

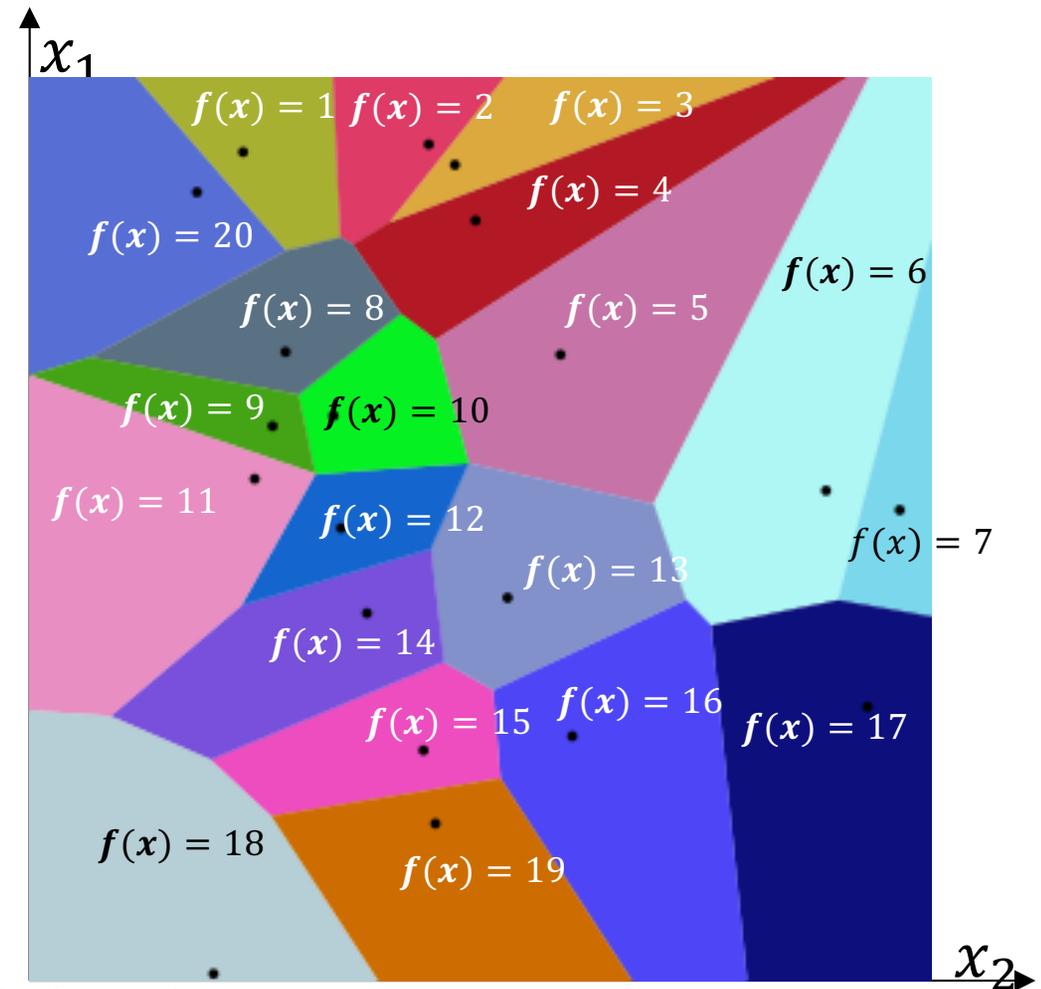
CS440/ECE448

Lecture 13:

Softmax

Mark Hasegawa-Johnson

These slides are in the public domain. Re-use, remix, redistribute at will.



By Balu Ertl - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=38534275>

Outline

- Sigmoid review
- Multi-class linear classifiers
- One-hot vectors
- Softmax nonlinearity
- Derivative of the log softmax

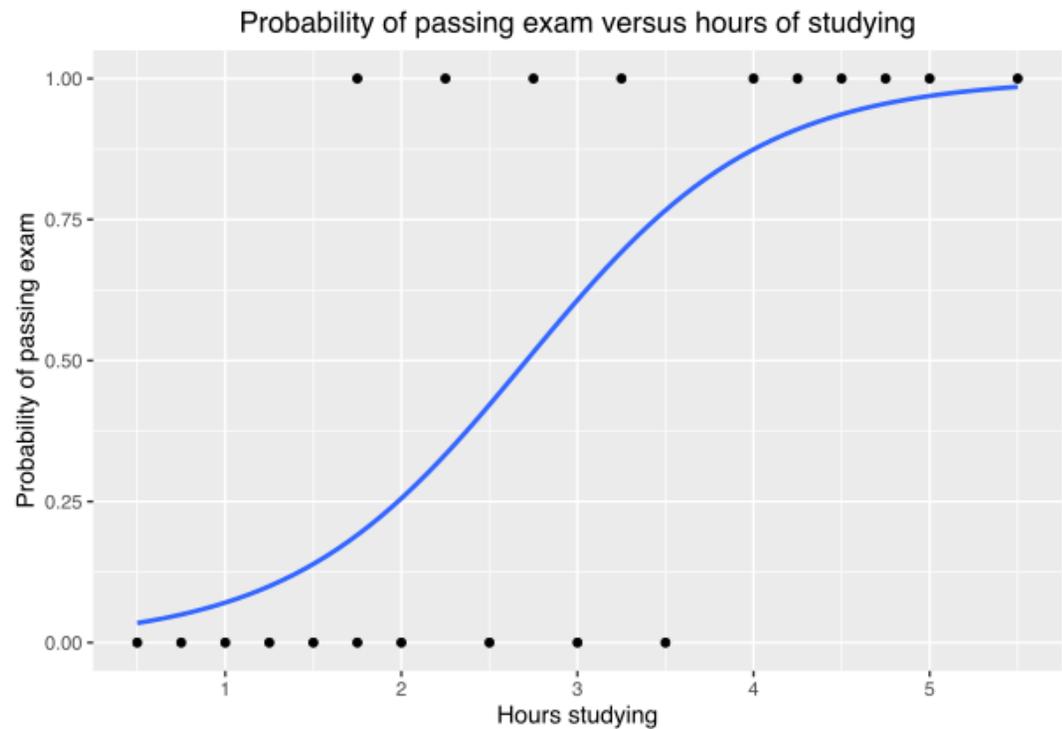
Logistic regression: Output is a probability

- The label is either $y = 0$ or $y = 1$
- Model:

$$P(Y = 1|\mathbf{x}) = f(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w})$$

- Logistic sigmoid function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Loss function: Binary cross entropy

- Loss function

$$\begin{aligned}\mathcal{L} &= - \sum_{i=1}^n \log P(Y = y_i | \mathbf{x}_i) \\ &= \sum_{y_i=1} -\log \sigma(\mathbf{x}_i^T \mathbf{w}) + \sum_{y_i=0} -\log(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))\end{aligned}$$

- Derivative of the sigmoid

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

Loss Gradient has an interesting simplification

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= -\frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^n \log P(Y = y_i | \mathbf{x}_i) \\ &= \sum_{y_i=1} -\frac{\partial}{\partial \mathbf{w}} \log \sigma(\mathbf{x}_i^T \mathbf{w}) + \sum_{y_i=0} -\frac{\partial}{\partial \mathbf{w}} \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \\ &= \sum_{y_i=1} -(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i + \sum_{y_i=0} \sigma(\mathbf{x}_i^T \mathbf{w}) \mathbf{x}_i \\ &= -\sum_{i=1}^n (y_i - \sigma(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i\end{aligned}$$

Outline

- Sigmoid review
- Multi-class linear classifiers
- One-hot vectors
- Softmax nonlinearity
- Derivative of the log softmax

Linear classifier: Notation

- The observation $\mathbf{x}^T = [x_1, \dots, x_d]$ is a real-valued vector (d is the number of feature dimensions)
- The class label $y \in \mathcal{Y}$ is drawn from some finite set of class labels.
- Usually the output vocabulary, \mathcal{Y} , is some set of strings. For convenience, though, we usually map the class labels to a sequence of integers, $\mathcal{Y} = \{1, \dots, v\}$, where v is the vocabulary size.

Linear classifier: Definition

A linear classifier is defined by

$$f(\mathbf{x}) = \operatorname{argmax} \mathbf{W}\mathbf{x}$$

where:

$$\mathbf{W}\mathbf{x} = \begin{bmatrix} w_{1,1} & \cdots & w_{1,d} \\ \vdots & \ddots & \vdots \\ w_{v,1} & \cdots & w_{v,d} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \vdots \\ \mathbf{w}_v^T \mathbf{x} \end{bmatrix}$$

\mathbf{w}_k is the weight vector corresponding to class k, and the argmax function finds the element of the vector $\mathbf{W}\mathbf{x}$ with the largest value.

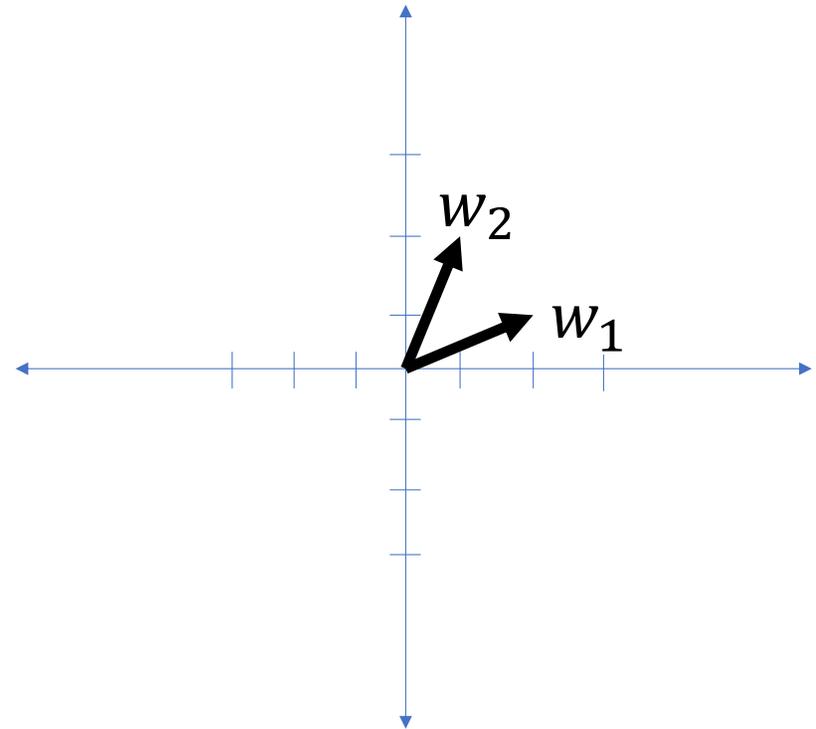
There are a total of vd trainable parameters: the elements of the matrix \mathbf{W} .

Example

Consider a two-class classification problem, i.e., $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2]^T$, where:

$$\mathbf{w}_1^T = [w_{1,1}, w_{1,2}] = [2, 1]$$

$$\mathbf{w}_2^T = [w_{2,1}, w_{2,2}] = [1, 2]$$



Example

Notice that in the two-class case, the equation

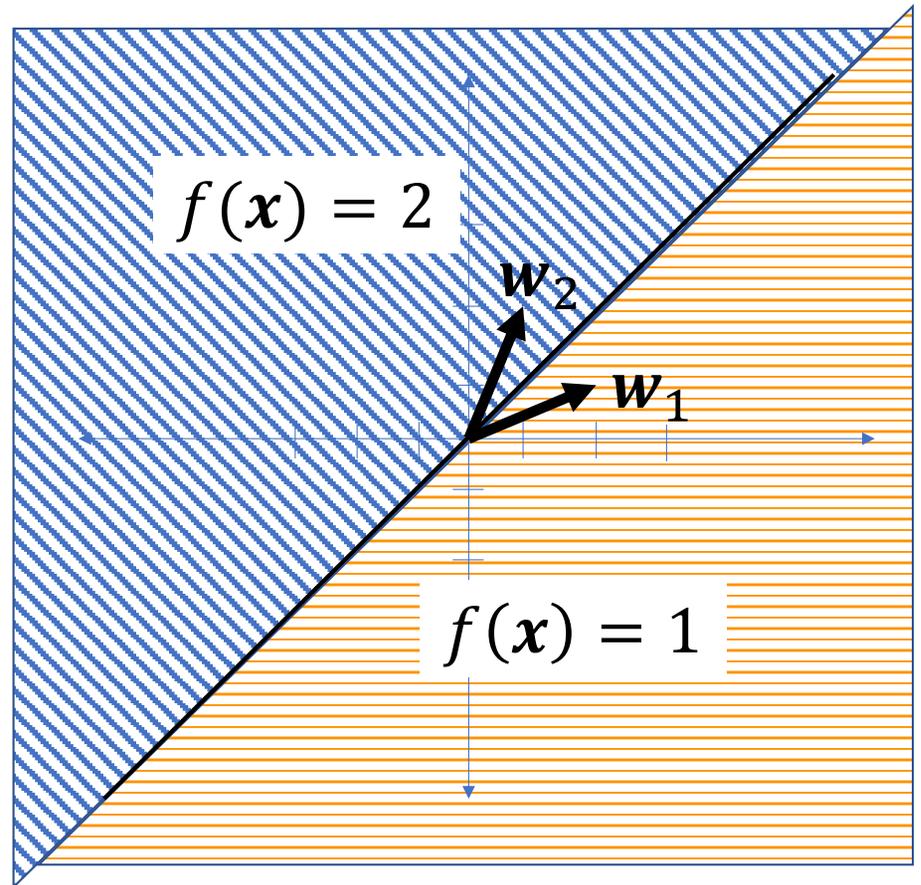
$$f(\mathbf{x}) = \operatorname{argmax} \mathbf{W}\mathbf{x}$$

Simplifies to

$$f(\mathbf{x}) = \begin{cases} 1 & \mathbf{w}_1^T \mathbf{x} > \mathbf{w}_2^T \mathbf{x} \\ 2 & \mathbf{w}_1^T \mathbf{x} < \mathbf{w}_2^T \mathbf{x} \end{cases}$$

The class boundary is the line whose equation is

$$(\mathbf{w}_2 - \mathbf{w}_1)^T \mathbf{x} = 0$$



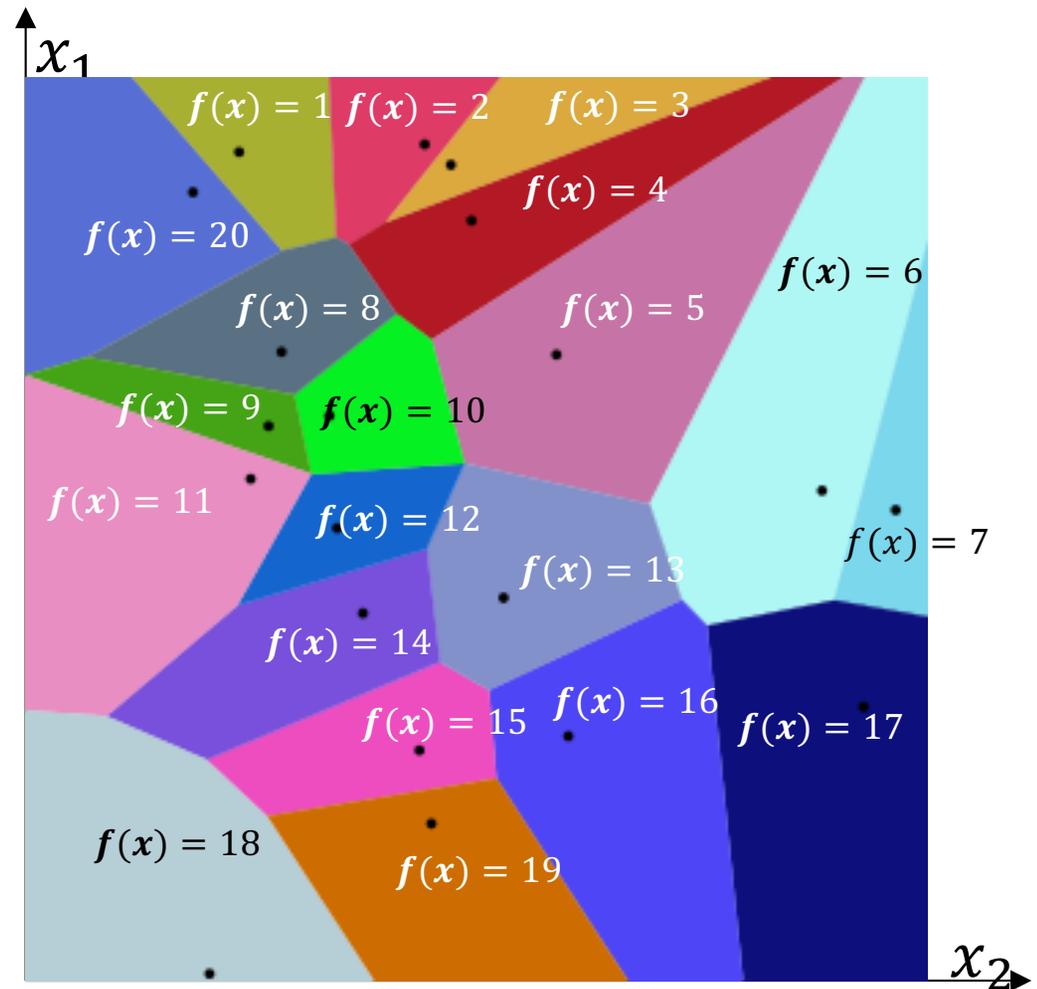
Multi-class linear classifier

In a general multi-class linear classifier,

$$f(\mathbf{x}) = \operatorname{argmax} \mathbf{W}\mathbf{x}$$

The boundary between class k and class l is the line (or plane, or hyperplane) given by the equation

$$(\mathbf{w}_k - \mathbf{w}_l)^T \mathbf{x} = 0$$

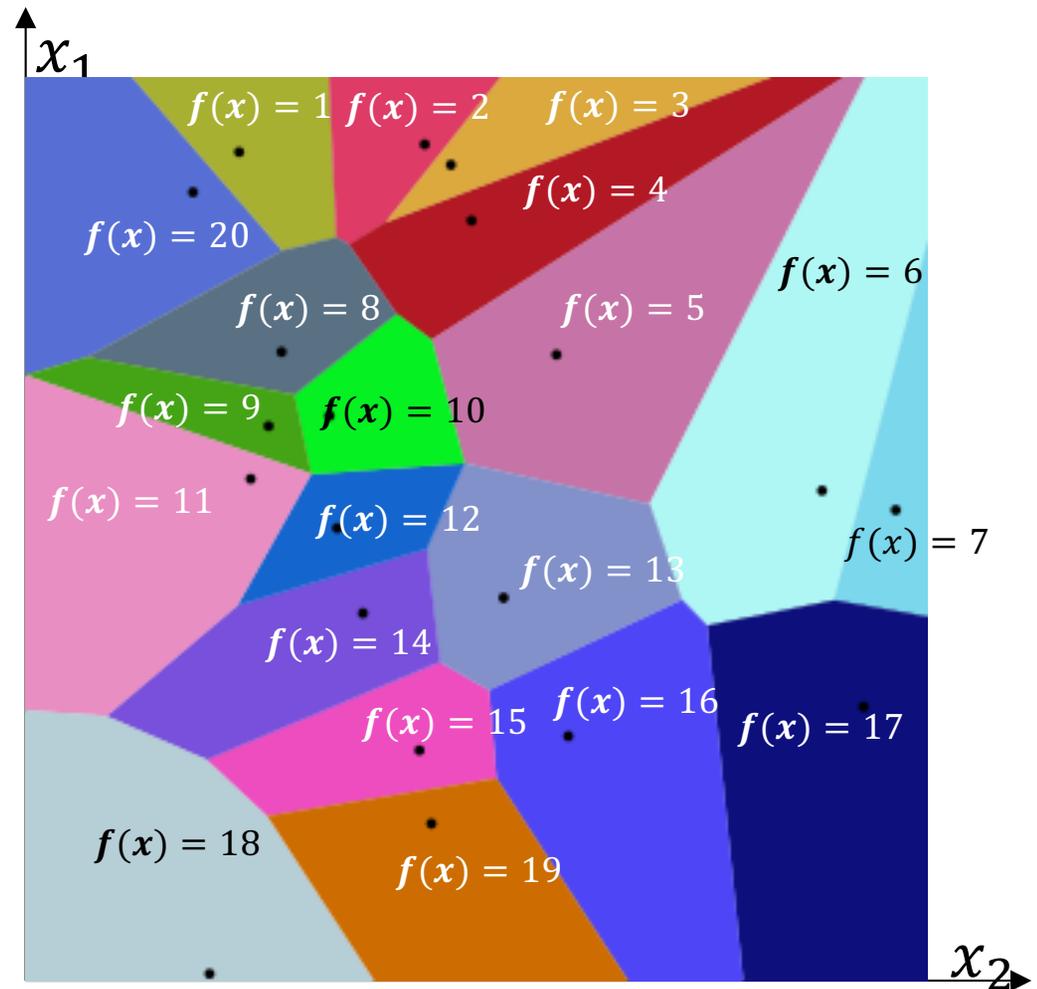


Voronoi regions

The classification regions in a linear classifier are called Voronoi regions.

A **Voronoi region** is a region that is

- Convex (if x_1 and x_2 are points in the region, then every point on the line segment connecting them is also in the region)
- Bounded by piece-wise linear boundaries



Outline

- Sigmoid review
- Multi-class linear classifiers
- One-hot vectors
- Softmax nonlinearity
- Derivative of the log softmax

One-hot vectors

A **one-hot vector** is a binary vector in which all elements are 0 except for a single element that's equal to 1.

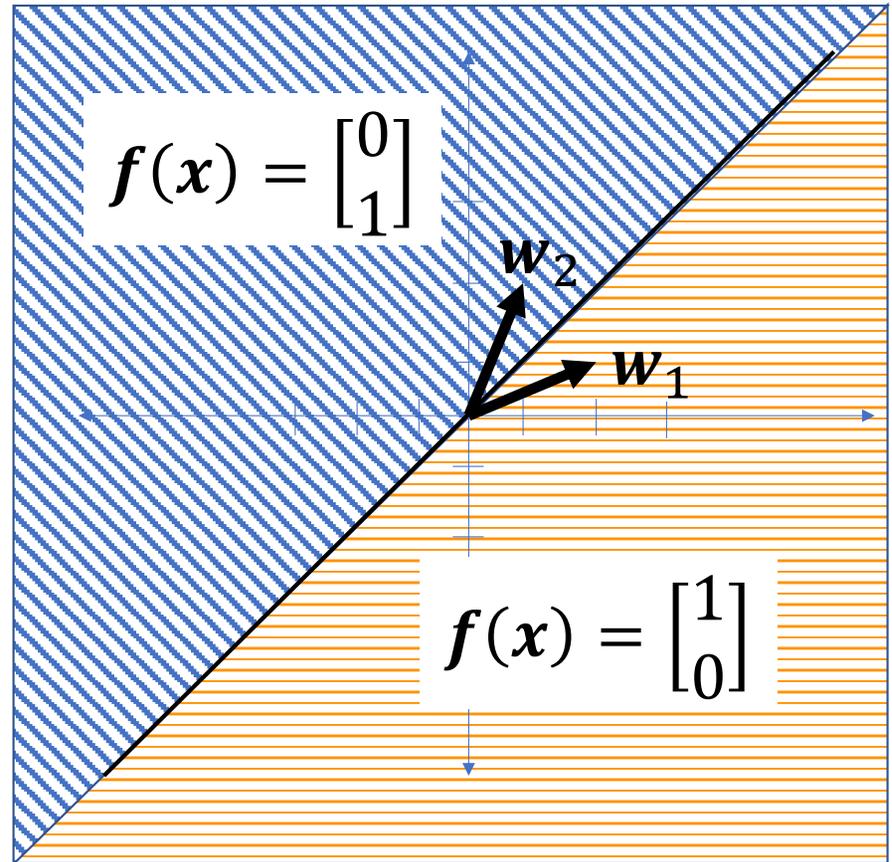
Example: Binary classifier

Consider the classifier

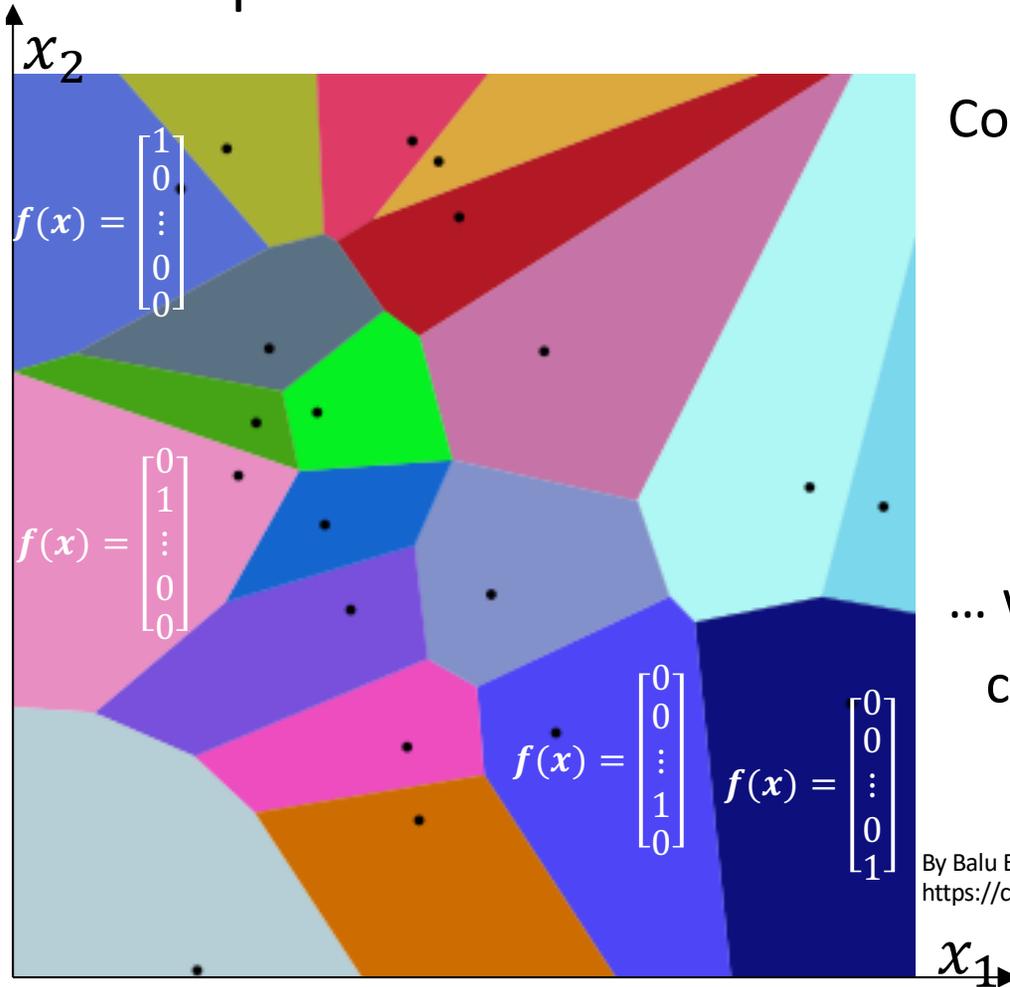
$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\text{argmax } Wx=1} \\ \mathbb{1}_{\text{argmax } Wx=2} \end{bmatrix}$$

...where $\mathbb{1}_P$ is called the “indicator function,” and it means:

$$\mathbb{1}_P = \begin{cases} 1 & P \text{ is true} \\ 0 & P \text{ is false} \end{cases}$$



Example: Multi-Class



Consider the classifier

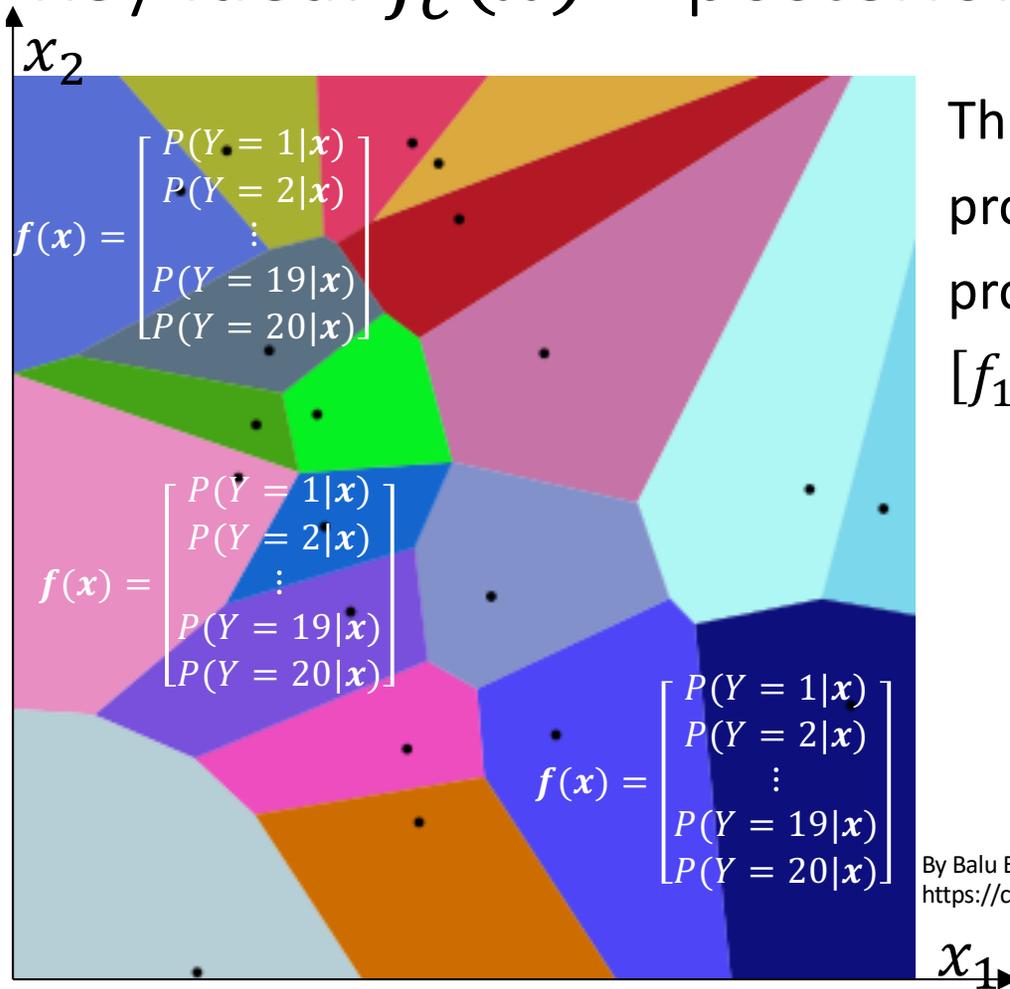
$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_v(x) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\arg\max Wx=1} \\ \vdots \\ \mathbb{1}_{\arg\max Wx=v} \end{bmatrix}$$

... with 20 classes. Then some of the classifications might look like this.

Outline

- Sigmoid review
- Multi-class linear classifiers
- One-hot vectors
- **Softmax nonlinearity**
- **Derivative of the log softmax**

Key idea: $f_c(\mathbf{x})$ = posterior probability of class c



The key idea of a softmax: Instead of producing a one-hot output, the neural net produces a vector $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_v(\mathbf{x})]^T$ such that

$$f_c(\mathbf{x}) = \Pr(Y = c|\mathbf{x})$$

Key idea: $f_c(\mathbf{x})$ = posterior probability of class c

A softmax computes $f_c(\mathbf{x}) = \Pr(Y = c|\mathbf{x})$. The conditions for this to be true are:

1. It needs to satisfy the axioms of probability:

$$0 \leq f_c(\mathbf{x}) \leq 1, \quad \sum_{c=1}^v f_c(\mathbf{x}) = 1$$

2. The weight matrix, \mathbf{W} , is trained using a loss function that encourages $\mathbf{f}(\mathbf{x})$ to approximate posterior probability of the labels on some training dataset:

$$f_c(\mathbf{x}) = \Pr(Y = c|\mathbf{x})$$

Softmax satisfies the axioms of probability

- Axiom #1, probabilities are non-negative ($f_k(\mathbf{x}) \geq 0$). There are many ways to do this, but one way that works is to choose:

$$f_c(\mathbf{x}) \propto \exp(\mathbf{w}_c^T \mathbf{x})$$

- Axiom #2, probabilities should sum to one ($\sum_{k=1}^V f_k(\mathbf{x}) = 1$). This can be done by normalizing:

$$f_c(\mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{k=0}^{V-1} \exp(\mathbf{w}_k^T \mathbf{x})}$$

The softmax function

This is called the softmax function:

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_v(\mathbf{x})]^T$$

$$f_c(\mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{k=1}^v \exp(\mathbf{w}_k^T \mathbf{x})}$$

...where \mathbf{w}_k^T is the k^{th} row of the matrix \mathbf{W} .

Quiz

Go to PrairieLearn, try the quiz!

Outline

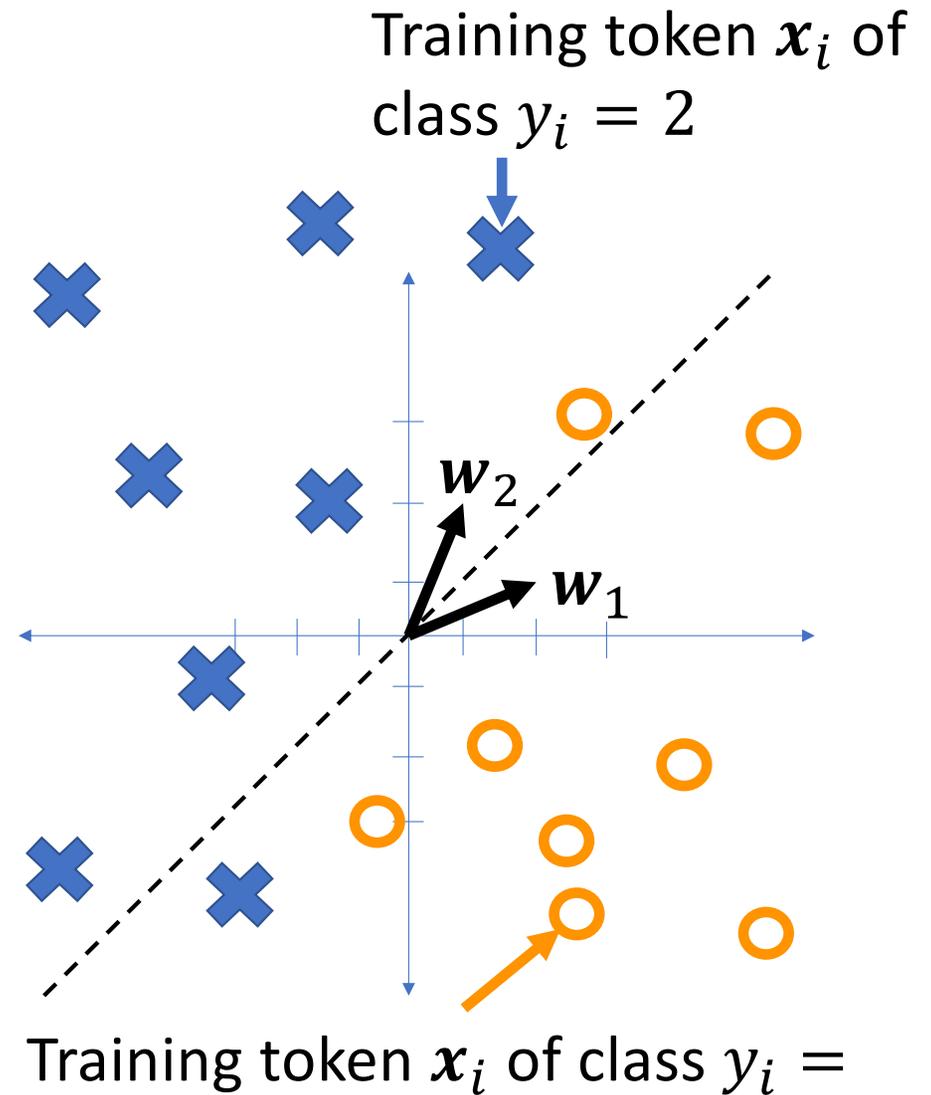
- Sigmoid review
- Multi-class linear classifiers
- One-hot vectors
- Softmax nonlinearity
- **Derivative of the log softmax**

Gradient descent

Suppose we have training tokens (x_i, y_i) , and we have some initial class vectors w_1 and w_2 . We want to update them as

$$w_1 \leftarrow w_1 - \eta \frac{\partial \mathcal{L}}{\partial w_1}$$
$$w_2 \leftarrow w_2 - \eta \frac{\partial \mathcal{L}}{\partial w_2}$$

...where \mathcal{L} is some loss function.
What loss function makes sense?



Cross Entropy Loss = negative log probability of the training set

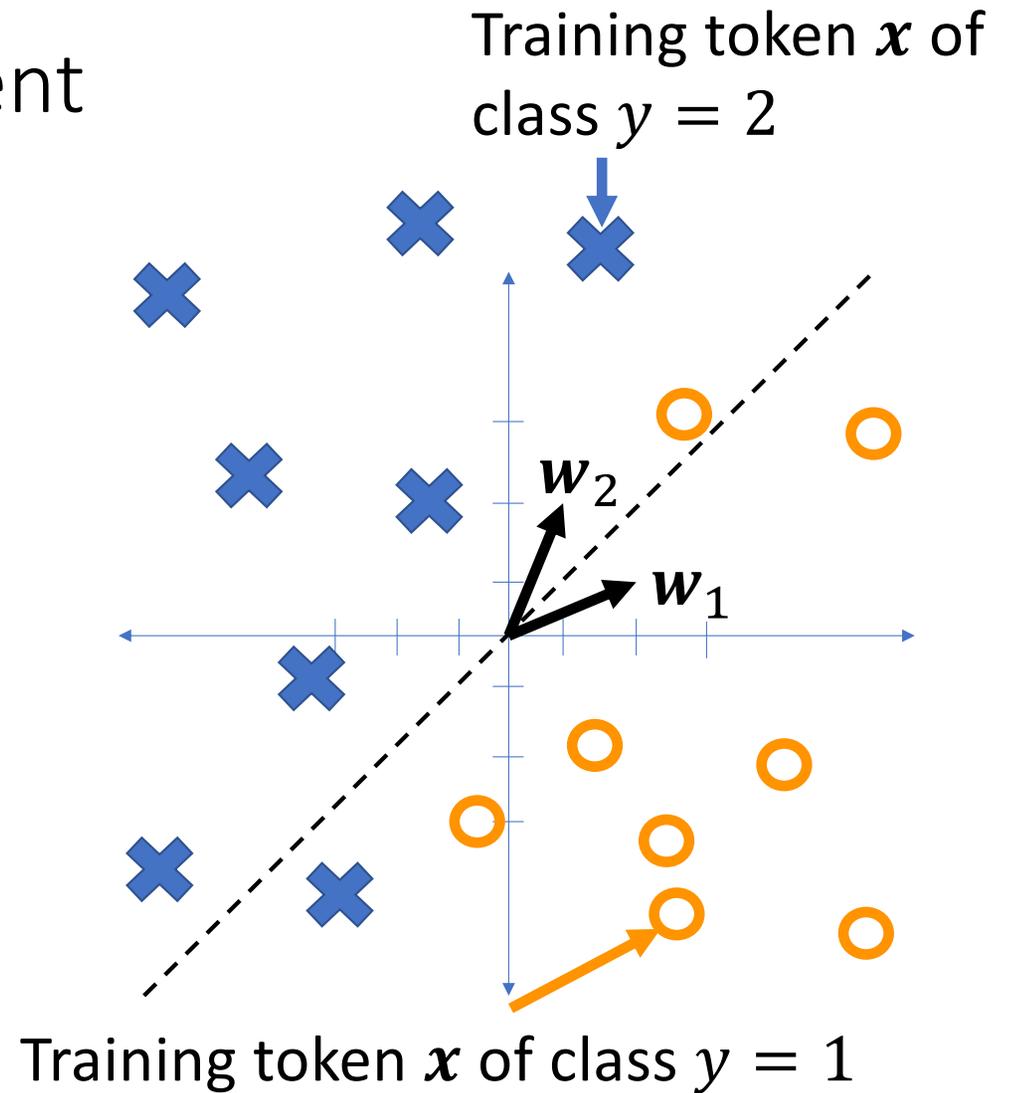
Suppose we have a softmax output, so we want $f_c(\mathbf{x}) = \Pr(Y = c|\mathbf{x})$. We can train this by learning \mathbf{W} to maximize the probability of the training corpus. If we assume all training tokens are independent, we get:

$$\begin{aligned}\mathbf{W} &= \operatorname{argmax}_{\mathbf{W}} \prod_{i=1}^n \Pr(Y = y_i|\mathbf{x}_i) \\ &= \operatorname{argmax}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \ln \Pr(Y = y_i|\mathbf{x}_i) \\ &= \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i, \quad \mathcal{L}_i = -\ln f_{y_i}(\mathbf{x}_i)\end{aligned}$$

Stochastic gradient descent

Suppose we have a training example (\mathbf{x}, y) . We want to find

$$\mathbf{w}_c \leftarrow \mathbf{w}_c - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}_c}$$



Stochastic gradient descent

Suppose we have a training example (\mathbf{x}, y) with the loss

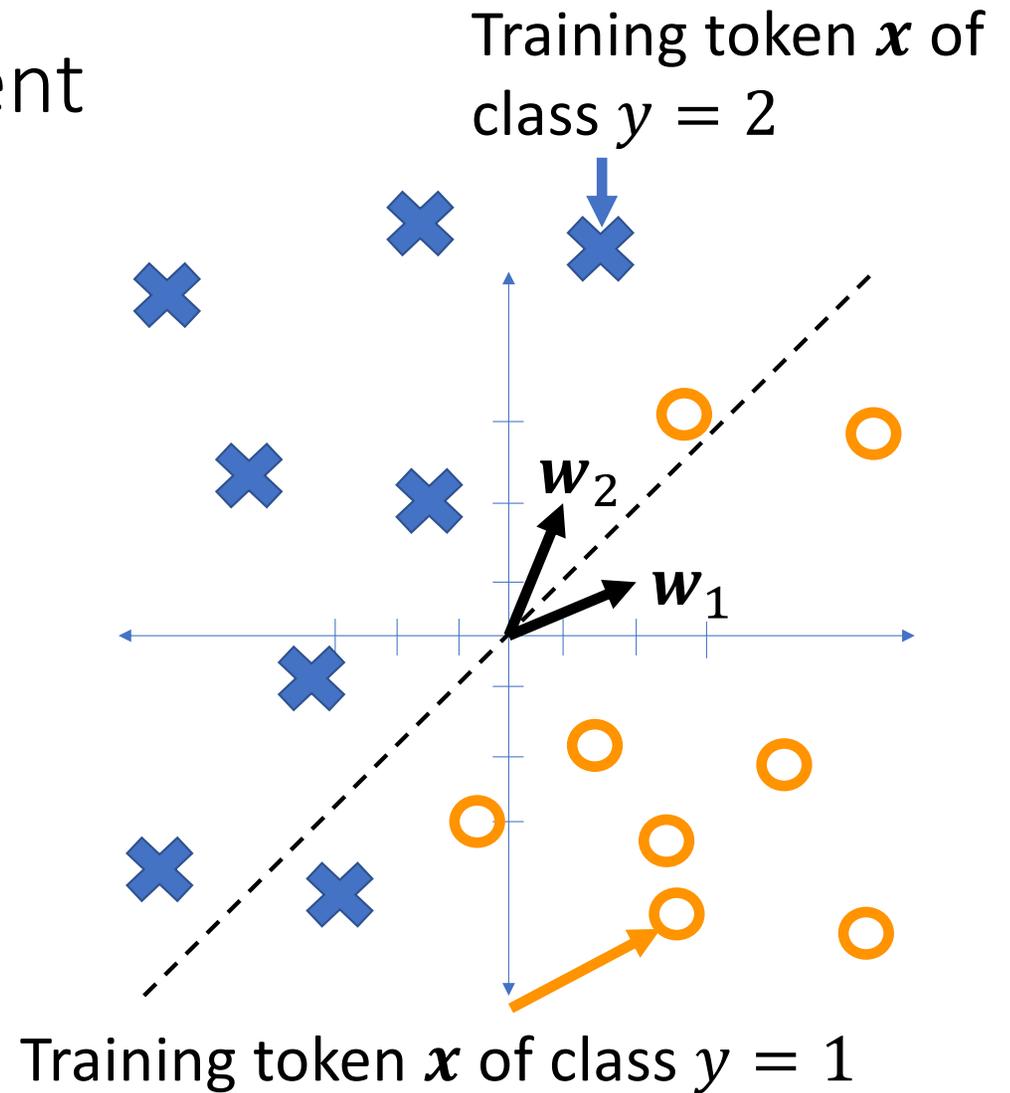
$$\mathcal{L} = -\ln f_y(\mathbf{x}) = -\ln \left(\frac{e^{\mathbf{w}_y^T \mathbf{x}}}{\sum_{k=1}^v e^{\mathbf{w}_k^T \mathbf{x}}} \right) = \ln \left(\sum_{k=1}^v e^{\mathbf{w}_k^T \mathbf{x}} \right) - \mathbf{w}_y^T \mathbf{x}$$

Its gradient with respect to \mathbf{w}_c is

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_c} &= \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{k=1}^v e^{\mathbf{w}_k^T \mathbf{x}}} \mathbf{x} - \mathbb{1}_{y=c} \mathbf{x} \\ &= (f_c(\mathbf{x}) - \mathbb{1}_{y=c}) \mathbf{x} \end{aligned}$$

Stochastic gradient descent

- Shift \mathbf{w}_y in the direction of \mathbf{x} :
$$\mathbf{w}_y \leftarrow \mathbf{w}_y + \eta (1 - f_y(\mathbf{x})) \mathbf{x}$$
- For all other classes, $c \neq y$, shift them in the direction opposite \mathbf{x} :
$$\mathbf{w}_c \leftarrow \mathbf{w}_c + \eta (0 - f_c(\mathbf{x})) \mathbf{x}$$



Summary

- Sigmoid review

$$f(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w}), \quad \mathbf{w} \leftarrow \mathbf{w} + \eta \sum_{i=1}^n (y_i - f(\mathbf{x})) \mathbf{x}_i$$

- Multi-class linear classifiers

$$f(\mathbf{x}) = \operatorname{argmax} \mathbf{W} \mathbf{x}$$

- One-hot vectors

$$f_c(\mathbf{x}) = \mathbb{1}_{\operatorname{argmax} \mathbf{W} \mathbf{x} = c}$$

- Softmax nonlinearity

$$f_c(\mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{k=1}^v \exp(\mathbf{w}_k^T \mathbf{x})}$$

- Derivative of the log softmax

$$\mathbf{w}_c \leftarrow \mathbf{w}_c + \eta \sum_{i=1}^n (\mathbb{1}_{y_i=c} - f_c(\mathbf{x})) \mathbf{x}_i$$