

# Multiclass Classification

# Logistics

---

- ▶ MidTerm Solution & grade: Friday (April 3)
- ▶ Reading Materials will be uploaded the day before
- ▶ Slides will be uploaded the just before the class

# This Lecture

---

- ▶ Multiclass fundamentals
- ▶ Feature extraction
- ▶ Multiclass logistic regression

# Multiclass Fundamentals

# Text Classification

---

## A Cancer Conundrum: Too Many Drug Trials, Too Few Patients

Breakthroughs in immunotherapy and a rush to develop profitable new treatments have brought a crush of clinical trials scrambling for patients.

By GINA KOLATA



## Yankees and Mets Are on Opposite Tracks This Subway Series

As they meet for a four-game series, the Yankees are playing for a postseason spot, and the most the Mets can hope for is to play spoiler.

By FILIP BONDY



# Text Classification

---

## A Cancer Conundrum: Too Many Drug Trials, Too Few Patients

Breakthroughs in immunotherapy and a rush to develop profitable new treatments have brought a crush of clinical trials scrambling for patients.

By GINA KOLATA



→ Health

## Yankees and Mets Are on Opposite Tracks This Subway Series

As they meet for a four-game series, the Yankees are playing for a postseason spot, and the most the Mets can hope for is to play spoiler.

By FILIP BONDY



→ Sports

# Text Classification

---

## A Cancer Conundrum: Too Many Drug Trials, Too Few Patients

Breakthroughs in immunotherapy and a rush to develop profitable new treatments have brought a crush of clinical trials scrambling for patients.

By GINA KOLATA

## Yankees and Mets Are on Opposite Tracks This Subway Series

As they meet for a four-game series, the Yankees are playing for a postseason spot, and the most the Mets can hope for is to play spoiler.

By FILIP BONDY



→ Health



→ Sports

~20 classes

# Image Classification

---



→ Dog



→ Car

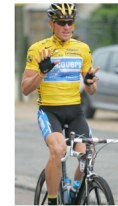
Thousands of classes (ImageNet)



# Entity Linking

---

Although he originally won the event , the United States Anti-Doping Agency announced in August 2012 that they had disqualified **Armstrong** from his seven consecutive Tour de France wins from 1999–2005.



Lance Edward Armstrong is an American former professional road cyclist



Armstrong County is a county in Pennsylvania...

# Entity Linking

Although he originally won the event , the United States Anti-Doping Agency announced in August 2012 that they had disqualified **Armstrong** from his seven consecutive Tour de France wins from 1999–2005.



Lance Edward Armstrong is an American former professional road cyclist



Armstrong County is a county in Pennsylvania...

- ▶ 4,500,000 classes (all articles in Wikipedia)

# Reading Comprehension

---

One day, James thought he would go into town and see what kind of trouble he could get into. He went to the grocery store and pulled all the pudding off the shelves and ate two jars. Then he walked to the fast food restaurant and ordered 15 bags of fries. He didn't pay, and instead headed home.

3) Where did James go after he went to the grocery store?

A) his deck

B) his freezer

C) a fast food restaurant

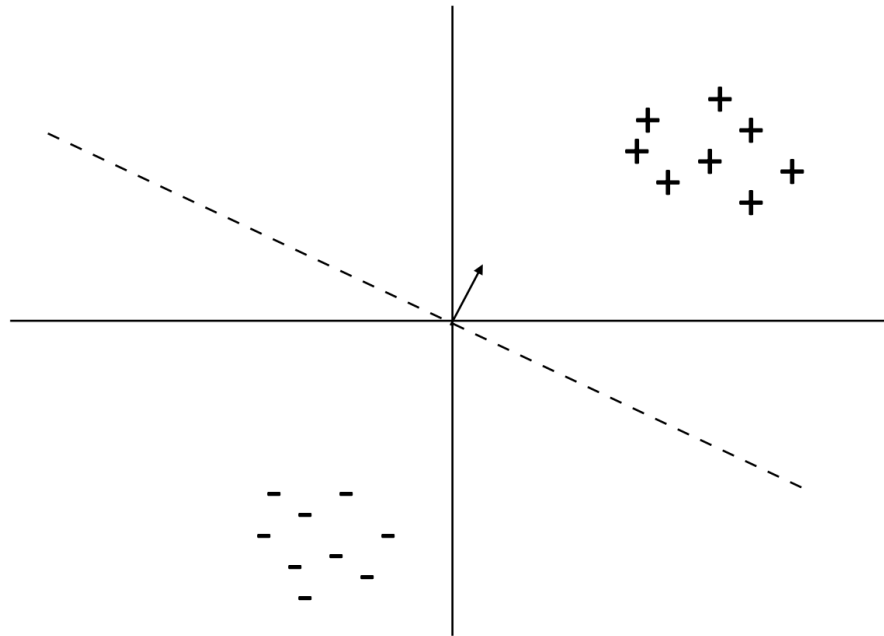
D) his room

► Multiple choice questions, 4 classes (but classes change per example)

# Binary Classification

---

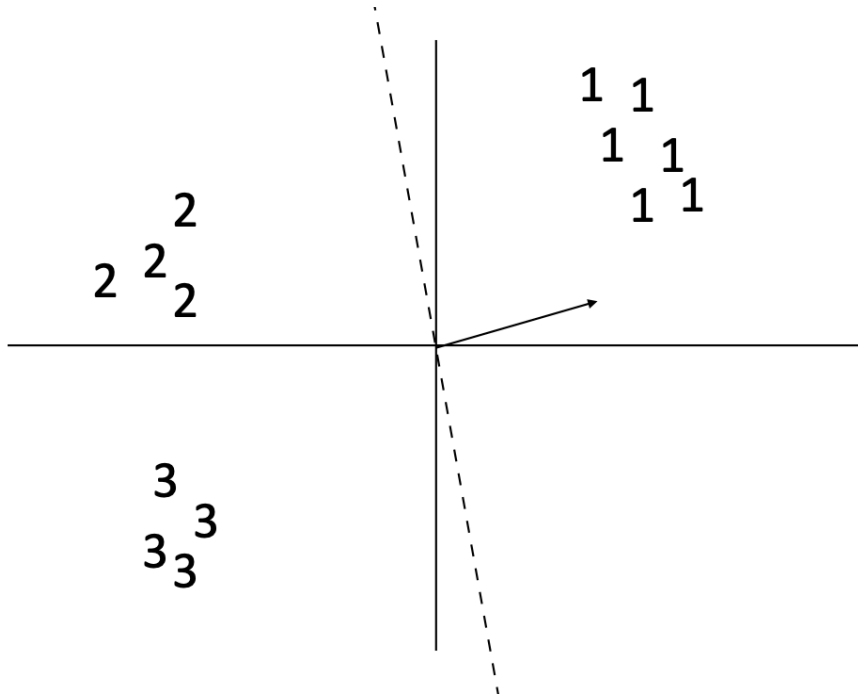
- ▶ Binary classification: one weight vector defines positive and negative classes



# Multiclass Classification

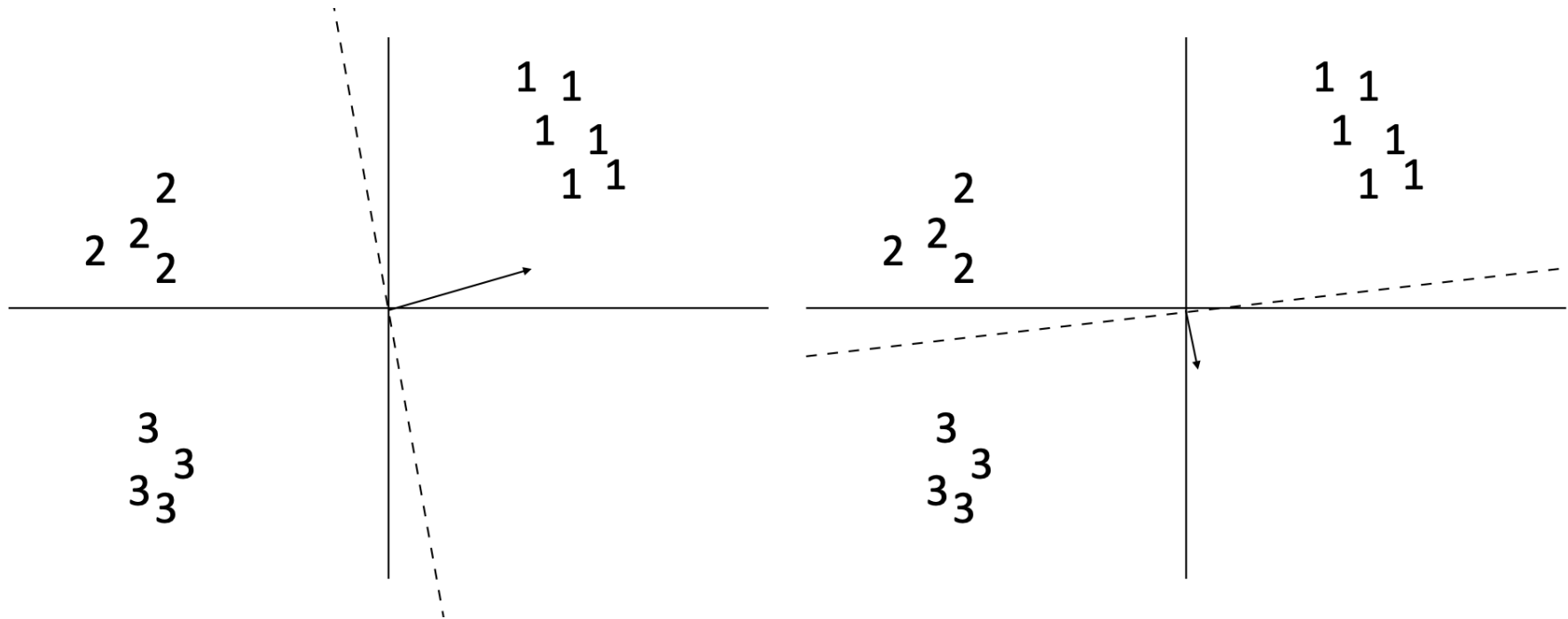
---

- Can we just use binary classifiers here?



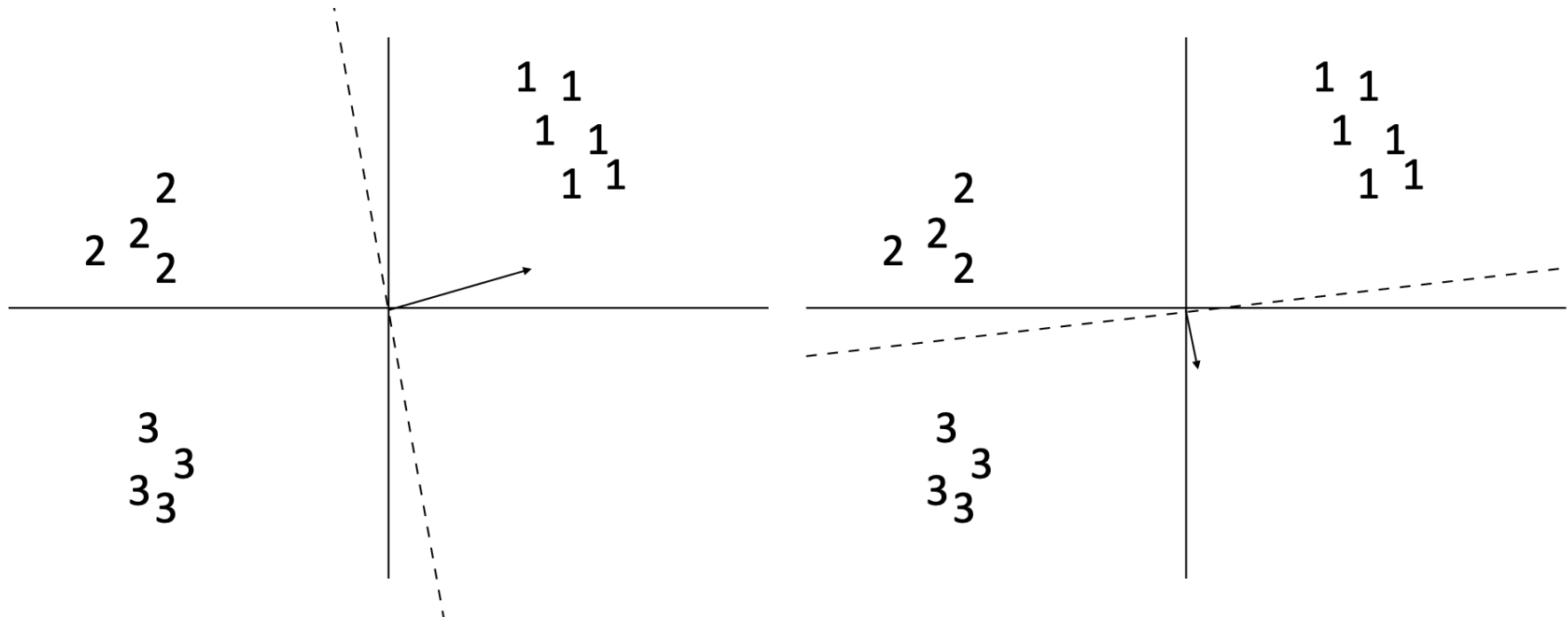
# Multiclass Classification

- One-vs-all: train  $k$  classifiers, one to distinguish each class from all the rest



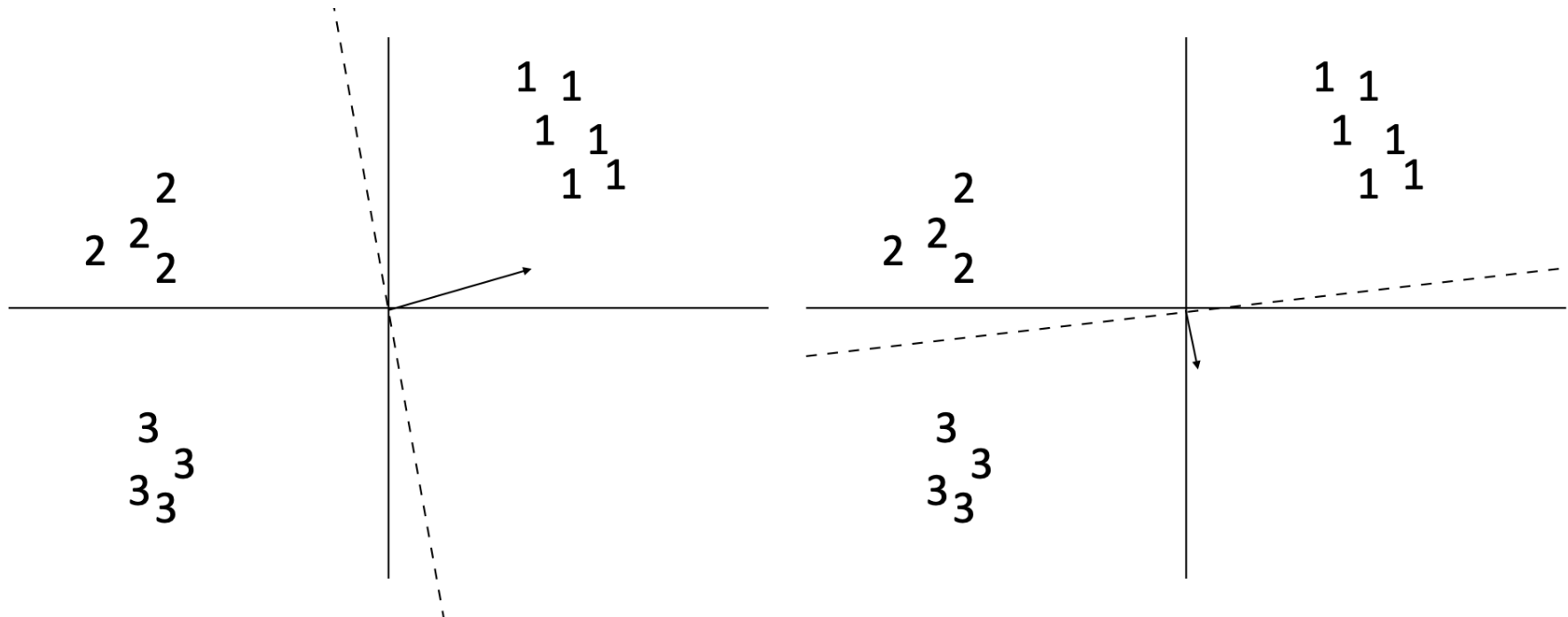
# Multiclass Classification

- ▶ One-vs-all: train  $k$  classifiers, one to distinguish each class from all the rest
- ▶ How do we reconcile multiple positive predictions?



# Multiclass Classification

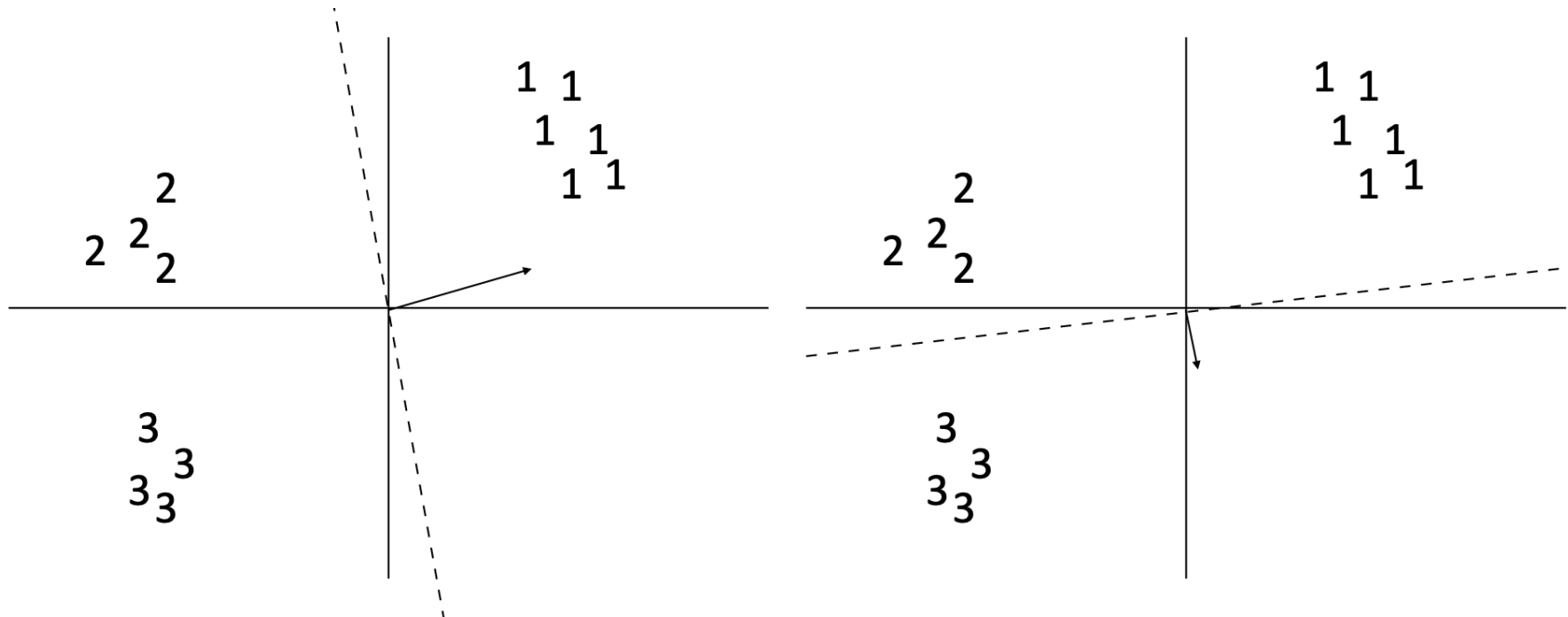
- ▶ One-vs-all: train  $k$  classifiers, one to distinguish each class from all the rest
- ▶ How do we reconcile multiple positive predictions? **Highest score**





# Multiclass Classification

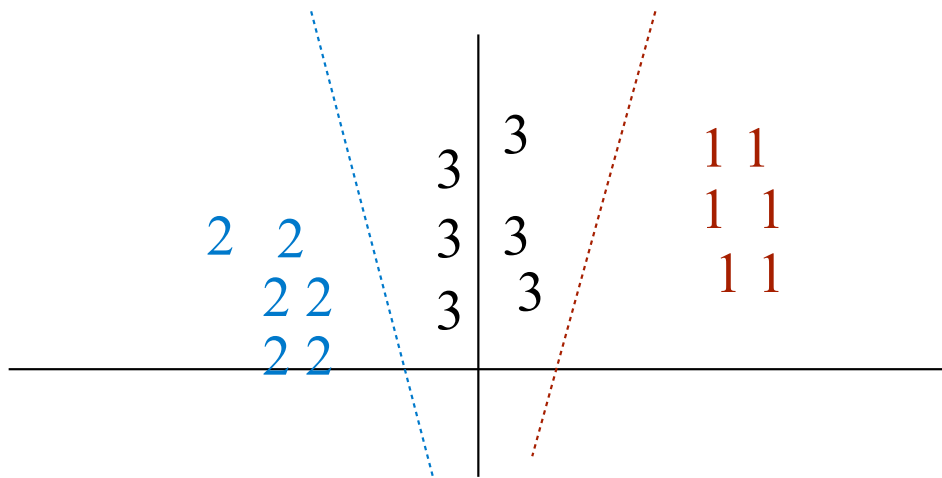
- ▶ All-vs-all: train  $n(n-1)/2$  classifiers to differentiate each pair of classes
- ▶ Again, how to reconcile?



# Multiclass Classification

---

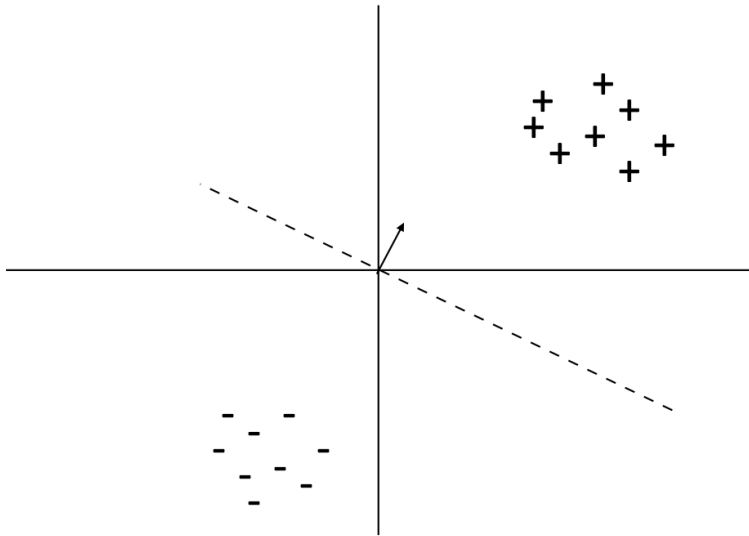
- ▶ Not all classes may even be separable using this approach



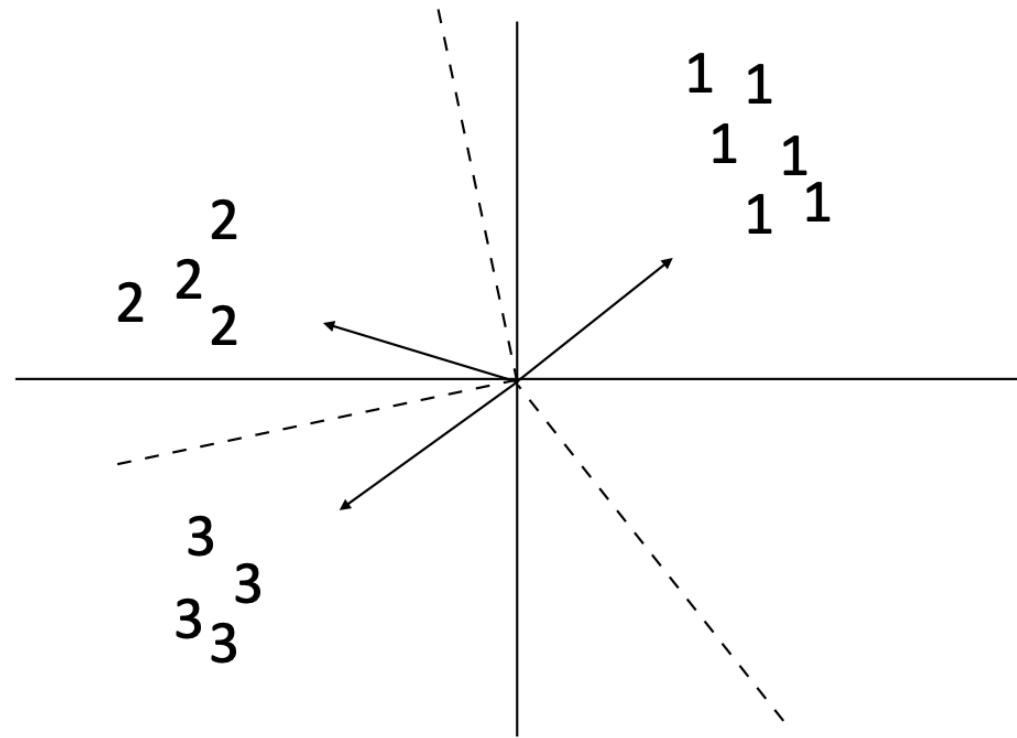
- ▶ Can separate 1 from 2+3 and 2 from 1+3 but not 3 from the others (with these features)

# Multiclass Classification

- ▶ Binary classification: one weight vector defines both classes



- ▶ Multiclass classification: different weights and/or features per class



# Multiclass Classification

---

- ▶ Formally: instead of two labels, we have an output space  $\mathcal{Y}$  containing a number of possible classes
-

# Multiclass Classification

---

- ▶ Formally: instead of two labels, we have an output space  $\mathcal{Y}$  containing a number of possible classes
- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

# Multiclass Classification

---

- ▶ Formally: instead of two labels, we have an output space  $\mathcal{Y}$  containing a number of possible classes
- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$ 
  - ▶ Multiple feature vectors, one weight vector

# Multiclass Classification

---

- ▶ Formally: instead of two labels, we have an output space  $\mathcal{Y}$  containing a number of possible classes

- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$ 
  - ▶ Multiple feature vectors, one weight vector

features depend on choice of label now! note: this isn't the gold label

# Multiclass Classification

---

- ▶ Formally: instead of two labels, we have an output space  $\mathcal{Y}$  containing a number of possible classes

- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

features depend on choice of label now! note: this isn't the gold label

- ▶ Can also have one weight vector per class:  $\operatorname{argmax}_{y \in \mathcal{Y}} w_y^\top f(x)$



# Feature Extraction

# Block Feature Vectors

---

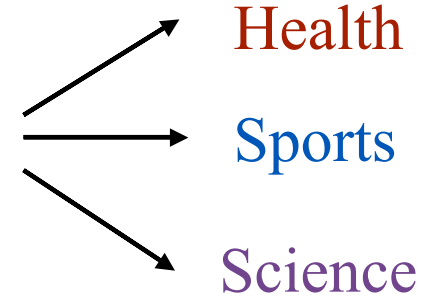
- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

# Block Feature Vectors

---

- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

*too many drug trials, too few patients*

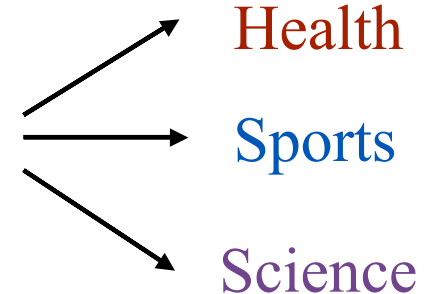


# Block Feature Vectors

---

- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

*too many drug trials, too few patients*



- ▶ Base feature function:

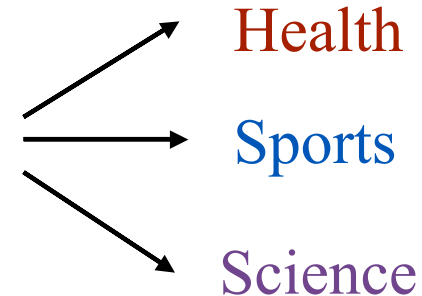
$$f(x) = \text{I}[\text{contains } \textit{drug}], \text{I}[\text{contains } \textit{patients}], \text{I}[\text{contains } \textit{baseball}] = [1, 1, 0]$$

# Block Feature Vectors

---

- ▶ Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

*too many drug trials, too few patients*



- ▶ Base feature function:

$$f(x) = \text{I}[\text{contains } drug], \text{I}[\text{contains } patients], \text{I}[\text{contains } baseball] = [1, 1, 0]$$

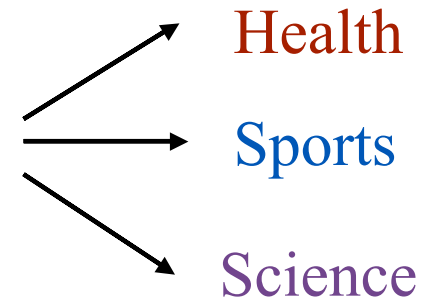
feature vector blocks for each label

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0] \text{I}[\text{contains } drug \text{ \& label} = \text{Health}]$$

# Block Feature Vectors

- Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

*too many drug trials, too few patients*



- Base feature function:

$$f(x) = \text{I}[\text{contains } drug], \text{I}[\text{contains } patients], \text{I}[\text{contains } baseball] = [1, 1, 0]$$

feature vector blocks for each label

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0] \text{I}[\text{contains } drug \text{ \& label = Health}]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0] \text{I}[\text{contains } drug \text{ \& label = Sports}]$$

# Block Feature Vectors

- Decision rule:  $\operatorname{argmax}_{y \in \mathcal{Y}} w^\top f(x, y)$

*too many drug trials, too few patients*

Health

Sports

Science

- Base feature function:

$$f(x) = \mathbb{I}[\text{contains drug}], \mathbb{I}[\text{contains patients}], \mathbb{I}[\text{contains baseball}] = [1, 1, 0]$$

feature vector blocks for each label

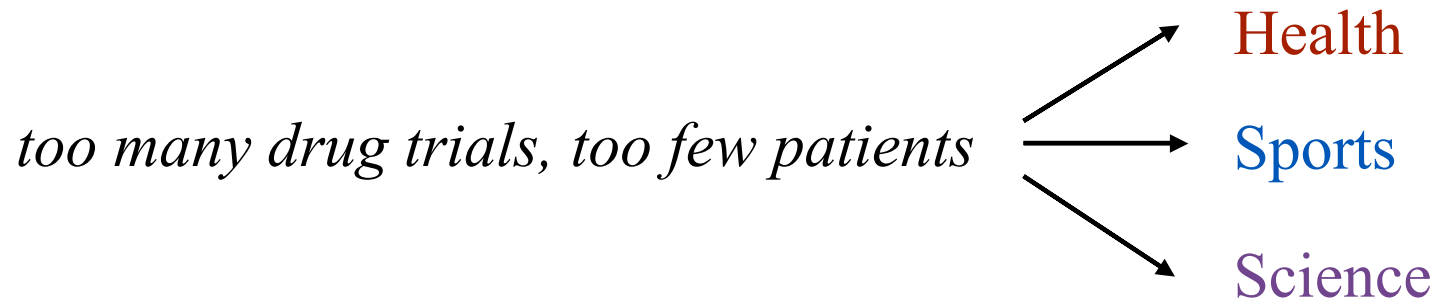
$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0] \quad \mathbb{I}[\text{contains drug} \ \& \ \text{label} = \text{Health}]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0] \quad \mathbb{I}[\text{contains drug} \ \& \ \text{label} = \text{Sports}]$$

- Equivalent to having three weight vectors in this case

# Making Decisions

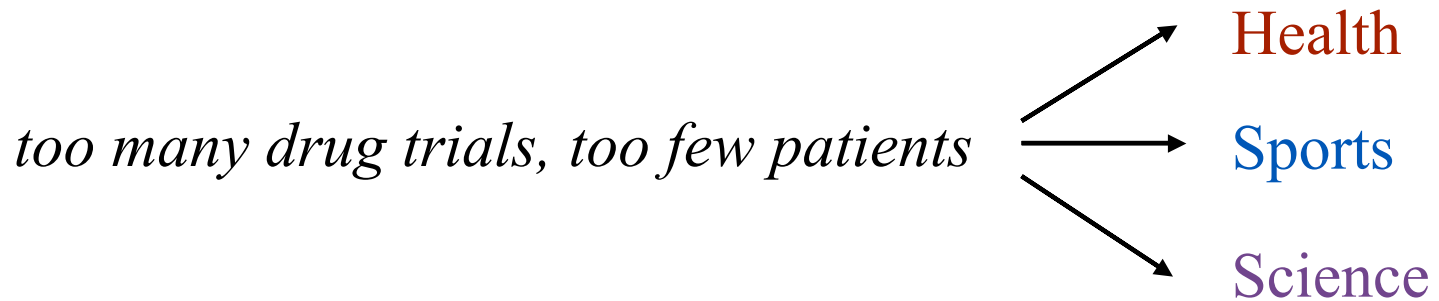
---





# Making Decisions

---



$f(x) = \text{I}[\text{contains } \textit{drug}], \text{I}[\text{contains } \textit{patients}], \text{I}[\text{contains } \textit{baseball}]$

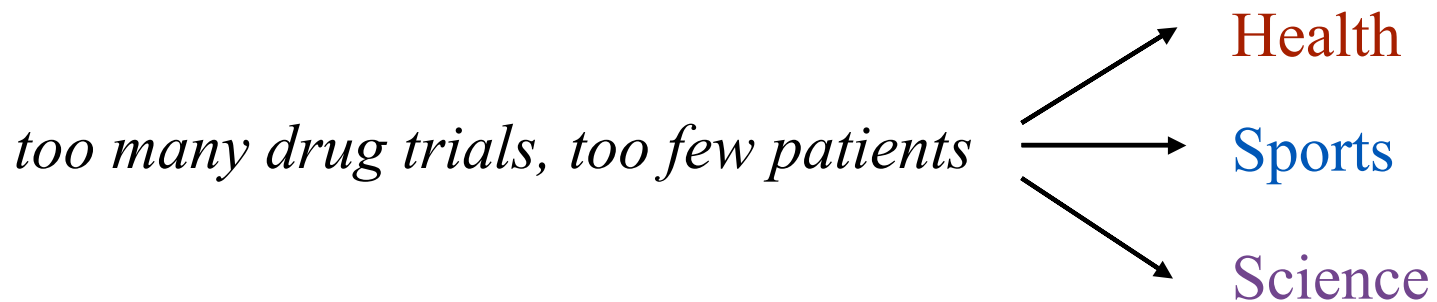
$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

“word drug in Science article” = +1.1

# Making Decisions

---



$f(x) = \text{I}[\text{contains } \textit{drug}], \text{I}[\text{contains } \textit{patients}], \text{I}[\text{contains } \textit{baseball}]$

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$$

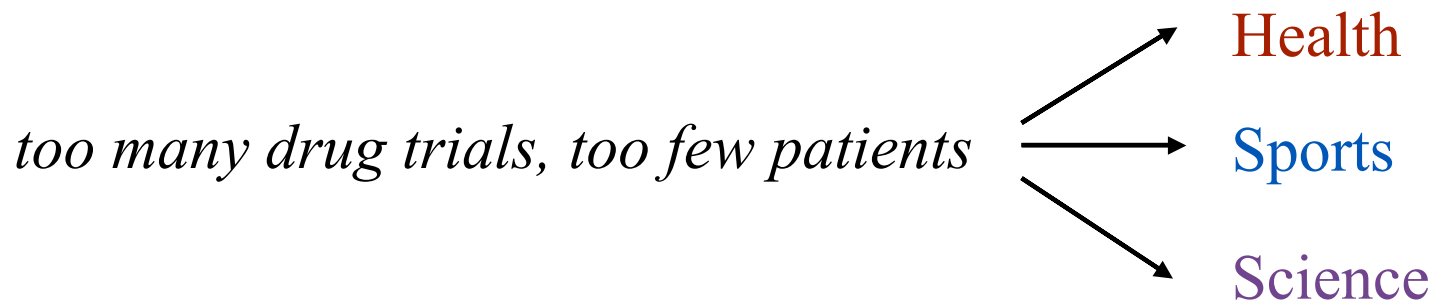
$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

$$w = [+2.1, +2.3, -5, -2.1, -3.8, 0, +1.1, -1.7, -1.3]$$

“word drug in Science article” = +1.1

# Making Decisions

---



$f(x) = \text{I}[\text{contains } \textit{drug}], \text{I}[\text{contains } \textit{patients}], \text{I}[\text{contains } \textit{baseball}]$

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

“word drug in Science article” = +1.1

$$w = [+2.1, +2.3, 0, -2.1, -3.8, 0, +1.1, -1.7, -1.3]$$

$$w^\top f(x, y) = \text{Health: } +4.4 \quad \text{Sports: } -5.9 \quad \text{Science: } -1.9$$

# Another example: POS tagging

---

- ▶ Classify *blocks* as one of 36 POS tags

*the router*

*blocks*

*the packets*

NNS

VBZ

NN

DT

...

## Another example: POS tagging

---

- ▶ Classify *blocks* as one of 36 POS tags      *the router*      *blocks*      *the packets*
- ▶ Example  $x$  : sentence with a word (in this case, *blocks*) highlighted

NNS

VBZ

NN

DT

...

# Another example: POS tagging

- ▶ Classify *blocks* as one of 36 POS tags
- ▶ Example  $x$  : sentence with a word (in this case, *blocks*) highlighted
- ▶ Extract features with respect to this word:

*the router blocks the packets*

NNS

VBZ

NN

DT

...

$f(x, y = \text{VBZ}) = \text{I}[\text{curr\_word} = \text{blocks} \ \& \ \text{tag} = \text{VBZ}],$

$\text{I}[\text{prev\_word} = \text{router} \ \& \ \text{tag} = \text{VBZ}]$

$\text{I}[\text{next\_word} = \text{the} \ \& \ \text{tag} = \text{VBZ}]$

$\text{I}[\text{curr\_suffix} = \text{s} \ \& \ \text{tag} = \text{VBZ}]$

not saying that *the* is tagged as VBZ! saying that *the* follows the VBZ word

# Multiclass Logistic Regression

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

↗  
sum over output  
space to normalize



# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

► Compare to binary:

$$P(y = 1|x) = \frac{\exp(w^\top f(x))}{1 + \exp(w^\top f(x))}$$

negative class implicitly had  
 $f(x, y=0) = \text{the zero vector}$

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$



Softmax  
function

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

Softmax  
function

↗  
sum over output  
space to normalize

# Multiclass Logistic Regression

---

Softmax  
function

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

Why? Interpret raw classifier scores as **probabilities**

# Multiclass Logistic Regression

Softmax  
function

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

Why? Interpret raw classifier scores as **probabilities**

*too many drug trials,  
too few patients*

Health: +4.4

Sports: -5.9 .

Science: -1.9

$w^\top f(x, y)$

# Multiclass Logistic Regression

Softmax  
function

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

Why? Interpret raw classifier scores as **probabilities**

*too many drug trials,  
too few patients*

Health: +4.4

Sports: -5.9

Science: -1.9

$w^\top f(x, y)$

exp →

81.45

0.002

0.014

# Multiclass Logistic Regression

Softmax  
function

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

Why? Interpret raw classifier scores as **probabilities**

too many drug trials,  
too few patients

Health: +4.4  
Sports: -5.9  
Science: -1.9

$w^\top f(x, y)$

exp

81.45  
0.002  
0.014

normalize

probabilities  
must sum to 1

0.99  
 $2.4 \times 10^{-5}$   
 $1.7 \times 10^{-4}$

probabilities



# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

↗  
sum over output  
space to normalize

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

↑  
sum over output  
space to normalize

► Training: maximize  $\mathcal{L}(x, y) = \sum_{j=1}^n \log P(y_j^* | x_j)$

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

↑  
sum over output  
space to normalize

► Training: maximize  $\mathcal{L}(x, y) = \sum_{j=1}^n \log P(y_j^* | x_j)$

$$= \sum_{j=1}^n \left( w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y)) \right)$$

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

↑  
sum over output  
space to normalize

► Training: maximize  $\mathcal{L}(x, y) = \sum_{j=1}^n \log P(y_j^* | x_j)$

$$= \sum_{j=1}^n \left( w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y)) \right)$$

# Multiclass Logistic Regression

---

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

sum over output  
space to normalize

► Training: maximize  $\mathcal{L}(x, y) = \sum_{j=1}^n \log P(y_j^* | x_j)$

$$= \sum_{j=1}^n \left( w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y)) \right)$$

# Training

---

► Multiclass logistic regression  $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$

► Likelihood  $\mathcal{L}(x_j, y_j^*) = w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y))$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \frac{\sum_y f_i(x_j, y) \exp(w^\top f(x_j, y))}{\sum_y \exp(w^\top f(x_j, y))}$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

# Training

---

► Multiclass logistic regression  $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$

► Likelihood  $\mathcal{L}(x_j, y_j^*) = w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y))$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \frac{\sum_y f_i(x_j, y) \exp(w^\top f(x_j, y))}{\sum_y \exp(w^\top f(x_j, y))}$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = \underbrace{f_i(x_j, y_j^*)}_{\text{gold feature value}} - \underbrace{\mathbb{E}_y[f_i(x_j, y)]}_{\text{model's expectation of feature value}}$$

# Training

► Multiclass logistic regression  $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$

► Likelihood  $\mathcal{L}(x_j, y_j^*) = w^\top f(x_j, y_j^*) - \log \sum_y \exp(w^\top f(x_j, y))$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \frac{\sum_y f_i(x_j, y) \exp(w^\top f(x_j, y))}{\sum_y \exp(w^\top f(x_j, y))}$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = \underbrace{f_i(x_j, y_j^*)}_{\text{gold feature value}} - \underbrace{\mathbb{E}_y[f_i(x_j, y)]}_{\text{model's expectation of feature value}}$$



# Training

---

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

# Training

---

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

*too many drug trials, too few patients*

$y^* = \text{Health}$

$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$

$P_w(y|x) = [0.21, 0.77, 0.02]$

$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$

# Training

---

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

*too many drug trials, too few patients*

$y^* = \text{Health}$

$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$

$P_w(y|x) = [0.21, 0.77, 0.02]$

$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$

gradient:  $[1, 1, 0, 0, 0, 0, 0, 0, 0] - 0.21 [1, 1, 0, 0, 0, 0, 0, 0, 0]$   
 $- 0.77 [0, 0, 0, 1, 1, 0, 0, 0, 0] - 0.02 [0, 0, 0, 0, 0, 0, 1, 1, 0]$   
 $= [0.79, 0.79, 0, -0.77, -0.77, 0, -0.02, -0.02, 0]$

# Training

---

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

*too many drug trials, too few patients*

$y^* = \text{Health}$

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$P_w(y|x) = [0.21, 0.77, 0.02]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

$$\begin{aligned} \text{gradient: } & [1, 1, 0, 0, 0, 0, 0, 0, 0] - 0.21 [1, 1, 0, 0, 0, 0, 0, 0, 0] \\ & - 0.77 [0, 0, 0, 1, 1, 0, 0, 0, 0] - 0.02 [0, 0, 0, 0, 0, 0, 1, 1, 0] \\ & = [0.79, 0.79, 0, -0.77, -0.77, 0, -0.02, -0.02, 0] \end{aligned}$$

update  $w^\top$ :

$$\begin{aligned} & [1.3, 0.9, -5, 3.2, -0.1, 0, 1.1, -1.7, -1.3] + [0.79, 0.79, 0, -0.77, -0.77, 0, -0.02, -0.02, 0] \\ & = [2.09, 1.69, 0, 2.43, -0.87, 0, 1.08, -1.72, 0] \quad \rightarrow \text{new } P_w(y|x) = [0.89, 0.10, 0.01] \end{aligned}$$

# Training

$$\frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) = f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y|x_j)$$

*too many drug trials, too few patients*

$y^* = \text{Health}$

$$f(x, y = \text{Health}) = [1, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$P_w(y|x) = [0.21, 0.77, 0.02]$$

$$f(x, y = \text{Sports}) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

$$\begin{aligned} \text{gradient: } & [1, 1, 0, 0, 0, 0, 0, 0, 0] - 0.21 [1, 1, 0, 0, 0, 0, 0, 0, 0] \\ & - 0.77 [0, 0, 0, 1, 1, 0, 0, 0, 0] - 0.02 [0, 0, 0, 0, 0, 0, 1, 1, 0] \\ & = [0.79, 0.79, 0, -0.77, -0.77, 0, -0.02, -0.02, 0] \end{aligned}$$

update  $w^\top$ :

$$\begin{aligned} & [1.3, 0.9, -5, 3.2, -0.1, 0, 1.1, -1.7, -1.3] + [0.79, 0.79, 0, -0.77, -0.77, 0, -0.02, -0.02, 0] \\ & = [2.09, 1.69, -5, 2.43, -0.87, 0, 1.08, -1.72, -1.3] \end{aligned}$$

$$P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$$

Softmax function

↪ new  $P_w(y|x) = [0.89, 0.10, 0.01]$

# Logistic Regression: Summary

---

- ▶ Model:  $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$
- ▶ Inference:  $\operatorname{argmax}_y P_w(y|x)$
- ▶ Learning: gradient ascent on the discriminative log-likelihood

$$f(x, y^*) - \mathbb{E}_y[f(x, y)] = f(x, y^*) - \sum_y [P_w(y|x) f(x, y)]$$

“towards gold feature value, away from expectation of feature value”