

# Machine Learning Crash Course

## Introduction to Classification

Dr. Henrich Matzinger and Dr. Greg Mayer  
Machine Learning and AI Crash Course  
December 2019

School of Mathematics  
Georgia Institute of Technology, Atlanta, GA

## **Objective**

Introduce a classification algorithm that we will use for text classification later in our crash course

## **Topics**

1. classification rules
2. training data
3. best classification rule on training data

# Catching Fish

Suppose:

- A boat on the sea catches only salmon and tuna
- A robot on the boat classifies each fish by its size: small, medium, or large
- The robot cannot tell whether fish are salmon or tuna

Given a size, how can we best predict the type of fish?

# A Probability Model

## Notation

- sizes: 1 = small, 2 = medium, 3 = large
- $X$  denotes size of fish
- $Y$  denotes the **class**, or type of fish

## Assumptions

- there is an underlying **probability model**, or “probability distribution” or in statistical parlance “the population distribution”
- we are given a **joint probability table**

	$X = 1$	$X=2$	$X = 3$
tuna	$P(Y = \text{tuna}, X = 1)$	$P(Y = \text{tuna}, X = 2)$	$P(Y = \text{tuna}, X = 3)$
salmon	$P(Y = \text{salmon}, X = 1)$	$P(Y = \text{salmon}, X = 2)$	$P(Y = \text{salmon}, X = 3)$

# Joint Probability Table Example

- Suppose the underlying probability model is:

	1	2	3
tuna	0.1	0.2	0.3
salmon	0.2	0.1	0.1

(1)

For example,  $P(Y = \text{salmon}, X = 1) = 20\%$ , meaning that in the waters in which we are fishing, 20% of the fish are salmon of size 1.

- Similarly,  $P(Y = \text{tuna}, X = 3) = 0.3$ , meaning that 30% of the fish in the water are *tuna* of size 3

What if we wanted to predict the type of fish given its size?

# A Decision Rule

- Our **decision rule** can be based on choosing the class which has highest probability given its size
- Suppose our robot classifies a fish as being size 2. This fish is then twice as likely to be a tuna:

$$P(Y = \textit{tuna} | X = 2) = \frac{0.2}{0.3} = \frac{2}{3}$$

and

$$P(Y = \textit{salmon} | X = 2) = \frac{0.1}{0.3} = \frac{1}{3}$$

So, if we catch a fish of size 2, our decision rule states that we classify that fish as tuna.

# Bayse Classifier

- In fancy machine learning parlance, classification rules are **classifiers**
- The set  $\{1, 2, 3\}$  would be the **feature space**
- *tuna* and *salmon* are the **classes**
- The best rule is denoted by  $g^*$  and is called a *Bayes classifier*

Formally the Bayes classifier can be defined by:

$$g^*(x) = \textit{tuna} \text{ if and only if } \frac{P(Y = \textit{tuna} | X = x)}{P(Y = \textit{salmon} | X = x)} \geq 1 \quad (2)$$

# Prior Probability

Let

- $\pi_{\text{tuna}}$  = probability that a given fish from the water is a tuna
- $\pi_{\text{salmon}}$  = probability that a given fish from the water is a salmon

Or,

$$\begin{aligned}\pi_{\text{tuna}} &:= P(Y = \text{tuna}) \\ \pi_{\text{salmon}} &:= P(Y = \text{salmon})\end{aligned}$$

Another way to rewrite the equations leading to the Bayse classifier, is obtained by using Bayse theorem and is:

$$g^*(x) = \text{tuna} \text{ if and only if } \frac{\pi_{\text{tuna}} \cdot P(X = x|Y = \text{tuna})}{\pi_{\text{salmon}} \cdot P(X = x|Y = \text{salmon})} \geq 1 \quad (3)$$



# Best Decision Rule

In the present case, the best classifier is given by

$$g^*(1) = \textit{salmon}, g^*(2) = \textit{tuna}, g^*(3) = \textit{tuna}$$

Our optimal decision rule is represented in our table as bold entries:

	1	2	3
tuna	0.1	<b>0.2</b>	<b>0.3</b>
salmon	<b>0.2</b>	0.1	0.1

(4)

# Probability of Missclassification

Our optimal decision being represented by the green entries:

	1	2	3
tuna	0.1	0.2	0.3
salmon	0.2	0.1	0.1

(5)

The **misclassification probability** is the sum of the red entries:

$$\begin{aligned}\text{misclassification probability of } g^* &= P(g^*(X) \neq Y) \\ &= 0.1 + 0.1 + 0.1 \\ &= 0.3\end{aligned}$$

This means that in the long run, our robot misclassifies 30% of the fish.

# Training Data

# Underlying Probability Model Usually Unknown

- on general, underlying probabilities are not exactly known
- but if we can catch some fish, label them manually as salmon or tuna and then estimate the probabilities given in Table 5.
- the data we use to determine a classification rule is called a **training sample**

# Estimating Probabilities

Suppose we catch 100 fish and obtain the frequency table

	1	2	3
tuna	4	10	45
salmon	15	6	20

(6)

Four fish are tuna of size 1, so we make the estimate:

$$\hat{P}(Y = \text{tuna}, X = 1) = 0.04$$

Fifteen salmon of size 1, which leads to our estimate:

$$\hat{P}(Y = \text{salmon}, X = 1) = 0.15$$

So, we have the estimated probabilities given in the table

	1	2	3
tuna	0.04	0.10	0.45
salmon	0.15	0.06	0.2

(7)

# Best classification rule on given data

The classification rule (classifier) which we chose is

$$g(1) = \text{salmon}, g(2) = \text{tuna}, g(3) = \text{tuna}$$

is based on the estimated probabilities given in the table

	1	2	3
tuna	0.04	0.10	0.45
salmon	0.15	0.06	0.2

(8)

**This is the rule that is best at classifying our data**

Essentially:

- we do not know true probabilities but we have their estimates.
- we act as if the estimates were the true probabilities and create a best prediction rule from them

# Best Classification Rules

# Overfitting

- When we have many more options for the feature and not enough data we encounter **overfitting**
- Suppose our robot can identify many more sizes and we collect the frequency table below

	1	1.2	1.4	1.6	1.8	2	2.2	2.4	2.6	2.8	3
tuna	0	<b>1</b>	0	<b>1</b>	<b>3</b>	<b>1</b>	0	0	<b>1</b>	0	0
salmon	0	0	<b>1</b>	0	0	0	<b>1</b>	<b>2</b>	0	<b>3</b>	<b>2</b>

- we can create a rule that classifies fish **in our training data** correctly 100 % of the time
- if we then catch another fish and use such a rule, what problems might we encounter?



# Formal Random Variables

Again, same frequency table:

	1	1.2	1.4	1.6	1.8	2	2.2	2.4	2.6	2.8	3
tuna	0	<b>1</b>	0	<b>1</b>	<b>3</b>	<b>1</b>	0	0	<b>1</b>	0	0
salmon	0	0	<b>1</b>	0	0	0	<b>1</b>	<b>2</b>	0	<b>3</b>	<b>2</b>

Consider the classification rule:

$$\text{size} \geq 2 \implies \text{salmon}$$

Based on our frequency table, this classifies correctly 13 out of 16

- Let  $X_i$  be the size of the  $i$ -th fish and let  $Y_i$  be its class.
- Our training sample is the data  $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$
- Let  $Z_i = 1$  if  $i$ th fish gets classified correctly with our rule and  $Z_i = 0$  otherwise. In our case

$$Z_1 + Z_2 + \dots + Z_{16} = 13$$

# Probability of Correct Classification

- Then, we can estimate the probability of correct classification:

$$\hat{P}(g(X) = Y) = \frac{Z_1 + Z_2 + \dots + Z_n}{n}$$

- This is estimating the parameter  $p$  in a binomial, which if  $n$  is not very small is quite precise.
- So we should be able to see which classification rule are best by using our estimate  $\hat{P}$

# Solution to Overfitting

- to avoid overfitting we can use a small number of rules to identify **best classification rule**
- In our data, if we only consider classification rules of the type

$$X \geq \text{constant} \implies \text{salmon}$$

what value for constant leads to least classification error in the training data?

	1	1.2	1.4	1.6	1.8	2	2.2	2.4	2.6	2.8	3
tuna	0	<b>1</b>	0	<b>1</b>	<b>3</b>	<b>1</b>	0	0	<b>1</b>	0	0
salmon	0	0	<b>1</b>	0	0	0	<b>1</b>	<b>2</b>	0	<b>3</b>	<b>2</b>

# Justification for Best Classification Rule

Assume that

1. size of tuna is normal with expectation  $\mu_t$
2. size of salmon normal with expectation  $\mu_s$
3. distributions both have standard deviation 1

Then, the optimal rule:

$$P(x | \text{tuna})\pi_t > P(x | \text{salmon})\pi_s \implies \text{tuna}$$

is equivalent to

$$\pi_t \cdot \exp(-(x - \mu_t)^2) > \pi_s \cdot \exp(-(x - \mu_s)^2)$$

solving for  $x$

$$x \geq \frac{\mu_s^2 - \mu_t^2}{\mu_t - \mu_s} + \log(\pi_s) - \log(\pi_t) = \text{constant}$$

Thus  $x > \text{constant}$  is a reasonable approach to classifying our data