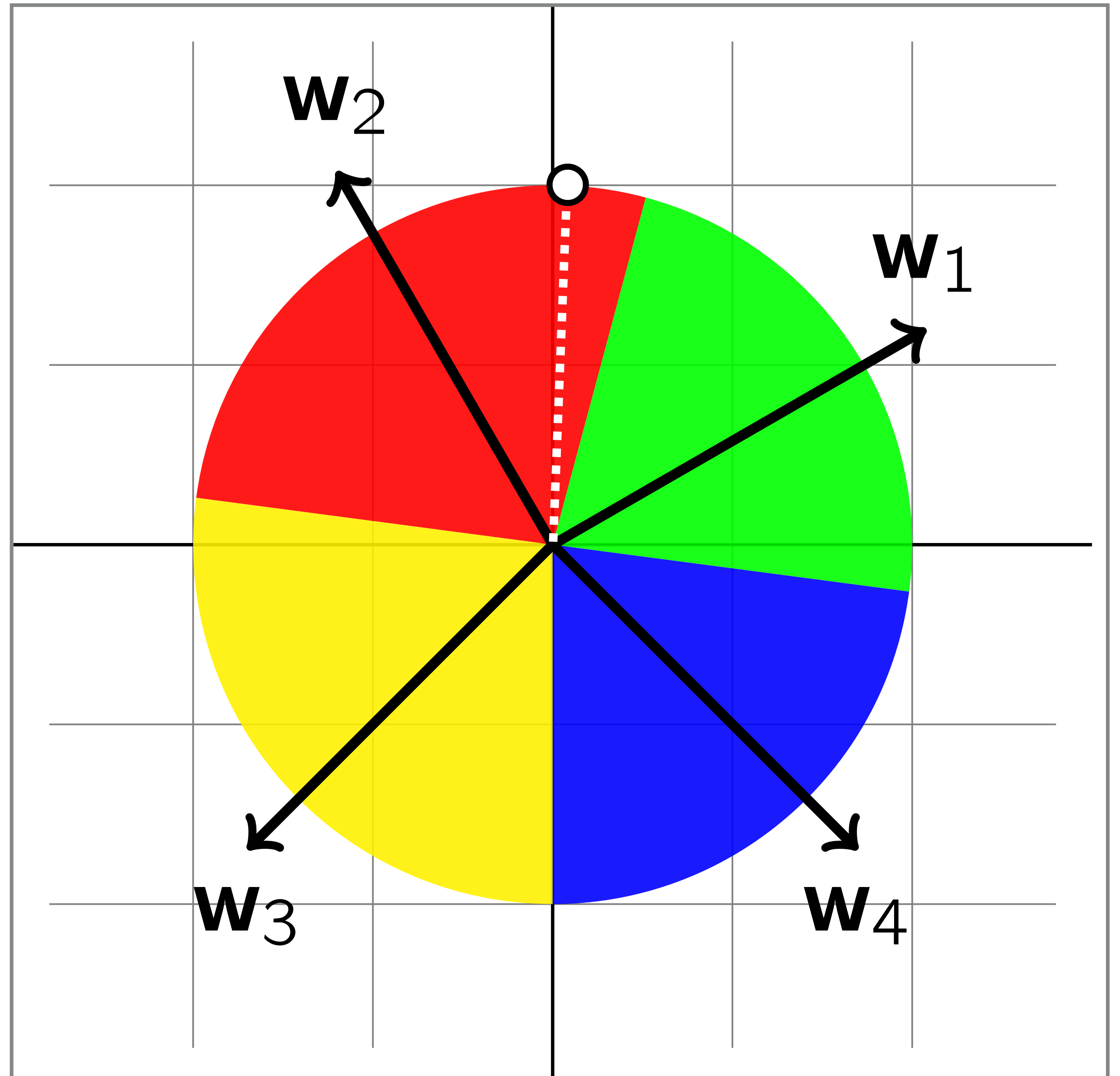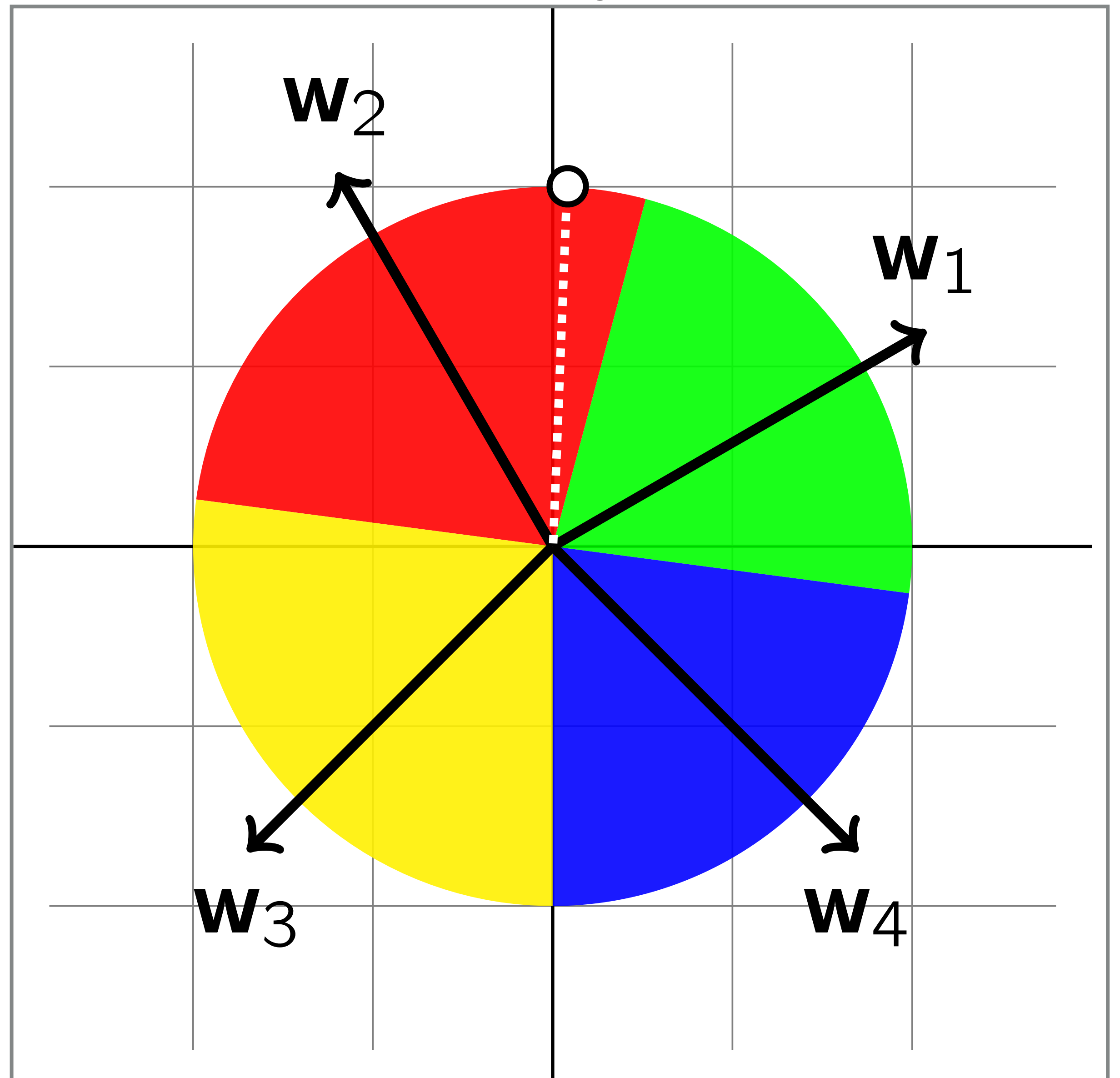$$\sphericalangle(\mathbf{w}, \mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

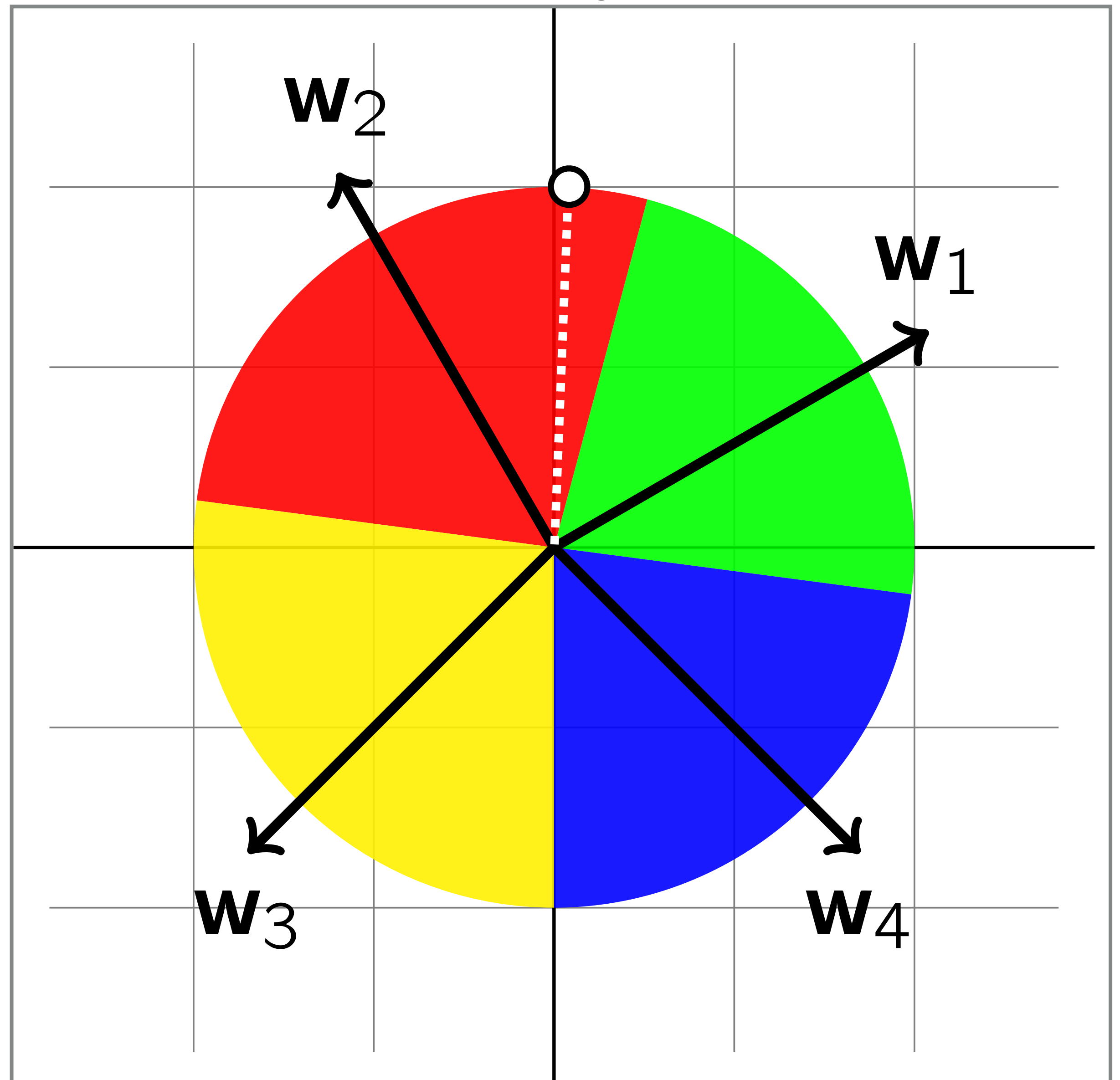# Multiclass Margin

$$W = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ -\mathbf{w}_3- \\ -\mathbf{w}_4- \end{bmatrix} \in \mathbf{R}^{4x2}$$

Assume : $\|\mathbf{w}_j\| = 1$  $\|x\| = 1$

$$\measuredangle(\mathbf{w}, \mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

$$\measuredangle(\mathbf{w_2}, \mathbf{x}) < \measuredangle(\mathbf{w_j}, \mathbf{x})$$
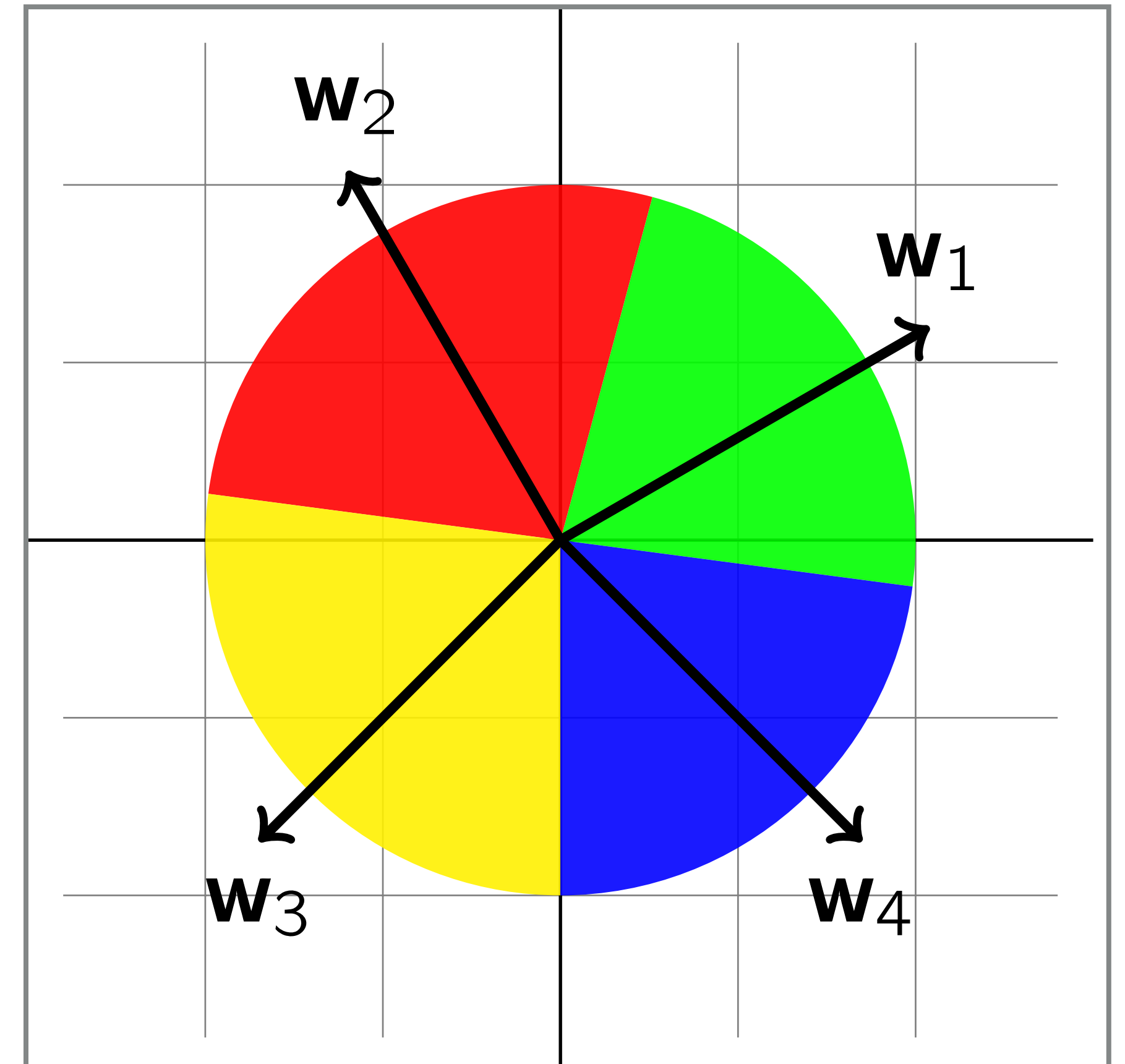
© 2020 Yoram Singer

11

# Multiclass Margin

For general vectors impose:

$$(\mathbf{x}, y) \quad \Rightarrow \quad \forall j \neq y : \quad \mathbf{w}_y \cdot \mathbf{x} > \mathbf{w}_j \cdot \mathbf{x}$$

In matrix-vector format:

$$(\mathbf{x}, y) \quad \Rightarrow \quad \forall j \neq y : \quad [W\mathbf{x}]_y > [W\mathbf{x}]_j$$

# Margin Loss

Predicted class: $\hat{y}(\mathbf{z}) = \arg\max\limits_{j=1}^{k} z_j$

Classification error:

$$\ell^{\text{MC}}(\mathbf{z}) = \mathbf{1}\big[\hat{y}(\mathbf{z}) \neq y\big]$$

Max-Margin Loss is difference in scores + penalty $\gamma$:

$$\ell^{\text{MM}}(\mathbf{z}) = \left[\gamma + \max\limits_{j \neq y} z_j - z_y\right]_+ \quad \text{where} \quad [z]_+ = \max\{0, z\}$$

Margin great than $\gamma \Rightarrow \ell^{\text{MC}} = \ell^{\text{MM}} = 0$



Margin $\in (0, \gamma) \Rightarrow \ell^{\text{MC}} = 0$ but $\ell^{\text{MM}} \geq 0$



Margin $< 0 \Rightarrow \ell^{\text{MC}} = 1$ and $\ell^{\text{MM}} \geq \gamma$

# Convexity of Max-Margin Loss*

Inner product $\mathbf{w}_j \cdot \mathbf{x}$ is linear in $\mathbf{w}_j$ $\Rightarrow$ $\mathbf{w}_j \cdot \mathbf{x}$ convex in $\mathbf{w}_j$ (and concave)

# Convexity of Max-Margin Loss*

Inner product $\mathbf{w_j} \cdot \mathbf{x}$ is linear in $\mathbf{w_j} \Rightarrow \mathbf{w_j} \cdot \mathbf{x}$ convex in $\mathbf{w_j}$ (and concave)

Maximum of convex function is convex

$$\max\{\mathbf{x}^2, |\mathbf{x} - a|\}$$

$$|\mathbf{x} - a|$$

$$\mathbf{x}^2$$

# Convexity of Max-Margin Loss*

Inner product $\mathbf{w_j} \cdot \mathbf{x}$ is linear in $\mathbf{w_j} \Rightarrow \mathbf{w_j} \cdot \mathbf{x}$ convex in $\mathbf{w_j}$ (and concave)

Maximum of convex function is convex

Therefore $\max\limits_{j \neq y} \mathbf{w_j} \cdot \mathbf{x}$ is convex in W

$$\max\{\mathbf{x}^2, |\mathbf{x} - a|\}$$
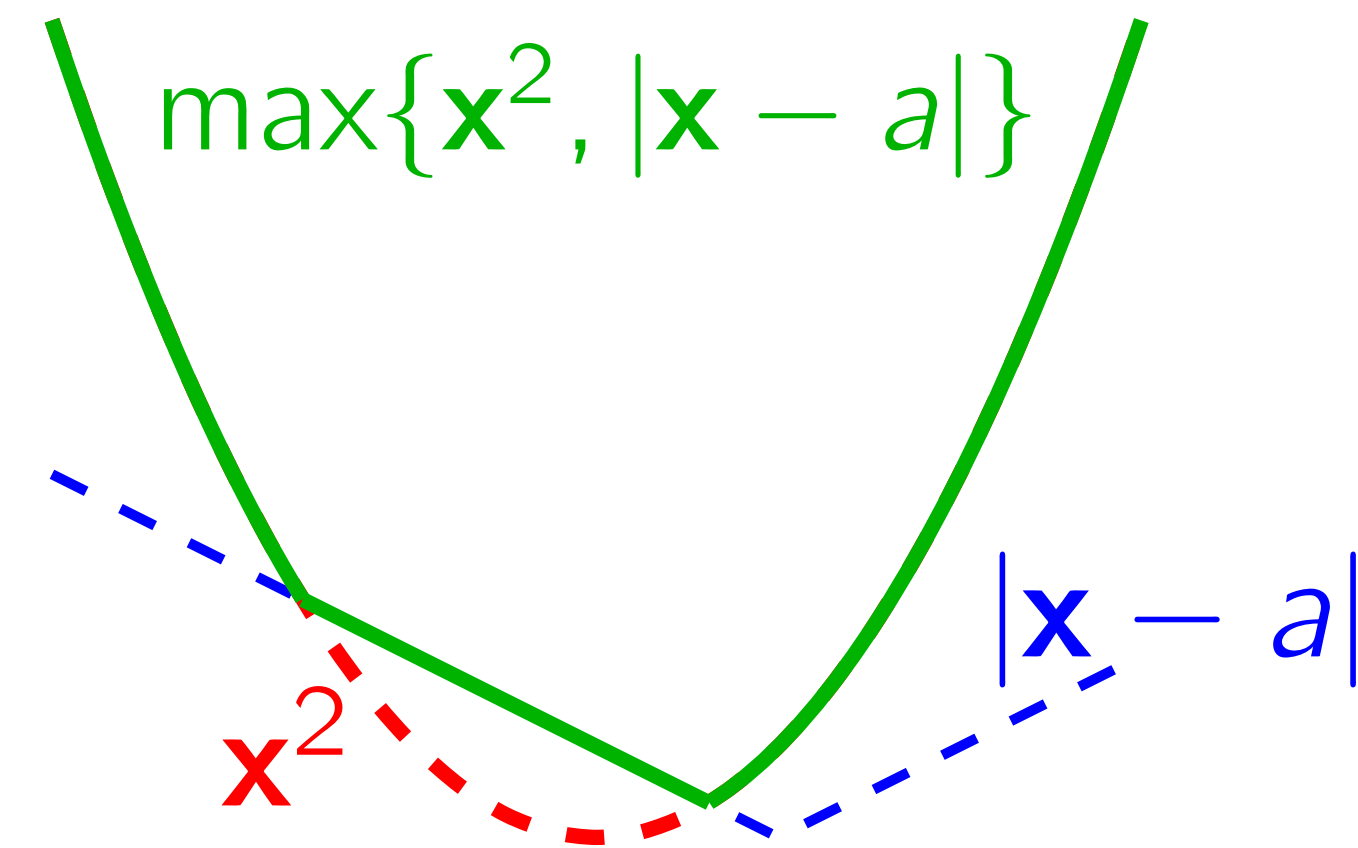
$$|\mathbf{x} - a|$$

$$\mathbf{x}^2$$

# Convexity of Max-Margin Loss*

Inner product $\mathbf{w}_j \cdot \mathbf{x}$ is linear in $\mathbf{w}_j \Rightarrow \mathbf{w}_j \cdot \mathbf{x}$ convex in $\mathbf{w}_j$ (and concave)

Maximum of convex function is convex

$\max\{\mathbf{x}^2, |\mathbf{x} - a|\}$

Therefore $\max\limits_{j \neq y} \mathbf{w}_j \cdot \mathbf{x}$ is convex in W

$|\mathbf{x} - a|$

$\mathbf{x}^2$

Sum of convex functions is convex $\Rightarrow \gamma + \max\limits_{j} \mathbf{w}_j \cdot \mathbf{x} - \mathbf{w}_y \cdot \mathbf{x}$ is convex in W
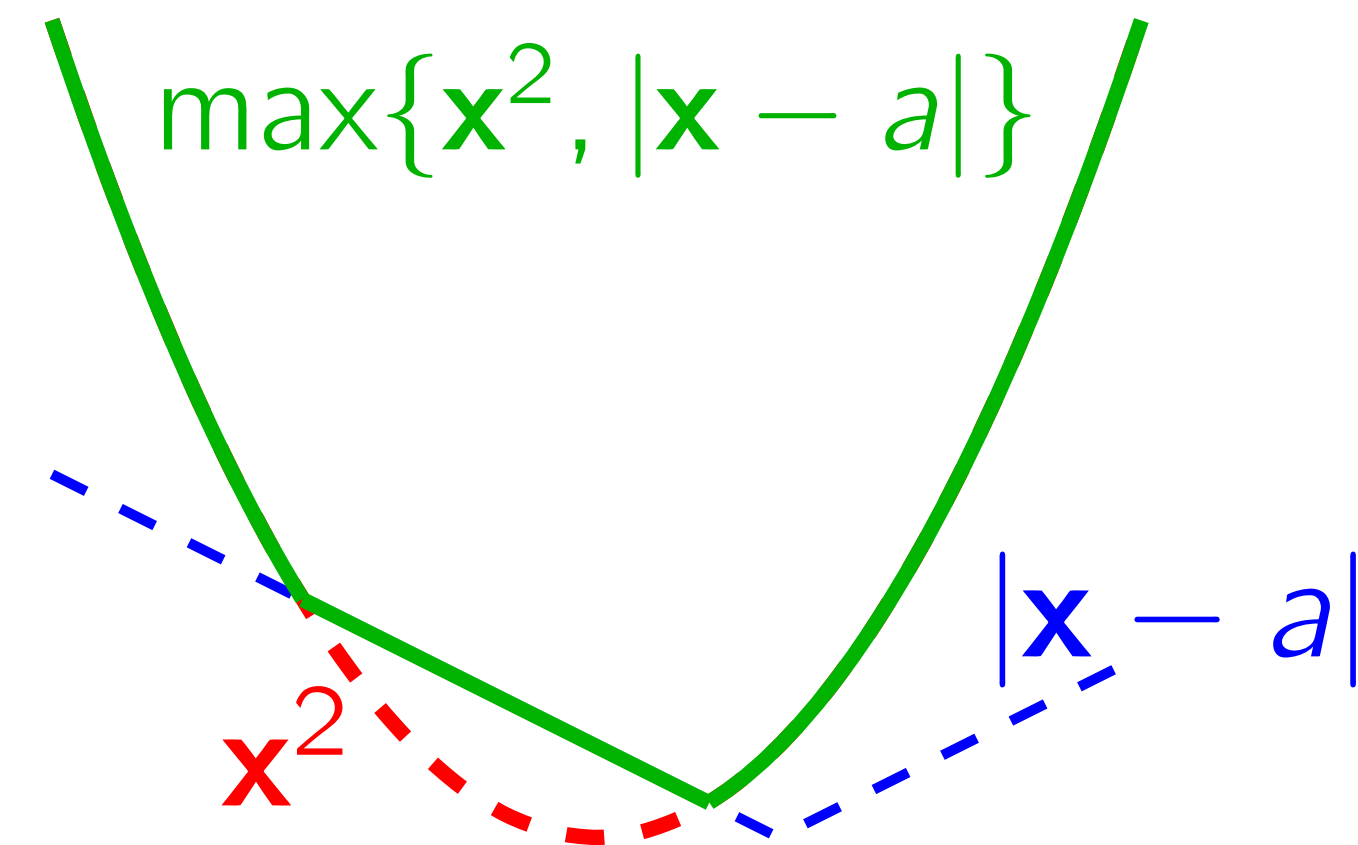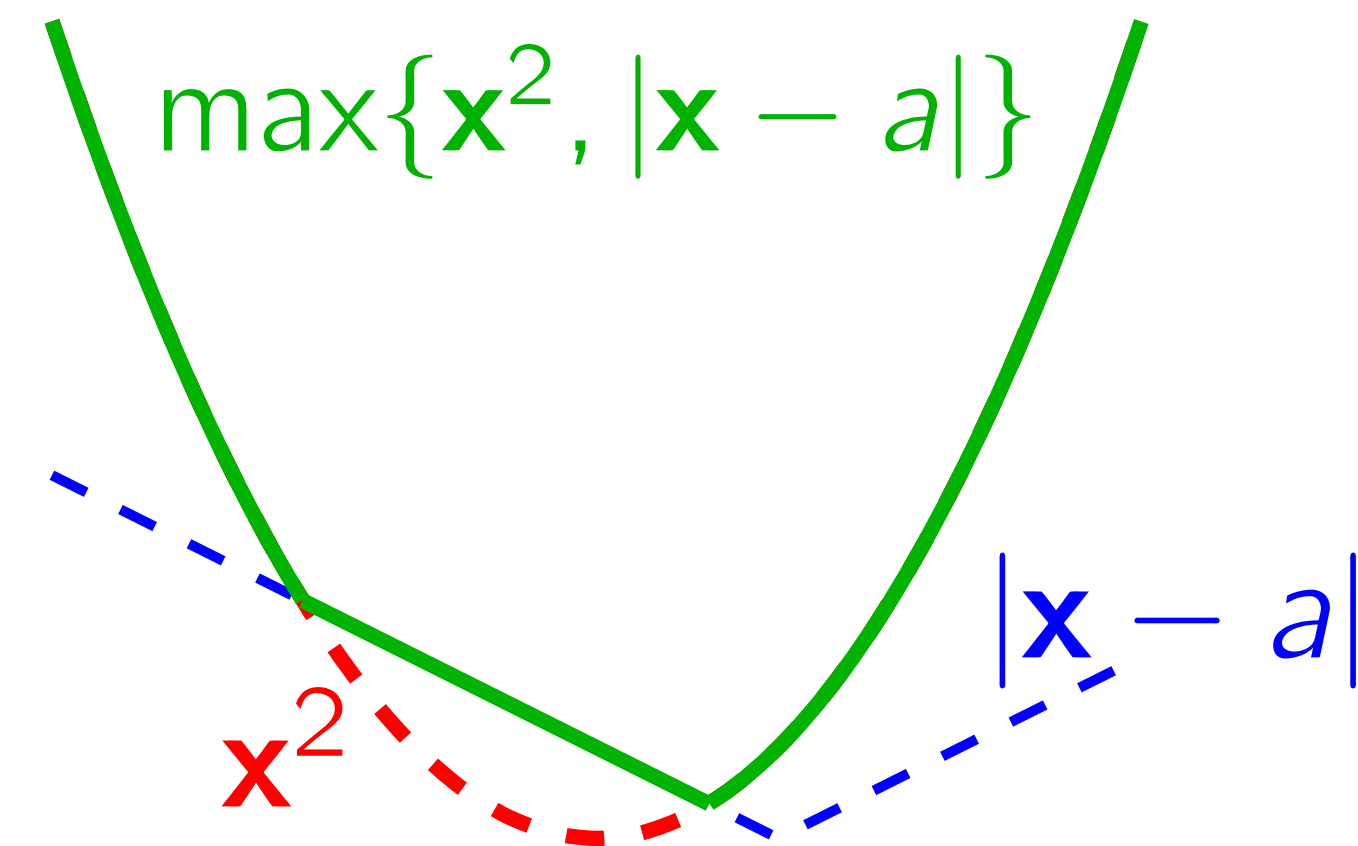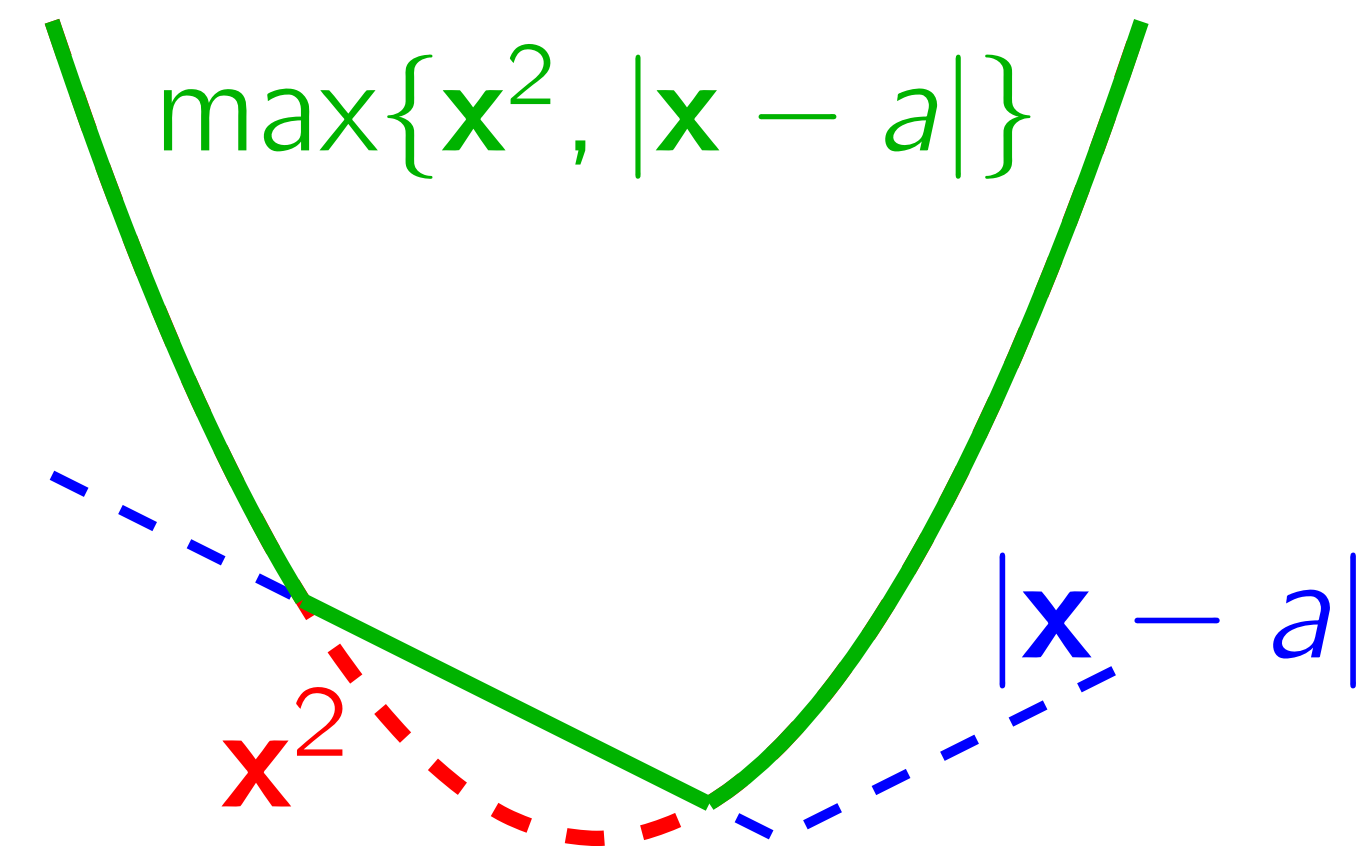
# Convexity of Max-Margin Loss*

Inner product $\mathbf{w_j} \cdot \mathbf{x}$ is linear in $\mathbf{w_j} \Rightarrow \mathbf{w_j} \cdot \mathbf{x}$ convex in $\mathbf{w_j}$ (and concave)

Maximum of convex function is convex

Therefore $\displaystyle\max_{j \neq y} \mathbf{w_j} \cdot \mathbf{x}$ is convex in W

Sum of convex functions is convex $\Rightarrow \gamma + \displaystyle\max_{j} \mathbf{w_j} \cdot \mathbf{x} - \mathbf{w_y} \cdot \mathbf{x}$ is convex in W

Using convexity of maximum again: $\ell^{\mathsf{MM}}(\mathbf{z}) = \max\left\{0, \gamma + \displaystyle\max_{j \neq y} z_j - z_y\right\}$

Implies that $\ell^{\mathsf{MM}}(\mathbf{z})$ is convex in W

$\max\{\mathbf{x}^2, |\mathbf{x} - a|\}$

$|\mathbf{x} - a|$

$\mathbf{x}^2$

# Multivariate Logistic Regression

As before $\mathbf{z} = W\mathbf{x}$

Define probability of class c to be $\mathbf{P}\left[c \mid \mathbf{z}\right] = \dfrac{e^{z_c}}{Z}$  where $Z = \displaystyle\sum_{j=1}^{k} e^{z_j}$

Loss: -log-probability of correct class  $\ell^{\mathsf{LR}}(\mathbf{z}) = -\log\Big(P\left[y \mid \mathbf{x}\right]\Big)$

$\hat{y} = \arg\displaystyle\max_{j=1}^{k} z_j \neq y$

$= \underbrace{\log\Big(\textstyle\sum_j \exp(z_j)\Big) - z_y}_{\text{SoftMax}}$

$\mathsf{SoftMax}(\mathbf{z}) \equiv \log\Big(\textstyle\sum_i e^{z_i}\Big) \geq \displaystyle\max_i z_i$

# MC Logistic Regression & Error

Multiclass prediction error:

$$\ell^{\mathsf{MC}}(\mathbf{y}, \mathbf{z}) = 1 \quad \Leftrightarrow \quad \hat{\mathbf{y}} = \arg \max_{j=1}^{k} z_j \neq \mathbf{y}$$

If $\hat{\mathbf{y}} \neq \mathbf{y}$ then $\ell^{\mathsf{LR}}(\mathbf{z}) = \log\left(\sum_j \exp(z_j)\right) - z_y$

$$\geq \log\left(\underbrace{\exp(z_y) + \exp(z_{\hat{y}})}_{\geq\, 2\exp(z_y)}\right) - z_y \geq \log(2)$$

Therefor $\ell^{\mathsf{MC}}(\mathbf{y}, \mathbf{z}) = 1 \quad \Rightarrow \quad \ell^{\mathsf{LR}}(\mathbf{y}, \mathbf{z}) \geq 2$

# Multiclass LR & Max-Margin

If $\ell^{\mathsf{MM}}(z) \geq \beta > 0 \quad \Rightarrow \quad \exists j : \gamma + z_j - z_y \geq \beta$

Then $\ell^{\mathsf{LR}}(z) = \log\left(\sum_j \exp(z_j)\right) - z_y \geq \log\left(\exp(z_y) + \exp(z_j)\right) - z_y$

Hence $\ell^{\mathsf{LR}}(z) \geq \log\left((\exp(z_y) + e^{\beta-\gamma}\exp(z_y)\right) - z_y$

$$\geq \log(1 + e^{\beta-\gamma}) \geq \beta - \gamma$$

In summary: $\ell^{\mathsf{MM}}(\mathbf{y}, \mathbf{z}) = \beta \quad \Rightarrow \quad \ell^{\mathsf{LR}}(\mathbf{y}, \mathbf{z}) \geq \beta - \gamma$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $W^t$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $W^t$

Loss $\mathscr{L}(W)$ is a function of $W \Rightarrow$ Gradient is a matrix $G^t = \nabla_W \mathscr{L}(W)$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $W^t$

Loss $\mathscr{L}(W)$ is a function of $W \Rightarrow$ Gradient is a matrix $G^t = \nabla_W \mathscr{L}(W)$

Gradient step without (or before) projection $W^{t+1} \leftarrow W^t - \eta_t G^t$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $\mathsf{W}^t$

Loss $\mathscr{L}(\mathsf{W})$ is a function of $\mathsf{W} \Rightarrow$ Gradient is a matrix $\mathsf{G}^t = \nabla_{\mathsf{W}} \mathscr{L}(\mathsf{W})$

Gradient step without (or before) projection $\mathsf{W}^{t+1} \leftarrow \mathsf{W}^t - \eta_t \mathsf{G}^t$

Projection adheres with matrix form:

$$
\begin{bmatrix}
\| - \mathbf{w}_1 - \| \leq r \\
\| - \mathbf{w}_2 - \| \leq r \\
\vdots \\
\| - \mathbf{w}_k - \| \leq r
\end{bmatrix}
$$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $W^t$

Loss $\mathscr{L}(W)$ is a function of $W \Rightarrow$ Gradient is a matrix $G^t = \nabla_W \mathscr{L}(W)$

Gradient step without (or before) projection $W^{t+1} \leftarrow W^t - \eta_t G^t$

Projection adheres with matrix form:

$$\begin{bmatrix} \| - \mathbf{w}_1 - \| \leq r \\ \| - \mathbf{w}_2 - \| \leq r \\ \vdots \\ \| - \mathbf{w}_k - \| \leq r \end{bmatrix} \qquad \mathbf{w}_j^{t+1/2} \leftarrow \mathbf{w}_j^t - \eta_t \mathbf{g}_j^t$$

# SGD for Multiclass

Instead of $\mathbf{w}_t$ maintain a matrix $W^t$

Loss $\mathscr{L}(W)$ is a function of W $\Rightarrow$ Gradient is a matrix $G^t = \nabla_W \mathscr{L}(W)$

Gradient step without (or before) projection $W^{t+1} \leftarrow W^t - \eta_t G^t$

Projection adheres with matrix form:

$$
\begin{bmatrix}
\| - \mathbf{w}_1 - \| \leq r \\
\| - \mathbf{w}_2 - \| \leq r \\
\vdots \\
\| - \mathbf{w}_k - \| \leq r
\end{bmatrix}
\qquad
\begin{aligned}
\mathbf{w}_j^{t+1/2} &\leftarrow \mathbf{w}_j^t - \eta_t \mathbf{g}_j^t \\
\mathbf{w}_j^{t+1} &\leftarrow \min\left\{1, r/\|\mathbf{w}_j^{t+1/2}\|\right\} \mathbf{w}_j^{t+1/2}
\end{aligned}
$$

# Multiclass Logistic Regression

For each example ($\mathbf{x}$,y) in mini-batch calculate $\mathbf{z} = \mathsf{W}\mathbf{x}$

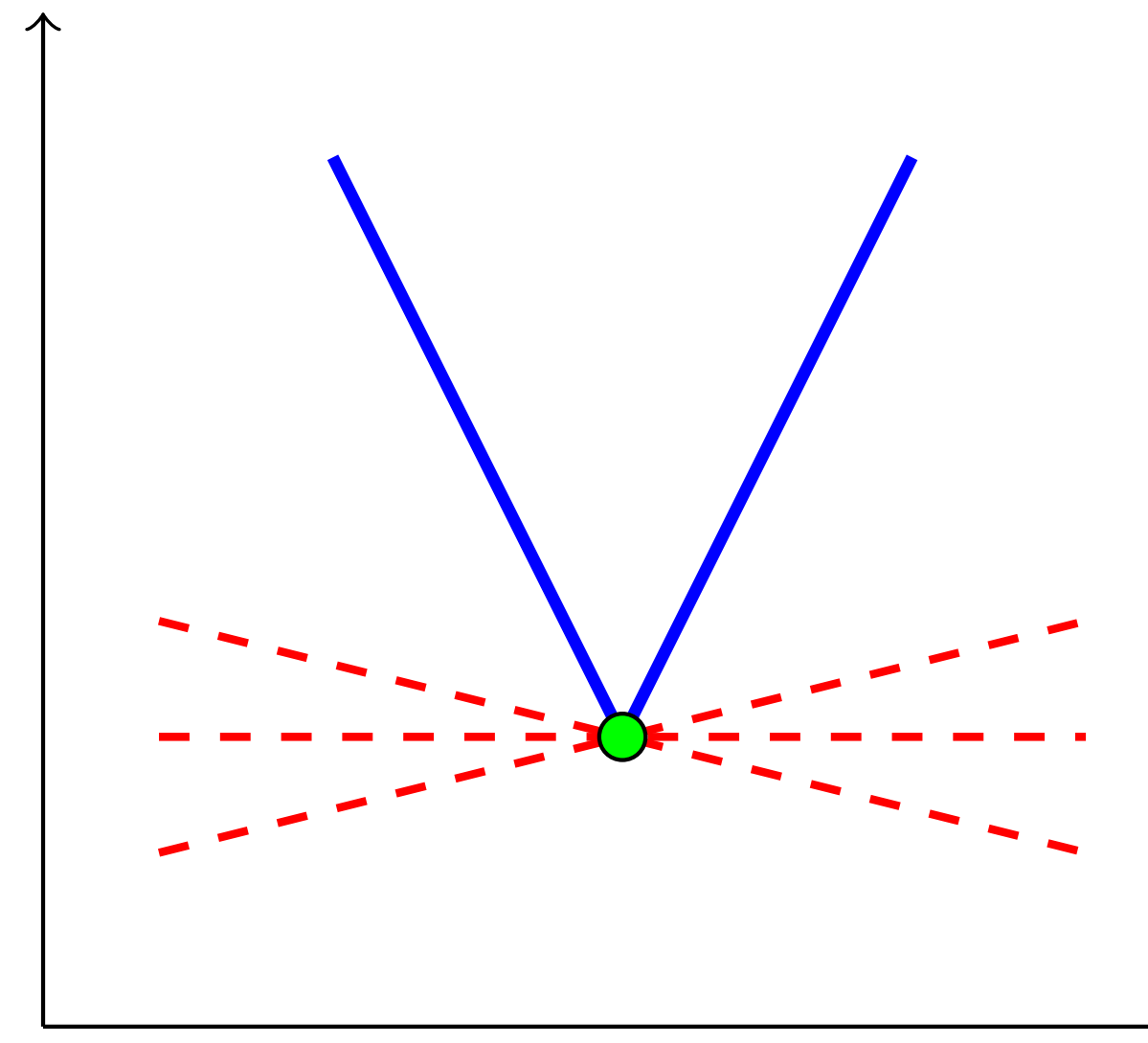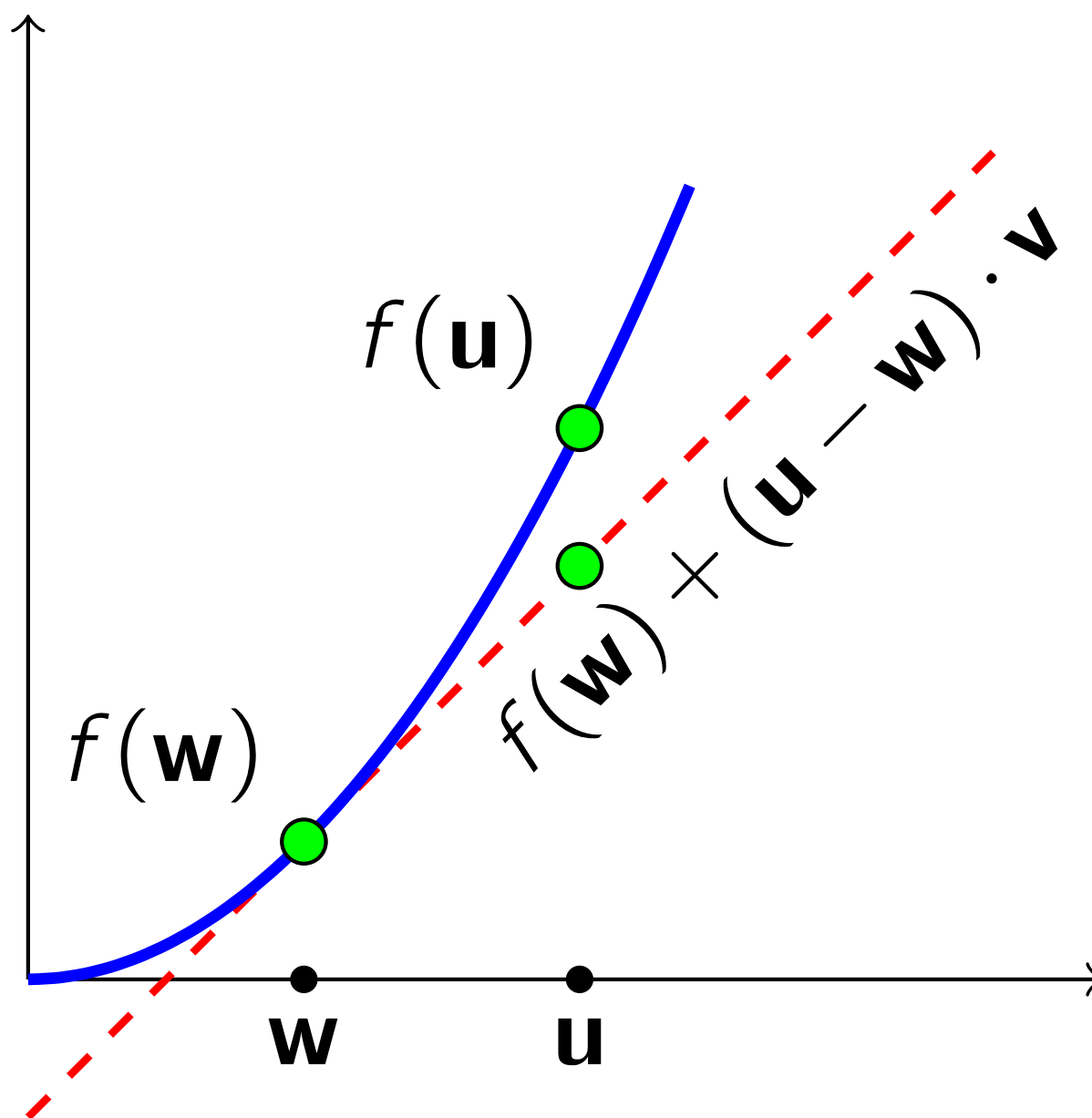Define: $v[j] = \dfrac{\exp(z_j)}{\sum_{i=1}^{k} \exp(z_i)} - \mathbf{1}[j=y]$

Gradient: $\mathbf{v}\,\mathbf{x}^\top = \begin{bmatrix} v[1]\,\mathbf{x} \\ v[2]\,\mathbf{x} \\ \vdots \\ v[k]\,\mathbf{x} \end{bmatrix}$ and for mini-batch: $\mathsf{G} = \dfrac{1}{|S|}\sum_{i \in S} \mathbf{v_i}\mathbf{x_i}^\top$

# Sub-gradients*

**w** is a sub-gradient of f at **w** if $\forall \mathbf{u}, \quad \mathbf{f(u)} \geq \mathbf{f(w)} + \mathbf{v} \cdot \mathbf{(u - w)}$

Differential set $\partial$f(**w**) is the set of sub-gradients of f at **w**

# Sub-gradient for Max Margin*

Set of labels with margin error $\Gamma = \{j \neq y \mid \gamma + z_j - z_y \geq 0\}$

Sub-gradients for MM loss are vectors **p** of the form:

$$p[y] = -1 \; ; \text{ for } \; j \notin \Gamma : p[j] = 0 \; ; \; \sum_{j \in \Gamma} p[j] = 1 \; (p[j] \geq 0)$$

Example: y = 2    **z** = [-2  3  2.5  1  7  4  1.9]    $\gamma = 1$

$\Gamma = \{3, 5, 6\}$    **p** = [0 -1 1 0 0 0 0] or **p** = [0 -1 0.1 0 0.4 0.5 0] or ...

# Families of Updates

$$y = 3 \quad \gamma = 2$$

z | -1 | 2.5 | 3 | 4 | 1 | 2.5

For all forms of updates:  $p[y] = 1$

Max only:  $p[\hat{y}] = -1$

p | 0 | 0 | -1 | 1 | 0 | 0

Uniform:  $\forall j \in \Gamma : p[j] = \dfrac{1}{|\Gamma|}$

p | 0 | 1/3 | -1 | 1/3 | 0 | 1/3

Margin-based:  $\forall j \in \Gamma : p[j] = \dfrac{z_j - z_y}{Z}$  where $Z = \sum\limits_{j \in \Gamma} z_j - z_y$

p | 0 | 1.5 | -1 | 3 | 0 | 1.5   /Z=6

p | 0 | 1/4 | -1 | 1/2 | 0 | 1/4

# Mini-Batch Max-Margin Subgradient*

For each $i \in S$:

1. Calculate predicted values: $\mathbf{z}_i = W \mathbf{x}_i$

2. Calculate margin-error sets: $\Gamma = \{j \neq y_i \,|\, \gamma + z_i[j] - z_i[y] \geq 0\}$

3. Form update vectors: $\mathbf{p}_i$

4. Gradient: $G = \dfrac{1}{|S|} \sum_{i \in S} \mathbf{p_i} \mathbf{x_i}^{\top}$

# Max-Margin vs. Soft-Max*

Both updates of the form: $W^{t+1} \leftarrow W^t - \eta_t \, \mathbf{p} \, \mathbf{x}^\top$

Both satisfy $\sum_j \mathsf{p}[j] = 0 \;\; ; \;\; \sum_{j \neq y} \mathsf{p}[j] \leq 1$

If $\Gamma \neq \varnothing$ then for MM $\mathsf{p}[y] = -1$ and for LR $\mathsf{p}[y] > -1$

LR is a dense update $|\{j : \mathsf{p}[j] > 0\}| = k - 1$

MM is a sparse update $|\{j : \mathsf{p}[j] > 0\}| \leq |\Gamma|$

# Cost-Sensitive Multiclass*

Classes often have semantic meaning and similarities

Image classification: Ape ≈ Baboon but Ape ≉ Subaru



Cost of confusing class y with class y': C(y,y')>0   [and C(y,y)=0]

Replace a fixed margin of $\gamma$ with label-dependent margin C(y,y')

# Cost Sensitive Multiclass*

Proxy for bounding $C(y, \hat{y})$

$$C(y, \hat{y}) \leq C(y, \hat{y}) + z_{\hat{y}} - z_y$$

$$\leq \underbrace{\max_r C(y, r) + z_r - z_y}_{\equiv \ell(y, \mathbf{z})}$$

# Usage: Hierarchical Classification*
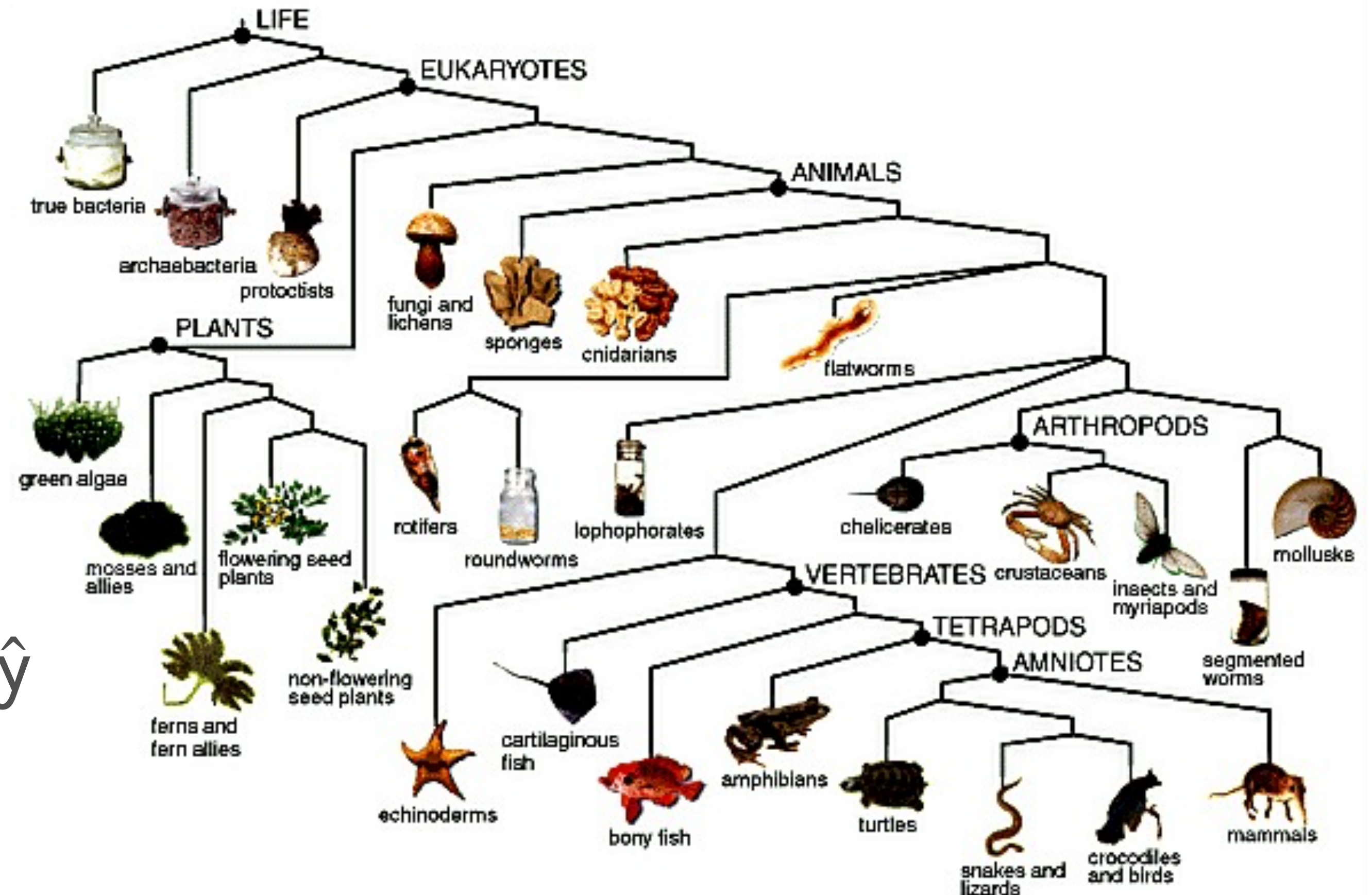

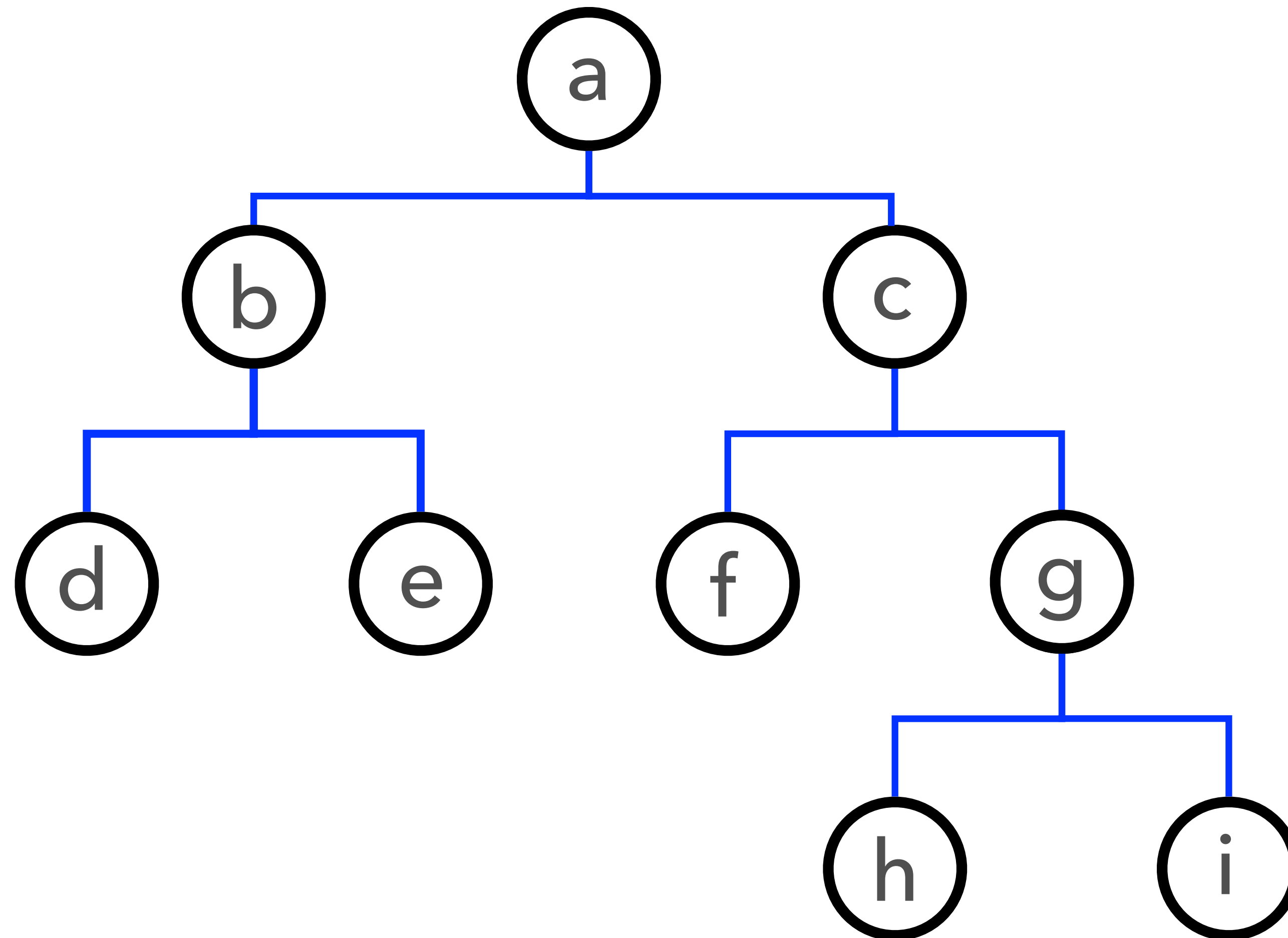
Classes organized in a hierarchy

Cost of $C(y, \hat{y})$:

Length of (unique) path from y to $\hat{y}$

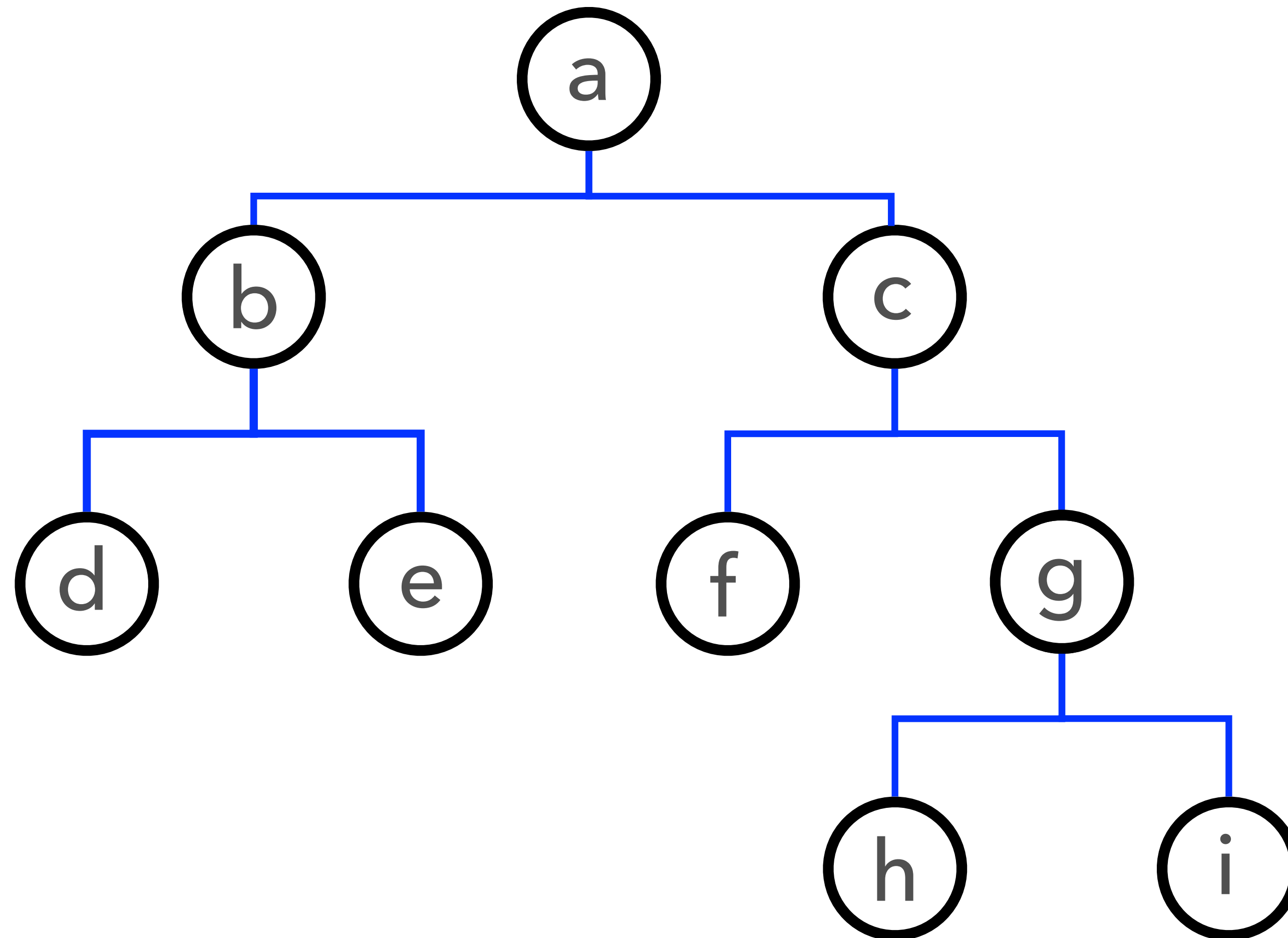C(turtles, snakes) = 1
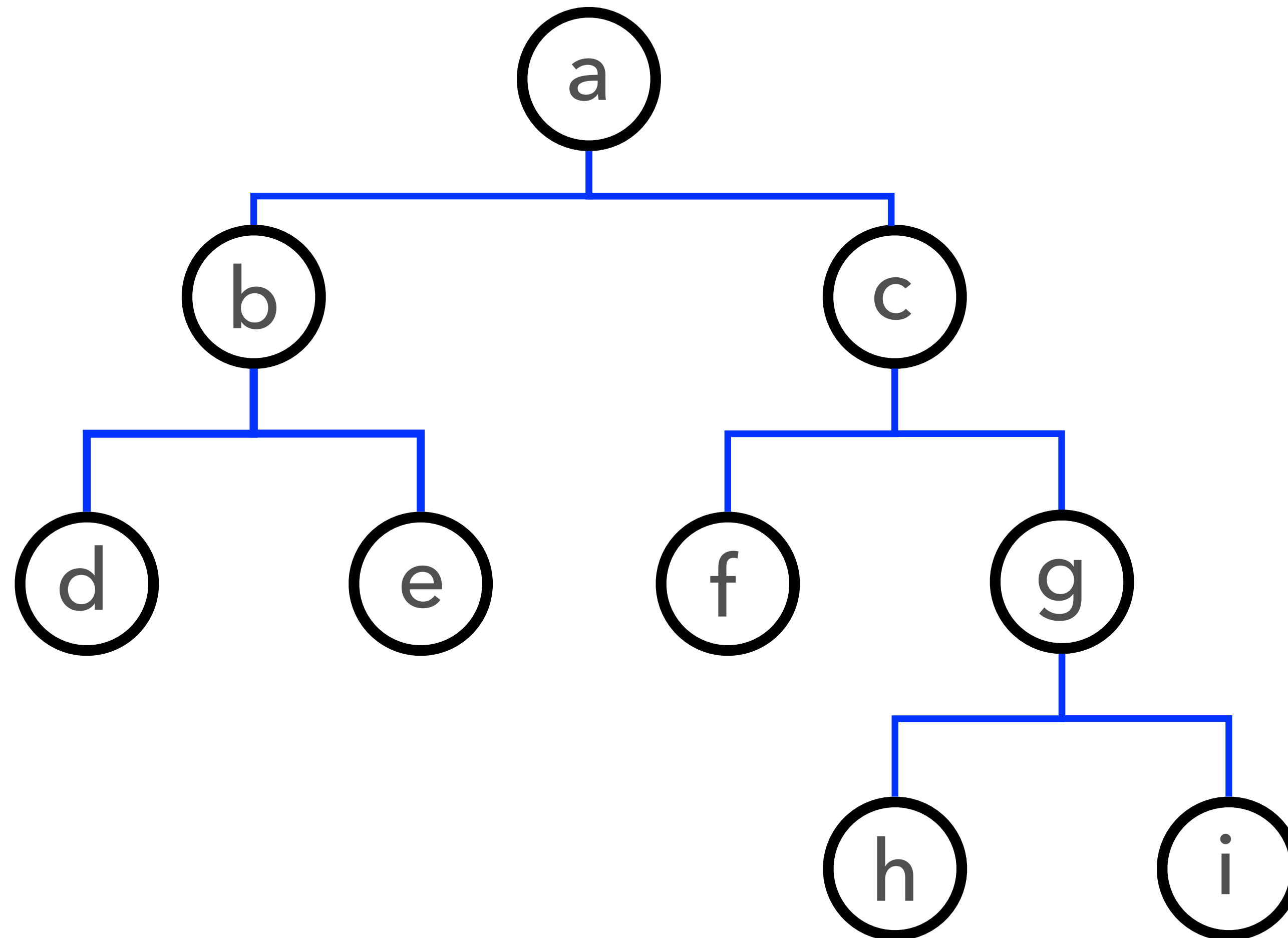C(bacteria, mammals) = 14 ...

# Hierarchical Cost*

# Hierarchical Cost*

C(a,e) = ?

# Hierarchical Cost*

C(a,e) = 2

# Hierarchical Cost*



C(a,e) = 2

C(b,h) = 4