

# **18-661 Introduction to Machine Learning**

## Multi-class Classification

---

Spring 2023

ECE – Carnegie Mellon University

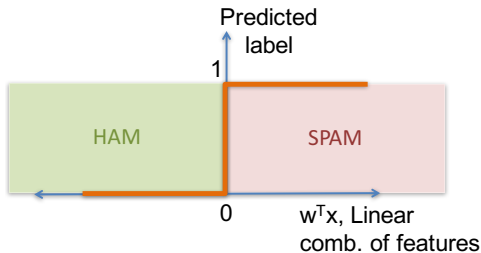
1. Review of Logistic Regression
2. Non-linear Decision Boundaries
3. Multi-class Classification
  - Multi-class Naïve Bayes
  - Multi-class Logistic Regression
4. Evaluating Classification Methods

# Review of Logistic Regression

---

# Visualizing a Linear Classifier

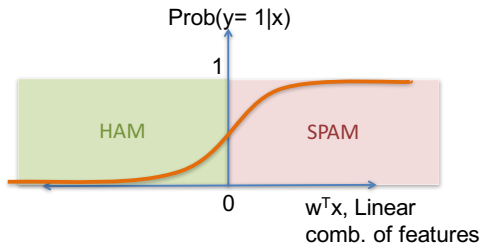
- $x_1$  = # of times 'lottery' appears in an email
- $x_2$  = # of times 'meet' appears in an email
- Define feature vector  $\mathbf{x} = [1, x_1, x_2]$
- Learn the decision boundary  $w_0 + w_1x_1 + w_2x_2 = 0$  such that
  - If  $\mathbf{w}^\top \mathbf{x} \geq 0$  declare  $y = 1$  (spam)
  - If  $\mathbf{w}^\top \mathbf{x} < 0$  declare  $y = 0$  (ham)



$y = 1$  for spam,  $y = 0$  for ham

# Intuition: Logistic Regression

- Suppose we want to output the **probability** of an email being spam/ham instead of just 0 or 1
- This gives information about the confidence in the decision
- Use a function  $\sigma(\mathbf{w}^\top \mathbf{x})$  that maps  $\mathbf{w}^\top \mathbf{x}$  to a value between 0 and 1



Probability that predicted label is 1 (spam)

**Key Problem:** Finding optimal weights  $\mathbf{w}$  that accurately predict this probability for a new email

# Formal Setup: Binary Logistic Classification

- Input/features:  $\mathbf{x} = [1, x_1, x_2, \dots, x_D] \in \mathbb{R}^{D+1}$
- Output:  $y \in \{0, 1\}$
- Training data:  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$
- Model:

$$p(y = 1|\mathbf{x}; \mathbf{w}) = \sigma[g(\mathbf{x})]$$

where

$$g(\mathbf{x}) = w_0 + \sum_d w_d x_d = \mathbf{w}^\top \mathbf{x}$$

and  $\sigma[\cdot]$  stands for the **sigmoid** function

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

# How to Optimize $w$ ?

- Probability of a single training sample  $(\mathbf{x}_n, y_n)$

$$P(y_n|\mathbf{x}_n; \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^\top \mathbf{x}_n) & \text{if } y_n = 1 \\ 1 - \sigma(\mathbf{w}^\top \mathbf{x}_n) & \text{otherwise} \end{cases}$$

- Compact expression, exploiting that  $y_n$  is either 1 or 0

$$P(y_n|\mathbf{x}_n; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_n)^{y_n} [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]^{1-y_n}$$

- Minimize the negative log-likelihood of the whole training data  $\mathcal{D}$ ,  
i.e., **the cross-entropy error function**

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

# Cross-Entropy as a Loss Function

## Supervised learning

We aim to build a function  $h(\mathbf{x})$  to predict the true value  $y$  associated with  $\mathbf{x}$ . If we make a mistake, we incur a **loss**

$$\ell(h(\mathbf{x}), y)$$

Cross-entropy is also a sum over all data samples:

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

where  $h(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x})$ .

What is the loss function?

$$\ell(h(\mathbf{x}_n), y) = -\{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$



# Gradient Descent for Logistic Regression

- We want to minimize the cross-entropy error function:

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

- Simple fact: derivatives of  $\sigma(a)$  have a nice form:

$$\frac{d}{da} \sigma(a) = \sigma(a)[1 - \sigma(a)]$$

- Gradient of cross-entropy loss is then

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_n \underbrace{\{\sigma(\mathbf{w}^\top \mathbf{x}_n) - y_n\}}_{:=e_n} \mathbf{x}_n$$

- $e_n = \{\sigma(\mathbf{w}^\top \mathbf{x}_n) - y_n\}$  is called the **error** for the  $n$ th training sample.

## Gradient descent for logistic regression

- Choose a proper step size  $\eta > 0$
- Iteratively update the parameters following the negative gradient to minimize the error function

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \sum_n \left\{ \sigma(\mathbf{w}^{(t)\top} \mathbf{x}_n) - y_n \right\} \mathbf{x}_n$$

## Stochastic gradient descent for logistic regression

- Choose a proper step size  $\eta > 0$
- Draw a sample  $n$  uniformly at random
- Iteratively update the parameters following the negative gradient to minimize the error function

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \left\{ \sigma(\mathbf{w}^{(t)\top} \mathbf{x}_n) - y_n \right\} \mathbf{x}_n$$

# Naïve Bayes vs. Logistic Regression

Both classification models are linear functions of features

## **Joint vs. conditional distribution**

Naive Bayes models the **joint** distribution:  $P(X, Y) = P(Y)P(X|Y)$

Logistic regression models the **conditional** distribution:  $P(Y|X)$

## **Correlated vs. independent features**

Naive Bayes assumes independence of features and multiple occurrences

Logistic Regression implicitly captures correlations when training weights

## **Generative vs. Discriminative**

NB is a **generative** model, LR is a **discriminative** model

# Logistic Regression vs. Linear Regression

	Logistic regression	Linear regression
Training data	$(\mathbf{x}_n, y_n), y_n \in \{0, 1\}$	$(\mathbf{x}_n, y_n), y_n \in \mathbb{R}$
Loss function	cross-entropy	RSS
Interpretation of $y_n   \mathbf{x}_n, \mathbf{w}$	$\sim \text{Ber}(\sigma(\mathbf{w}^\top \mathbf{x}_n))$	$\sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$
Gradient per sample	$(\sigma(\mathbf{x}_n^\top \mathbf{w}) - y_n) \mathbf{x}_n$	$(\mathbf{x}_n^\top \mathbf{w} - y_n) \mathbf{x}_n$

**Cross-entropy loss function (logistic regression):**

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

**RSS loss function (linear regression):**

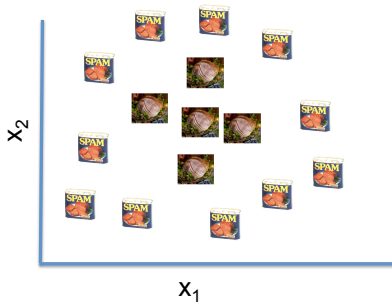
$$RSS(\mathbf{w}) = \frac{1}{2} \sum_n (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

1. Review of Logistic Regression
2. Non-linear Decision Boundaries
3. Multi-class Classification
  - Multi-class Naïve Bayes
  - Multi-class Logistic Regression
4. Evaluating Classification Methods

# Non-linear Decision Boundaries

---

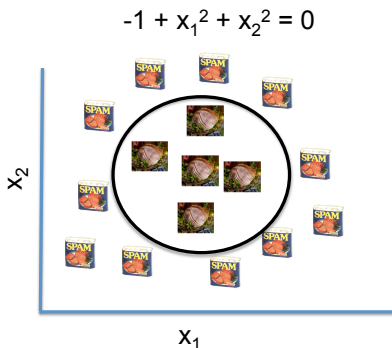
# How to Handle More Complex Decision Boundaries?



- This data is not linearly separable...
- Use **non-linear basis functions** to add more features.

# Adding Polynomial Features

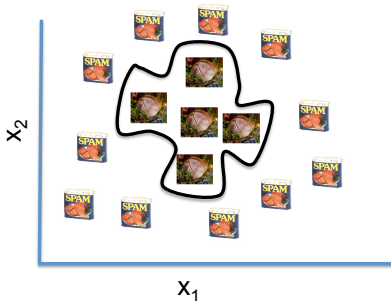
- New feature vector is  $\mathbf{x} = [1, x_1, x_2, x_1^2, x_2^2]$
- $\Pr(y = 1|\mathbf{x}) = \sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2)$
- If  $\mathbf{w} = [-1, 0, 0, 1, 1]$ , the boundary is  $-1 + x_1^2 + x_2^2 = 0$ 
  - If  $-1 + x_1^2 + x_2^2 \geq 0$  declare spam
  - If  $-1 + x_1^2 + x_2^2 < 0$  declare ham





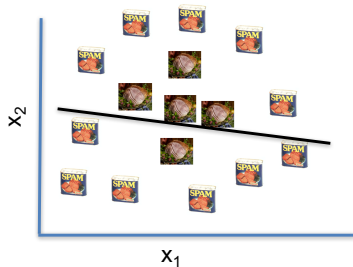
# Adding Polynomial Features

- What if we add many more features and define  $\mathbf{x} = [1, x_1, x_2, x_1^2, x_2^2, x_1^3, x_2^3, \dots]$ ?
- We get a complex decision boundary

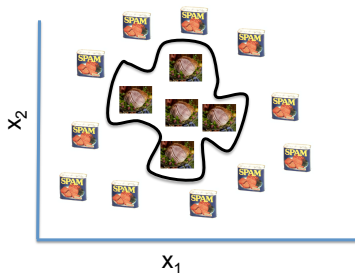


Can result in overfitting and bad generalization to new data points.

# Concept-check: Bias-Variance Trade-off



high bias



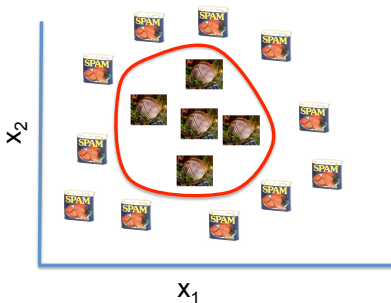
high variance

# Solution to Overfitting: Regularization

- Add regularization term to be cross entropy loss function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1-y_n) \log [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\} + \underbrace{\frac{1}{2} \lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

- Perform gradient descent on this regularized function
- Often, we do **NOT** regularize the bias term  $w_0$



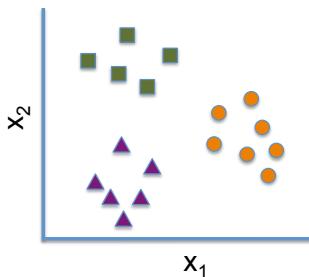
1. Review of Logistic Regression
2. Non-linear Decision Boundaries
3. Multi-class Classification
  - Multi-class Naïve Bayes
  - Multi-class Logistic Regression
4. Evaluating Classification Methods

# Multi-class Classification

---

# What If There Are More than 2 Classes?

- Dog vs. cat. vs crocodile
- Movie genres (action, horror, comedy, ...)
- Part of speech tagging (verb, noun, adjective, ...)
- ...



**Predict multiple classes/outcomes  $C_1, C_2, \dots, C_M$ :**

- Weather prediction: sunny, cloudy, raining, etc
- Optical character recognition: 10 digits + 26 characters (lower and upper cases) + special characters, etc.

$M$  = number of classes

**Methods we've studied for binary classification:**

- Naïve Bayes
- Logistic regression

Do they generalize to multi-class classification?

# Naïve Bayes Is Already Multi-class!

## Formal Definition

Given a random vector  $\mathbf{X} \in \mathbb{R}^K$  and a dependent variable  $Y \in [C]$ , the Naïve Bayes model defines the joint distribution

$$P(\mathbf{X} = \mathbf{x}, Y = c) = P(Y = c)P(\mathbf{X} = \mathbf{x} | Y = c) \quad (1)$$

$$= P(Y = c) \prod_{k=1}^K P(\text{word}_k | Y = c)^{x_k} \quad (2)$$

$$= \pi_c \prod_{k=1}^K \theta_{ck}^{x_k} \quad (3)$$

where  $x_k$  is the number of occurrences of the  $k$ th word,  $\pi_c$  is the prior probability of class  $c$  (which allows multiple classes!), and  $\theta_{ck}$  is the weight of the  $k$ th word for the  $c$ th class.



# Learning Multi-class Naïve Bayes

## Training data

$$\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N \rightarrow \mathcal{D} = \{(\{x_{nk}\}_{k=1}^K, y_n)\}_{n=1}^N$$

## Our goal

Learn  $\pi_c, c = 1, 2, \dots, C$ , and  $\theta_{ck}, \forall c \in [C], k \in [K]$  under the constraints:

$$\sum_c \pi_c = 1$$

and

$$\sum_k \theta_{ck} = \sum_k P(\text{word}_k | Y = c) = 1$$

as well as  $\pi_c, \theta_{ck} \geq 0$ .

# Our Hammer: Maximum Likelihood Estimation

- Find the log-likelihood of the training data

$$\begin{aligned}\mathcal{L} &= \log P(\mathcal{D}) = \log \prod_{n=1}^N \pi_{y_n} P(\mathbf{x}_n | y_n) \\ &= \log \prod_{n=1}^N \left( \pi_{y_n} \prod_k \theta_{y_n k}^{x_{nk}} \right) \\ &= \sum_n \left( \log \pi_{y_n} + \sum_k x_{nk} \log \theta_{y_n k} \right) \\ &= \sum_n \log \pi_{y_n} + \sum_{n,k} x_{nk} \log \theta_{y_n k}\end{aligned}$$

- Optimize it!

$$(\pi_c^*, \theta_{ck}^*) = \arg \max \sum_n \log \pi_{y_n} + \sum_{n,k} x_{nk} \log \theta_{y_n k}$$

# Our Hammer: Maximum Likelihood Estimation

## Optimization Problem

$$(\pi_c^*, \theta_{ck}^*) = \arg \max \left( \sum_n \log \pi_{y_n} + \sum_{n,k} x_{nk} \log \theta_{y_n k} \right)$$

## Solution

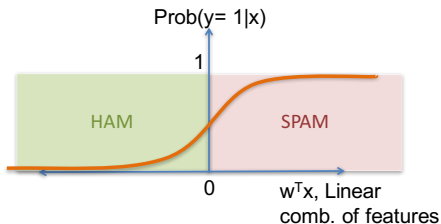
$$\theta_{ck}^* = \frac{\text{\#of times word } k \text{ shows up in data points labeled as } c}{\text{\#total trials for data points labeled as } c}$$

$$\pi_c^* = \frac{\text{\#of data points labeled as } c}{N}$$

1. Review of Logistic Regression
2. Non-linear Decision Boundaries
3. Multi-class Classification
  - Multi-class Naïve Bayes
  - Multi-class Logistic Regression
4. Evaluating Classification Methods

# Logistic Regression for Predicting Multiple Classes?

- The linear decision boundary that we optimized was specific to binary classification.
  - If  $\sigma(\mathbf{w}^\top \mathbf{x}) \geq 0.5$  declare  $y = 1$  (spam)
  - If  $\sigma(\mathbf{w}^\top \mathbf{x}) < 0.5$  declare  $y = 0$  (ham)
- How to extend it to multi-class classification?

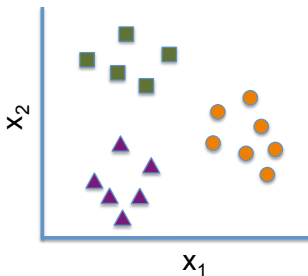


$y = 1$  for spam,  $y = 0$  for ham

Idea: Express as multiple binary classification problems

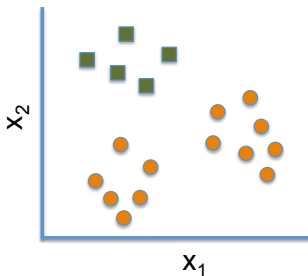
# The One-versus-Rest or One-versus-All Approach

- For each class  $c$ , change the problem into binary classification
  1. Relabel training data with label  $c$ , into POSITIVE (or '1').
  2. Relabel all the rest data into NEGATIVE (or '0').
- Repeat this multiple times: Train  $C$  binary classifiers, using logistic regression to differentiate the two classes each time.



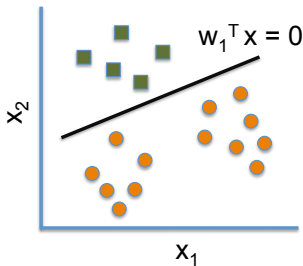
# The One-versus-Rest or One-versus-All Approach

- For each class  $c$ , change the problem into binary classification
  1. Relabel training data with label  $c$ , into POSITIVE (or '1')
  2. Relabel all the rest data into NEGATIVE (or '0')
- Repeat this multiple times: Train  $C$  binary classifiers, using logistic regression to differentiate the two classes each time.



# The One-versus-Rest or One-versus-All Approach

- For each class  $c$ , change the problem into binary classification
  1. Relabel training data with label  $c$ , into POSITIVE (or '1')
  2. Relabel all the rest data into NEGATIVE (or '0')
- Repeat this multiple times: Train  $C$  binary classifiers, using logistic regression to differentiate the two classes each time.

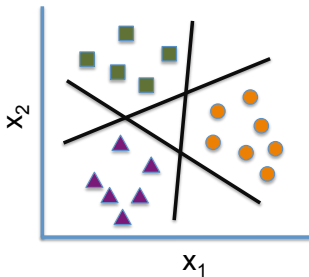




# The One-versus-Rest or One-versus-All Approach

How to combine these linear decision boundaries?

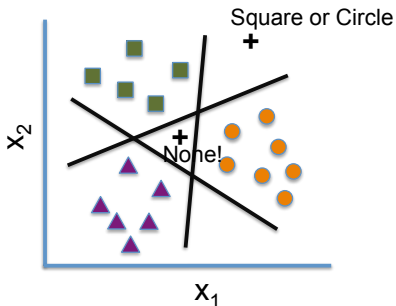
- There is ambiguity in some of the regions (the 4 triangular areas).



# The One-versus-Rest or One-versus-All Approach

How to combine these linear decision boundaries?

- There is ambiguity in some of the regions (the 4 triangular areas).
- How do we resolve this?

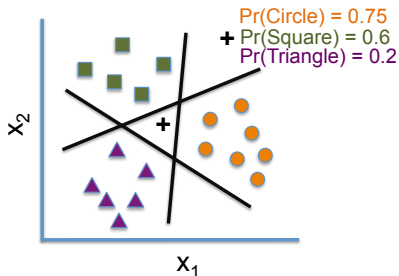


# The One-versus-Rest or One-versus-All Approach

How to combine these linear decision boundaries?

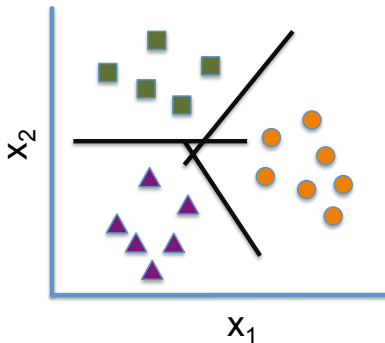
- Use the **confidence estimates**  $\Pr(y = 1|\mathbf{x}) = \sigma(\mathbf{w}_1^\top \mathbf{x})$ ,  
...  $\Pr(y = C|\mathbf{x}) = \sigma(\mathbf{w}_C^\top \mathbf{x})$
- Declare class  $c^*$  that maximizes

$$c^* = \arg \max_{c=1,\dots,C} \Pr(y = c|\mathbf{x}) = \sigma(\mathbf{w}_c^\top \mathbf{x})$$



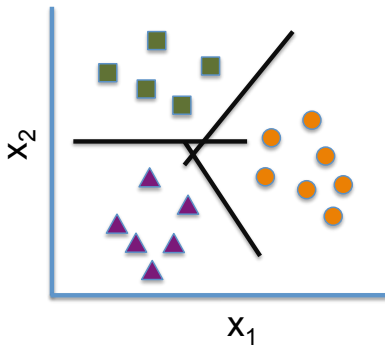
# The One-versus-One Approach

- For each **pair** of classes  $c$  and  $c'$ , change the problem into binary classification.
  1. Relabel training data with label  $c$ , into POSITIVE (or '1')
  2. Relabel training data with label  $c'$  into NEGATIVE (or '0')
  3. **Disregard** all other data



# The One-versus-One Approach

- How many binary classifiers for  $C$  classes?  $C(C - 1)/2$
- How to combine their outputs?
- Given  $\mathbf{x}$ , count the  $C(C - 1)/2$  votes from outputs of all binary classifiers and declare the winner as the predicted class.
- Use confidence scores to resolve ties.



# Contrast These Approaches

## Number of binary classifiers to be trained

- **One-versus-All:**  $C$  classifiers.
- **One-versus-One:**  $C(C - 1)/2$  classifiers – bad if  $C$  is large

## Effect of relabeling and splitting training data

- **One-versus-All:** imbalance in the number of positive and negative samples can cause bias in each trained classifier.
- **One-versus-One:** each classifier trained on a small subset of data (only data in two classes), which can result in high variance.

## Any other ideas?

- **Hierarchical classification** – we will see this in decision trees
- **Multinomial logistic regression** – directly output probabilities of  $y$  being in each of the  $C$  classes.

So, can we define the following conditional model?

$$P(y = c|\mathbf{x}) = \sigma[\mathbf{w}_c^\top \mathbf{x}].$$

This would **not** work because:

$$\sum_c P(y = c|\mathbf{x}) = \sum_c \sigma[\mathbf{w}_c^\top \mathbf{x}] \neq 1,$$

so each summand can be any number (independently) between 0 and 1.

**But we are close!**

Learn the  $C$  linear models jointly to ensure this property holds!

# Multinomial Logistic Regression

- **Model:** For each class  $c$ , we have a parameter vector  $\mathbf{w}_c$  and model the posterior probability as:

$$P(c|\mathbf{x}) = \frac{e^{\mathbf{w}_c^\top \mathbf{x}}}{\sum_{c'} e^{\mathbf{w}_{c'}^\top \mathbf{x}}} \quad \leftarrow \quad \textit{This is called the softmax function.}$$

- **Decision boundary:** Assign  $\mathbf{x}$  with the label that is the maximum of posterior:

$$\arg \max_c P(c|\mathbf{x}) \rightarrow \arg \max_c \mathbf{w}_c^\top \mathbf{x}.$$



# How Does the Softmax Function Behave?

Suppose we have

$$\mathbf{w}_1^\top \mathbf{x} = 100, \quad \mathbf{w}_2^\top \mathbf{x} = 50, \quad \mathbf{w}_3^\top \mathbf{x} = -20.$$

We would pick the **winning** class label 1.

**Softmax translates these scores into well-formed conditional probabilities**

$$P(y = 1|\mathbf{x}) = \frac{e^{100}}{e^{100} + e^{50} + e^{-20}} < 1$$

- Preserves relative ordering of scores.
- Maps scores to values between 0 and 1 that also sum to 1.

Multinomial model reduces to binary logistic regression when  $C = 2$ .

$$\begin{aligned} P(1|\mathbf{x}) &= \frac{e^{\mathbf{w}_1^\top \mathbf{x}}}{e^{\mathbf{w}_1^\top \mathbf{x}} + e^{\mathbf{w}_2^\top \mathbf{x}}} = \frac{1}{1 + e^{-(\mathbf{w}_1 - \mathbf{w}_2)^\top \mathbf{x}}} \\ &= \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \end{aligned}$$

when we define  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$ . Multinomial logistic regression thus generalizes the (binary) logistic regression to deal with multiple classes.

# Parameter Estimation for Multinomial Logistic Regression

**Discriminative approach:** Maximize conditional likelihood

$$\log P(\mathcal{D}) = \sum_n \log P(y_n | \mathbf{x}_n)$$

We will change  $y_n$  to  $\mathbf{y}_n = [y_{n1} \ y_{n2} \ \cdots \ y_{nC}]^\top$ , a  $C$ -dimensional vector using 1-of- $C$  encoding.

$$y_{nc} = \begin{cases} 1 & \text{if } y_n = c \\ 0 & \text{otherwise} \end{cases}$$

Ex: if  $y_n = 2$ , then,  $\mathbf{y}_n = [0 \ \mathbf{1} \ 0 \ 0 \ \cdots \ 0]^\top$ .

$$\Rightarrow \sum_n \log P(y_n | \mathbf{x}_n) = \sum_n \log \prod_{c=1}^C P(c | \mathbf{x}_n)^{y_{nc}} = \sum_n \sum_c y_{nc} \log P(c | \mathbf{x}_n)$$

# Cross-entropy Error Function

**Definition:** negative log-likelihood

$$\begin{aligned}\mathcal{E}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) &= - \sum_n \sum_c y_{nc} \log P(c|\mathbf{x}_n) \\ &= - \sum_n \sum_c y_{nc} \log \left( \frac{e^{\mathbf{w}_c^\top \mathbf{x}_n}}{\sum_{c'} e^{\mathbf{w}_{c'}^\top \mathbf{x}_n}} \right)\end{aligned}$$

## Properties of cross-entropy

- Convex in the  $\mathbf{w}$  vectors, therefore local minimum = global optimum
- Optimization requires numerical procedures, analogous to those used for binary logistic regression.

# Finding the Gradient

$$\begin{aligned}\mathcal{E}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) &= - \sum_n \sum_c y_{nc} \log P(c|\mathbf{x}_n) \\ &= - \sum_n \sum_c y_{nc} \log \left( \frac{e^{\mathbf{w}_c^\top \mathbf{x}_n}}{\sum_{c'} e^{\mathbf{w}_{c'}^\top \mathbf{x}_n}} \right)\end{aligned}$$

- Need to find the gradient w.r.t.  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C$  and update

$$\mathbf{w}_c \leftarrow \mathbf{w}_c - \eta \frac{\partial \mathcal{E}}{\partial \mathbf{w}_c}, \quad c = 1, \dots, C$$

Can you find the gradient? (Hint: what is the gradient of the softmax function?)

# Evaluating Classification Methods

---

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

- Easy to optimize!
- Average loss over the (training, validation, test) dataset
- ...but what does it mean?

# Interpretable Classification Metrics

True positive	False positive
False negative	True negative

- Measure the accuracy within each class
- Accounts for imbalance between classes

These metrics are **difficult to optimize directly**, but they have the advantage of being easily interpretable.

- **Sensitivity**: true positive rate

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- **Specificity**: true negative rate

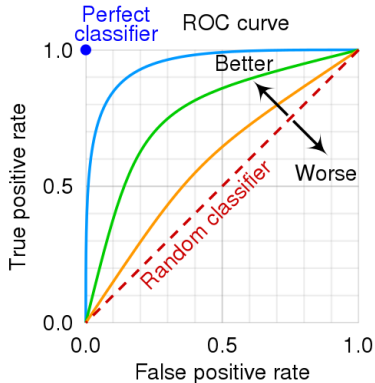
$$\text{TNR} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

- **Precision**: positive predictive value

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$



# Combining These Metrics: the ROC Curve



Receiver Operating Characteristic  
(ROC)

- Define a “threshold” for the positive/negative split
- Increasing the threshold: more samples are predicted to be positive
- **Area Under the ROC Curve:** want this as large as possible

# You Should Know

- How to generalize logistic regression to handle nonlinear decision boundaries.
- How to handle multiclass classification: one-versus-all, one-versus-one, multinomial regression.
- How to measure classification accuracy