

Classification: Probabilistic Generative Model

Classification



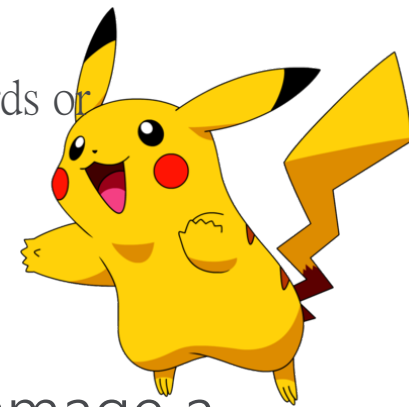
- Credit Scoring
 - Input: income, savings, profession, age, past financial history
 - Output: accept or refuse
- Medical Diagnosis
 - Input: current symptoms, age, gender, past medical history
 - Output: which kind of diseases
- Handwritten character recognition
- Face recognition
 - Input: image of a face, output: person

Input:  output: 金

Example Application Recognize the type of Pokémon.



$$f(\text{Pikachu}) = \text{Electric} \quad f(\text{Squirtle}) = \text{Water} \quad f(\text{Bulbasaur}) = \text{Grass}$$



Example Application

- **HP:** hit points, or health, defines how much damage a pokemon can withstand before fainting **35**
- **Attack:** the base modifier for normal attacks (eg. Scratch, Punch) **55**
- **Defense:** the base damage resistance against normal attacks **40**
- **SP Atk:** special attack, the base modifier for special attacks (e.g. fire blast, bubble beam) **50**
- **SP Def:** the base damage resistance against special attacks **90**

- **Speed** **rou**
- Can we predict the “type” of pokemon based on the information?

Ability value

The purpose is that if we know the type of the enemy Pokémon based on the ability value, we can choose the counter type of the Pokémon to battle.

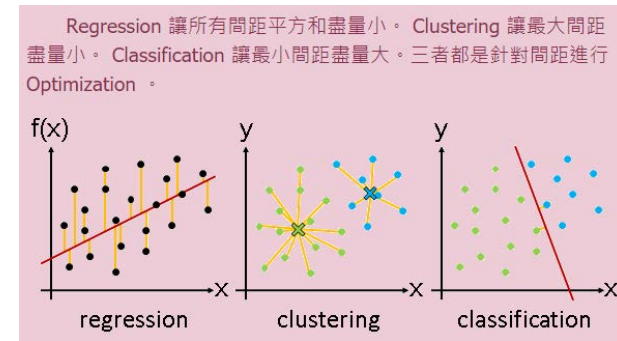
Example Application

×		防禦方的屬性																	
		一般	格鬥	飛行	毒	地面	岩石	蟲	幽靈	鋼	火	水	草	電	超能力	冰	龍	惡	妖精
攻擊方的屬性	一般	1×	1×	1×	1×	1×	1/2×	1×	0×	1/2×	1×	1×	1×	1×	1×	1×	1×	1×	1×
	格鬥	2×	1×	1/2×	1/2×	1×	2×	1/2×	0×	2×	1×	1×	1×	1×	1/2×	2×	1×	2×	1/2×
	飛行	1×	2×	1×	1×	1×	1/2×	2×	1×	1/2×	1×	1×	2×	1/2×	1×	1×	1×	1×	1×
	毒	1×	1×	1×	1/2×	1/2×	1/2×	1×	1/2×	0×	1×	1×	2×	1×	1×	1×	1×	1×	2×
	地面	1×	1×	0×	2×	1×	2×	1/2×	1×	2×	2×	1×	1/2×	2×	1×	1×	1×	1×	1×
	岩石	1×	1/2×	2×	1×	1/2×	1×	2×	1×	1/2×	2×	1×	1×	1×	1×	2×	1×	1×	1×
	蟲	1×	1/2×	1/2×	1/2×	1×	1×	1×	1/2×	1/2×	1/2×	1×	2×	1×	2×	1×	1×	2×	1/2×
	幽靈	0×	1×	1×	1×	1×	1×	1×	2×	1×	1×	1×	1×	1×	2×	1×	1×	1/2×	1×
	鋼	1×	1×	1×	1×	1×	2×	1×	1×	1/2×	1/2×	1/2×	1×	1/2×	1×	2×	1×	1×	2×
	火	1×	1×	1×	1×	1×	1/2×	2×	1×	2×	1/2×	1/2×	2×	1×	1×	2×	1/2×	1×	1×
	水	1×	1×	1×	1×	2×	2×	1×	1×	1×	2×	1/2×	1/2×	1×	1×	1×	1/2×	1×	1×
	草	1×	1×	1/2×	1/2×	2×	2×	1/2×	1×	1/2×	1/2×	2×	1/2×	1×	1×	1×	1/2×	1×	1×
	電	1×	1×	2×	1×	0×	1×	1×	1×	1×	1×	2×	1/2×	1/2×	1×	1×	1/2×	1×	1×
	超能力	1×	2×	1×	2×	1×	1×	1×	1×	1/2×	1×	1×	1×	1×	1/2×	1×	1×	0×	1×
	冰	1×	1×	2×	1×	2×	1×	1×	1×	1/2×	1/2×	1/2×	2×	1×	1×	1/2×	2×	1×	1×
	龍	1×	1×	1×	1×	1×	1×	1×	1×	1/2×	1×	1×	1×	1×	1×	1×	2×	1×	0×
	惡	1×	1/2×	1×	1×	1×	1×	1×	2×	1×	1×	1×	1×	1×	2×	1×	1×	1/2×	1/2×
妖精	1×	2×	1×	1/2×	1×	1×	1×	1×	1/2×	1/2×	1×	1×	1×	1×	1×	2×	2×	1×	

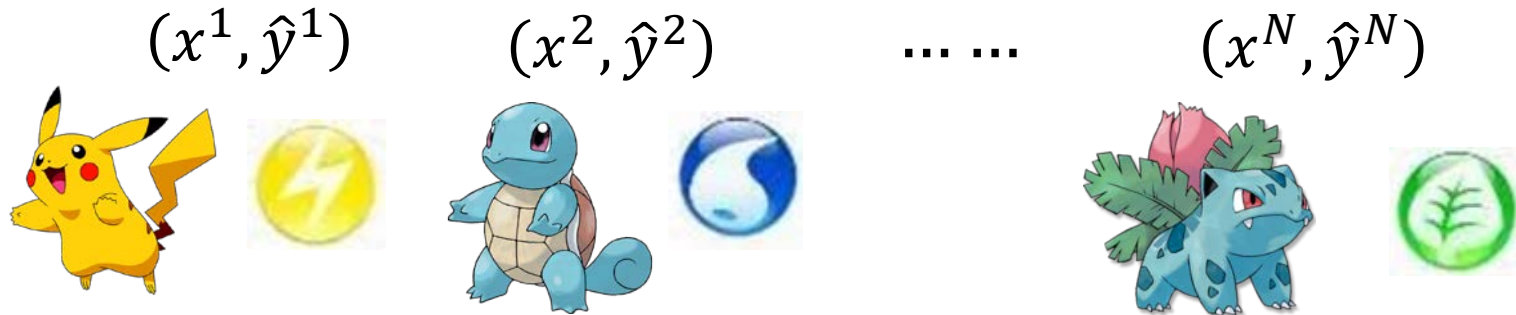
這些倍數適用於XY及之後的遊戲。

這些倍數適用於XY及之後的遊戲。

How to do Classification



- Training data for Classification



Classification as Regression?

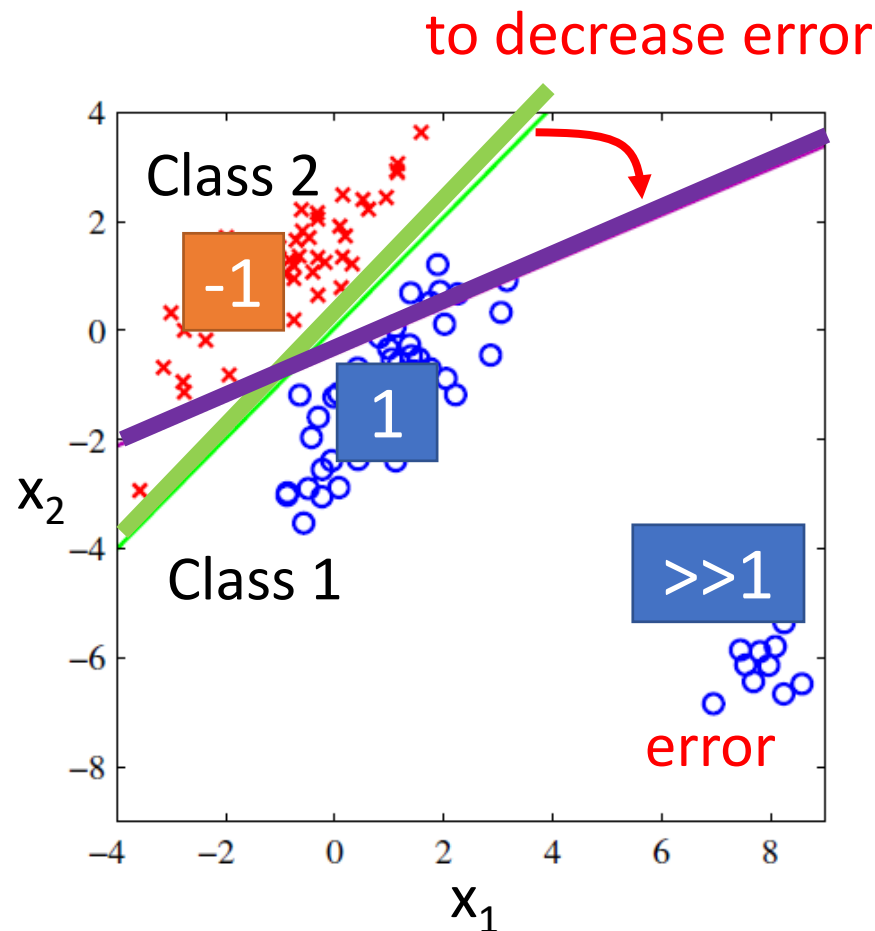
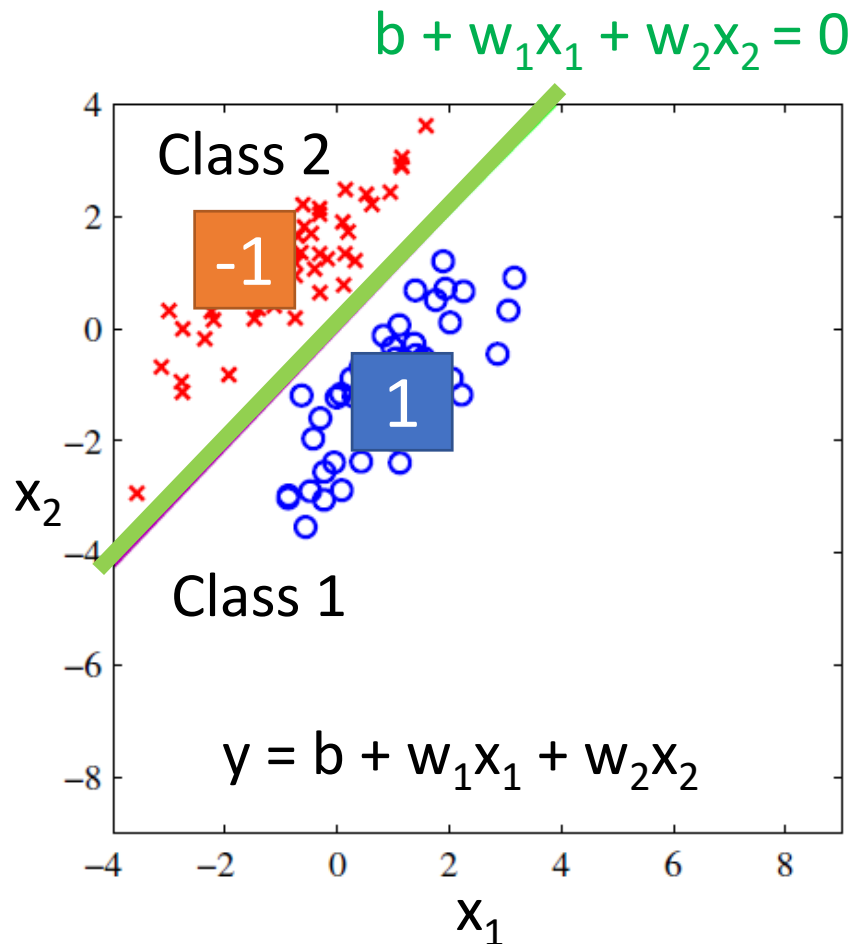
Binary classification as example

Use regression to train the classification problem.

⇒ Problematic (Next page)

Training: Class 1 means the target is 1; Class 2 means the target is -1

Testing: closer to 1 → class 1; closer to -1 → class 2



1° Penalize to the examples that are “too correct” ... (Bishop, P186)

- 2° • Multiple class: Class 1 means the target is 1; Class 2 means the target is 2; Class 3 means the target is 3 problematic

With square loss or cross entropy?
Discuss in the next chapter.

Ideal Alternatives

Change regression to classification.

→ Ideal: Add δ

Approximation: Add sigmoid

→ At the end of this chapter

- Function (Model):

x →

$g(x) > 0$
else

Output = class 1

Output = class 2

- Loss function:

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

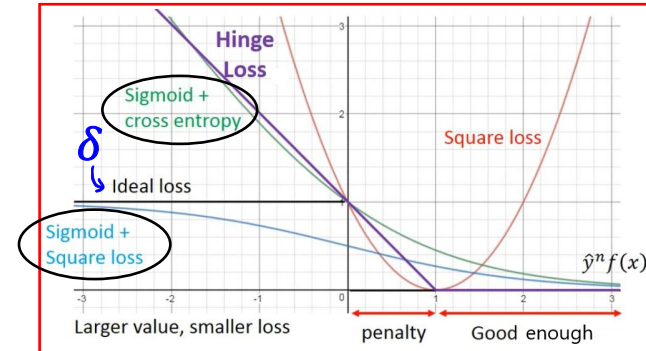
The number of times f get incorrect results on training data.

We can't differentiate this formula. ⇒ We can't use gradient descent.

- Find the best function: ⇒ Use sigmoid instead of δ

- Example: Perceptron, SVM

Not Today



Recall:

$$P(A) \text{ given } B: P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}$$

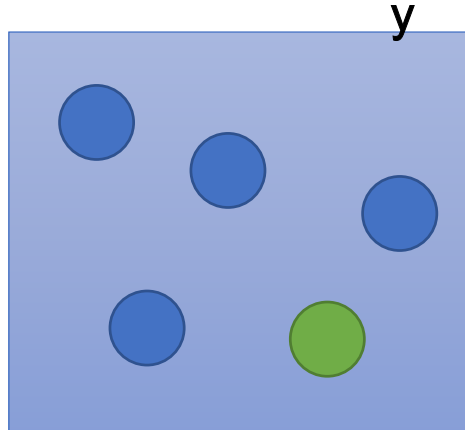
Two Boxes

We can regard a box as a class label, and regard a colorful ball as an object.

(B₁) Box 1

$$P(B_1) = 2/3$$

$$+ \\ P(B_2) \\ \parallel \\ 1$$

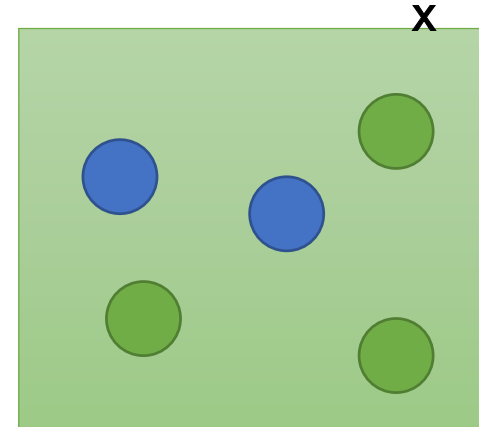


$$P(\text{Blue} | B_1) = 4/5$$

$$P(\text{Green} | B_1) = 1/5$$

(B₂) Box 2

$$P(B_2) = 1/3$$



$$P(\text{Blue} | B_2) = 2/5$$

$$P(\text{Green} | B_2) = 3/5$$

● from one of the boxes

Where does it come from?

(Distribution)

Likelihood

Prior



$$P(\text{Blue} | B_1) P(B_1)$$

Posterior

III

$$P(B_1 | \text{Blue}) =$$

$$\frac{P(\text{Blue} | B_1) P(B_1) + P(\text{Blue} | B_2) P(B_2)}{P(\text{Blue} | B_1) P(B_1) + P(\text{Blue} | B_2) P(B_2)} = P(B)$$

belongs to

Objective of
classification problem

Prior、Posterior 和 Likelihood 的理解与几种表达方式

1. 基本定义

Prior - 先验概率：预先知道事件发生的概率。

Posterior - 后验概率：Given evidence/experience/observation, Random event 发生的概率。

Likelihood - 似然：不符合概率的性质，但有 Probability 的含义。

2. 核心概念

$$\text{Posterior probability} \propto \text{Likelihood} \times \text{Prior probability.}$$

后验概率 \leftarrow 先验概率

3. 从两个事件的角度

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

4. 从假说和数据的角度

$$P(h|Data) = \frac{P(data|h)P(h)}{P(data)}$$

5. 从概率分布的参数估计的角度

$$\checkmark \quad p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$$

x : Observation
 θ : Probability 的参数

$p(x|\theta)$: Probability of x given θ .

Likelihood of θ given x observed.

Likelihood is a Function of θ .

6. 从估计 Sample 权重的角度

Predictions Matrix 表达 $y = X\omega$

X : matrix. 每一行是一个Sample, 每一行的每一个元素是一个 Feature.

Ω : weights

y : vector. 是X每一行与weights的乘积

Bayes 公式
求解 Posterior $p(w|X, y) \propto p(y|X, w)p(w)$

7. Likelihood的详细探讨

Likelihood : Given outcomes x, the likelihood of parameter θ .
Equal to the probability of observed x, given parameter θ .

$$\mathcal{L}(\theta|x) = P(x|\theta) \quad \mathcal{L}(\theta|x) \text{ 这个标记有时候容易引起误解}$$

Likelihood : Observed values x_1, x_2, \dots, x_n fixed, θ is variable.

$$\mathcal{L}(\theta; x_1, \dots, x_n) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta).$$

用分号分割避免误解

$$P(\theta | x_1, x_2, \dots, x_n) = \frac{f(x_1, x_2, \dots, x_n | \theta)P(\theta)}{P(x_1, x_2, \dots, x_n)}$$

$P(\theta)$ 是参数 θ 的先验分布。

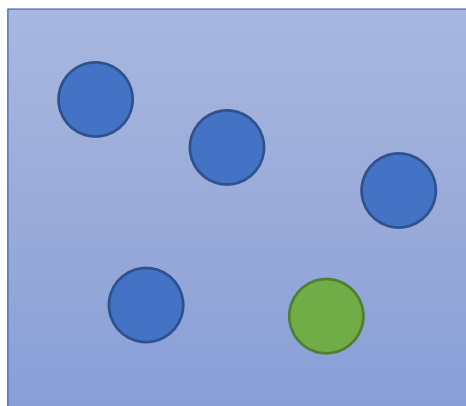
In this chapter, we only assume there are

Two Classes

Estimating the Probabilities
From training data

Class 1

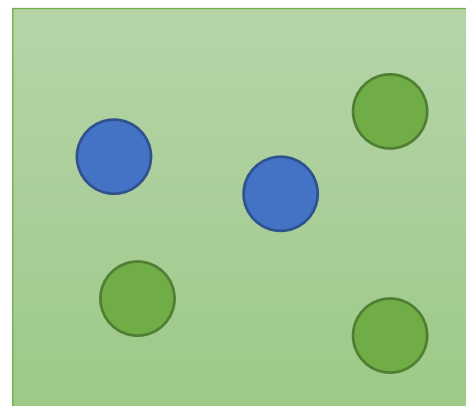
$P(C_1)$



$P(x|C_1)$

Class 2

$P(C_2)$



$P(x|C_2)$

Given an x , which class does it belong to

We can use prior probability
and likelihood (distribution) to
generate posterior probability
or object's probability.

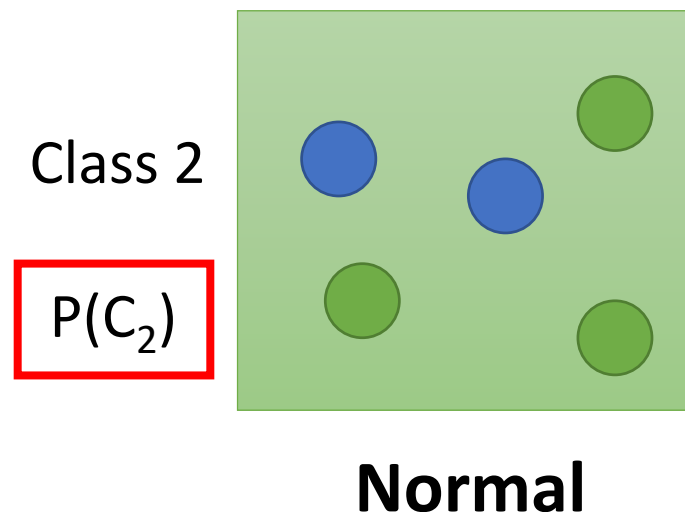
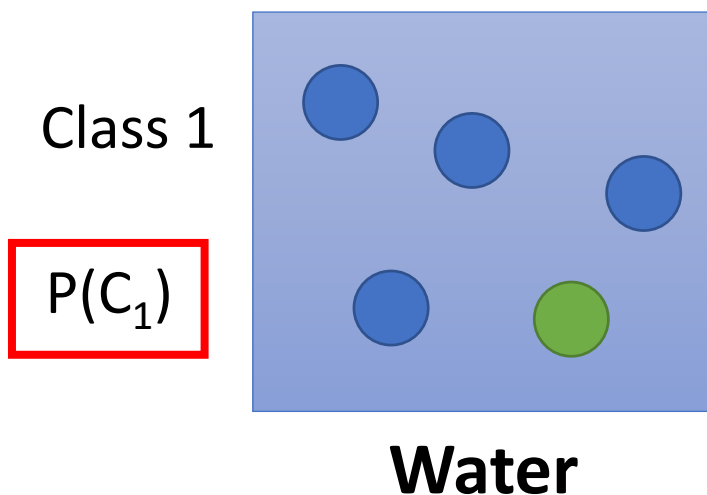
$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

Generative Model $P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$

$$P(x) = \sum P(x|C_i) P(C_i)$$

For Pokémon instead of colorful balls:

Prior



Water and Normal type with ID < 400 for training,
rest for testing

Training: 79 Water, 61 Normal

$$P(C_1) = 79 / (79 + 61) = 0.56$$

$$P(C_2) = 61 / (79 + 61) = 0.44$$

Probability from Class

There is no sea turtle in training data, so is the probability of sea turtle belongs to the water type equal to 0? **No** (Next page)

$$P(x | C_1) = ?$$

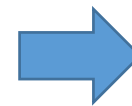
$$P($$



$$| \text{Water}) = ?$$

Likelihood

Each Pokémon is represented as a vector by its attribute.



feature

**Water
Type**

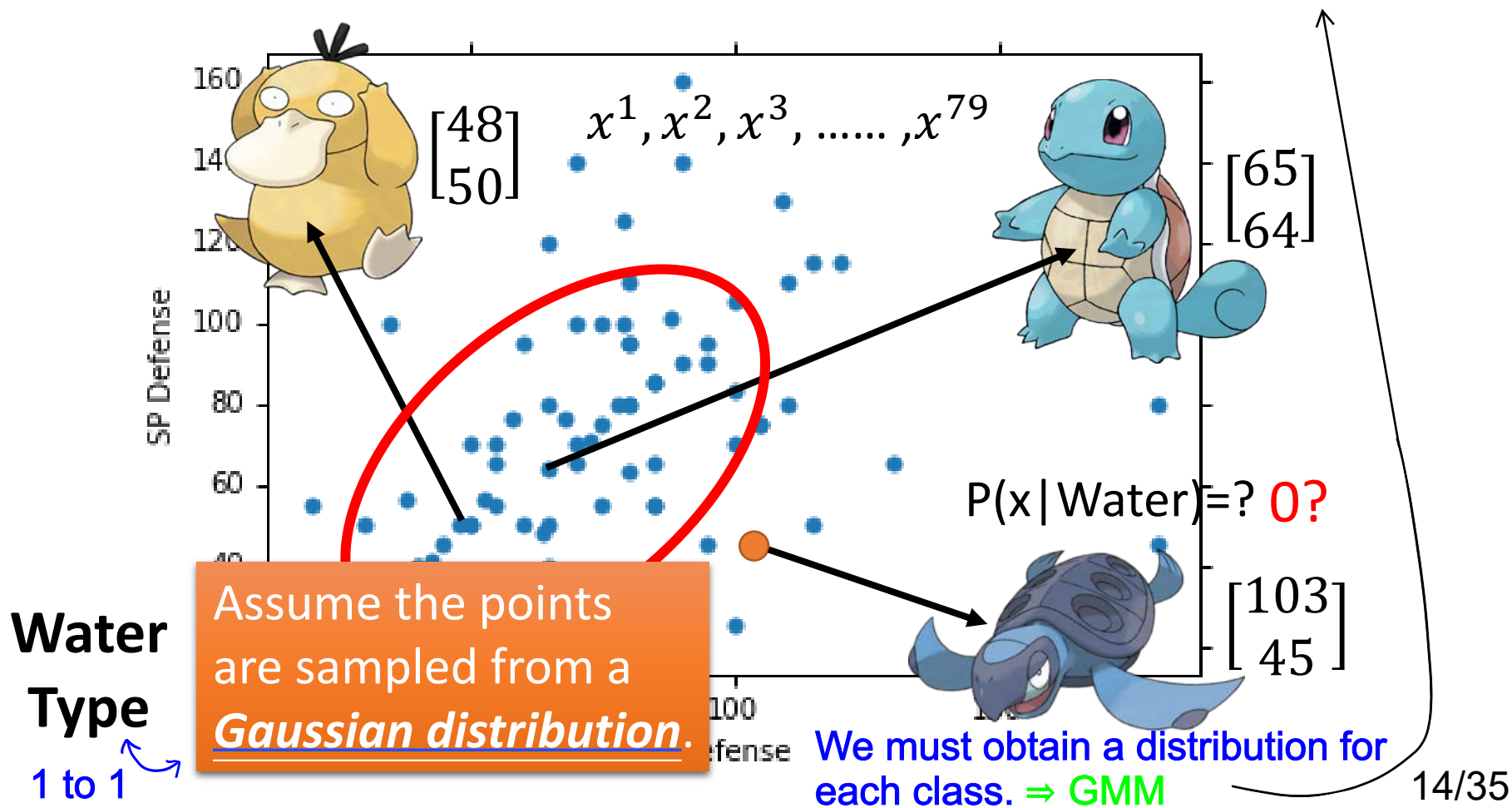


Probability from Class - Feature

We form the distribution of the water Pokémon based on their features.
(ex: Defense and SP Defense)

- Considering **Defense** and **SP Defense**

#Gaussian = #class



Gaussian Distribution

Be used to calculate likelihood

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Input: vector x , output: probability of sampling x

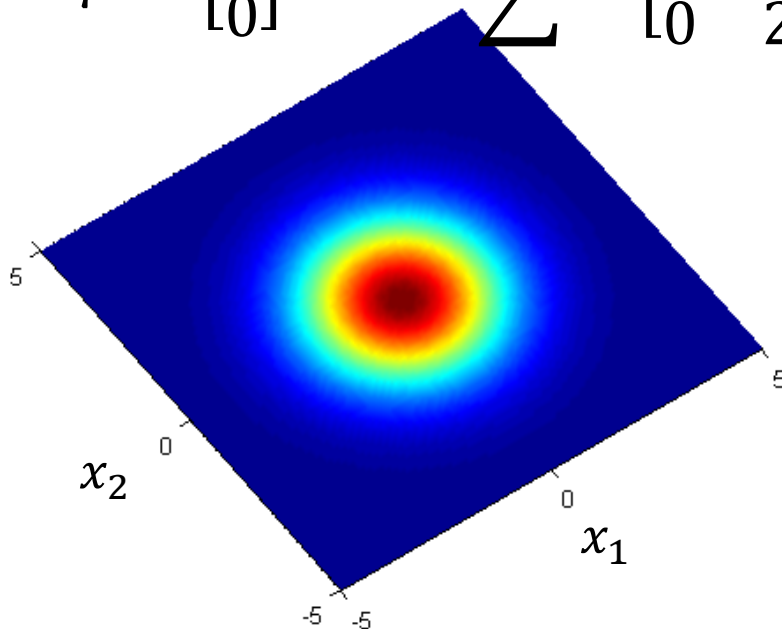
The shape of the function determines by **mean μ** and **covariance matrix Σ**

come from C_i

We can calculate $P(x|C_i)$ through the distribution of C_i . ($f_{\mu, \Sigma}(x) = P(x|C_i)$)

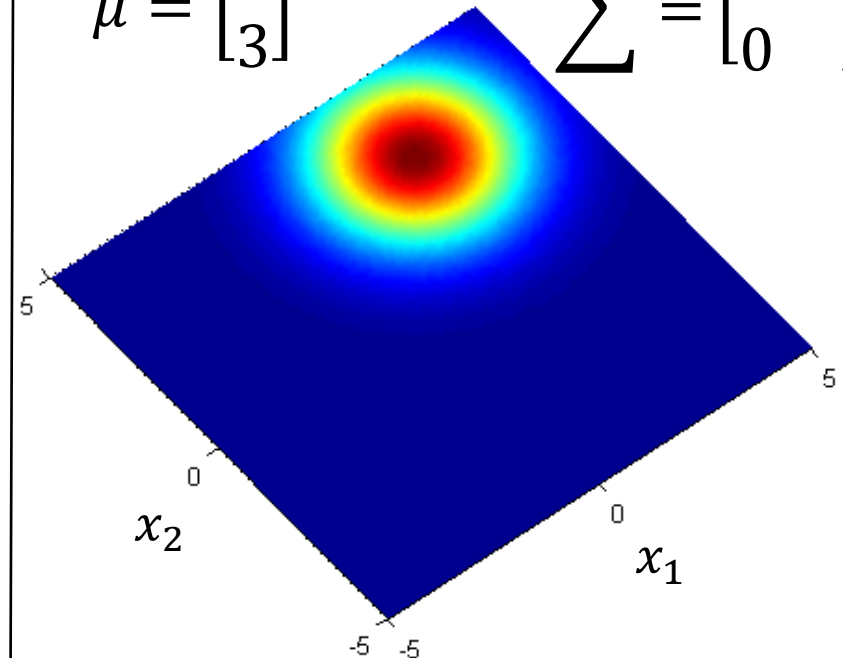
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

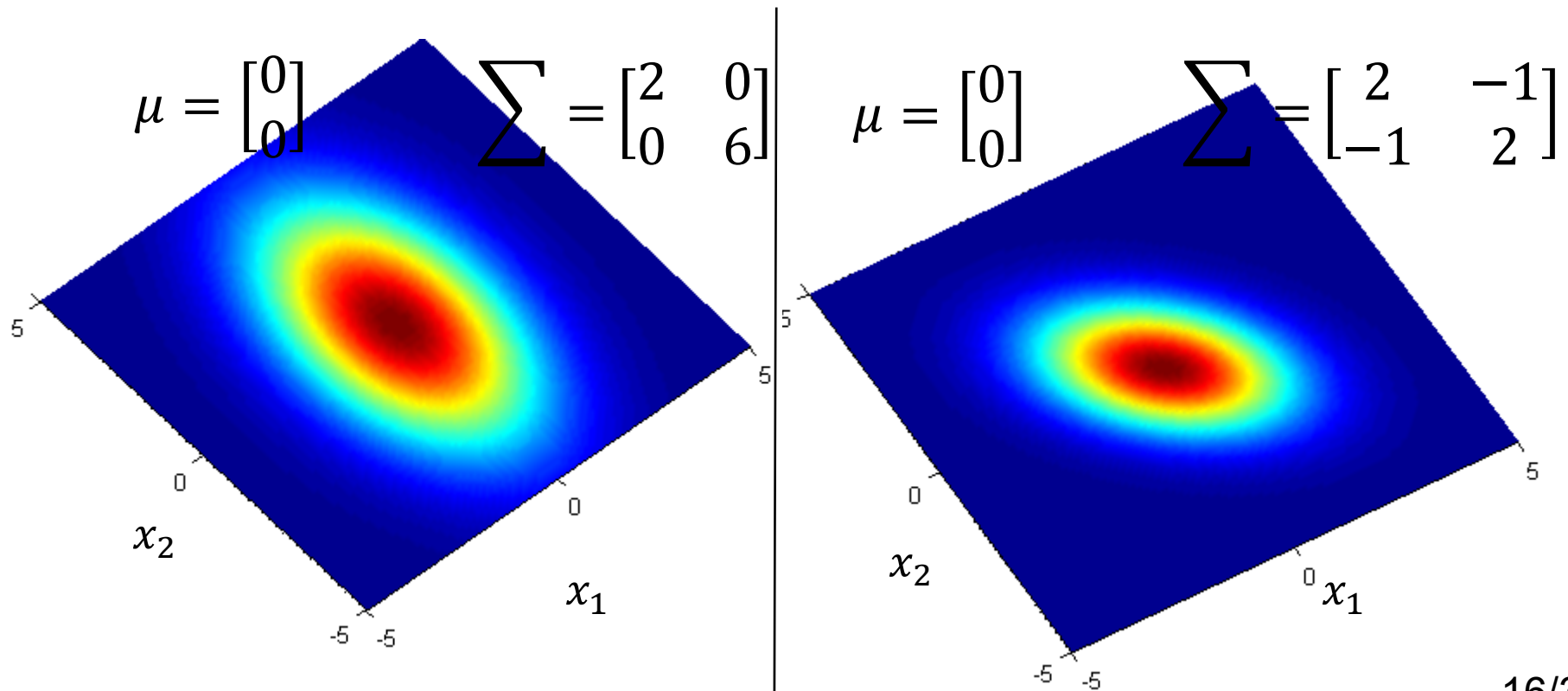


Gaussian Distribution

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Input: vector x , output: probability of sampling x

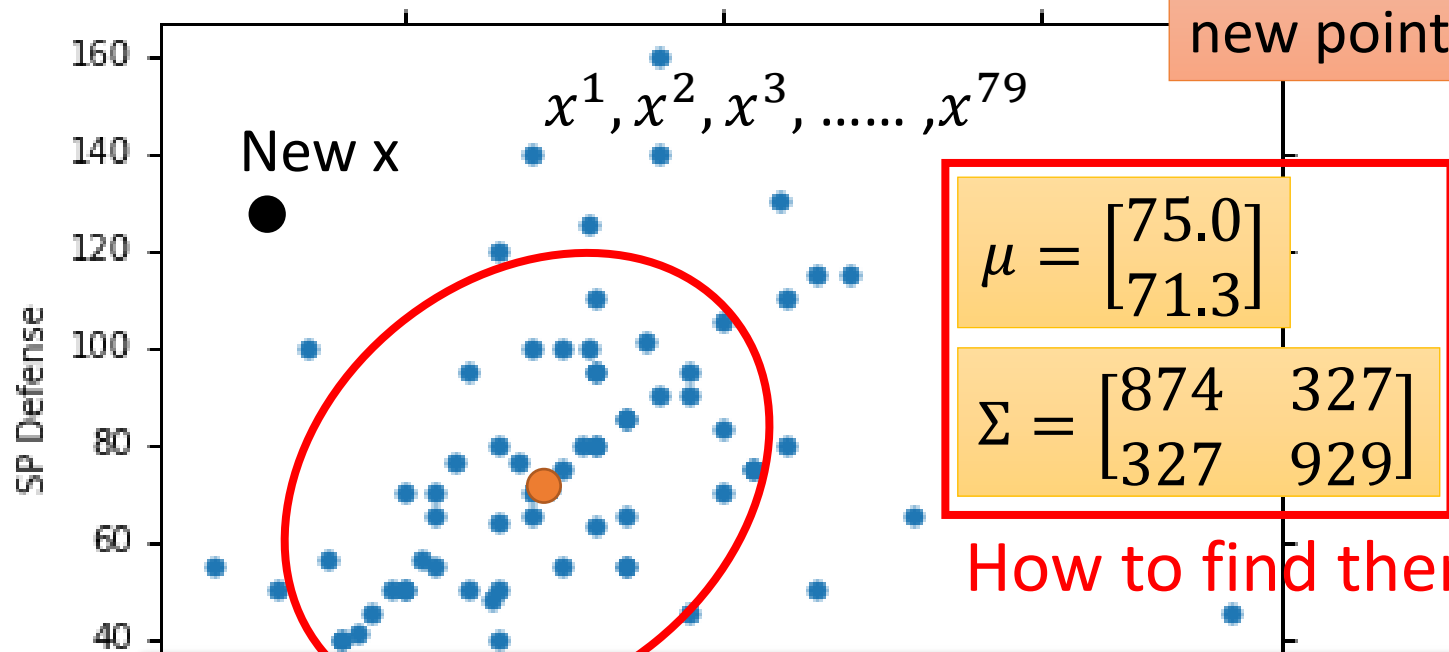
The shape of the function determines by **mean μ** and **covariance matrix Σ**



Probability from Class

Assume the points are sampled from a Gaussian distribution

Find the Gaussian distribution behind them → Probability for new points



How to find them? **MLE**

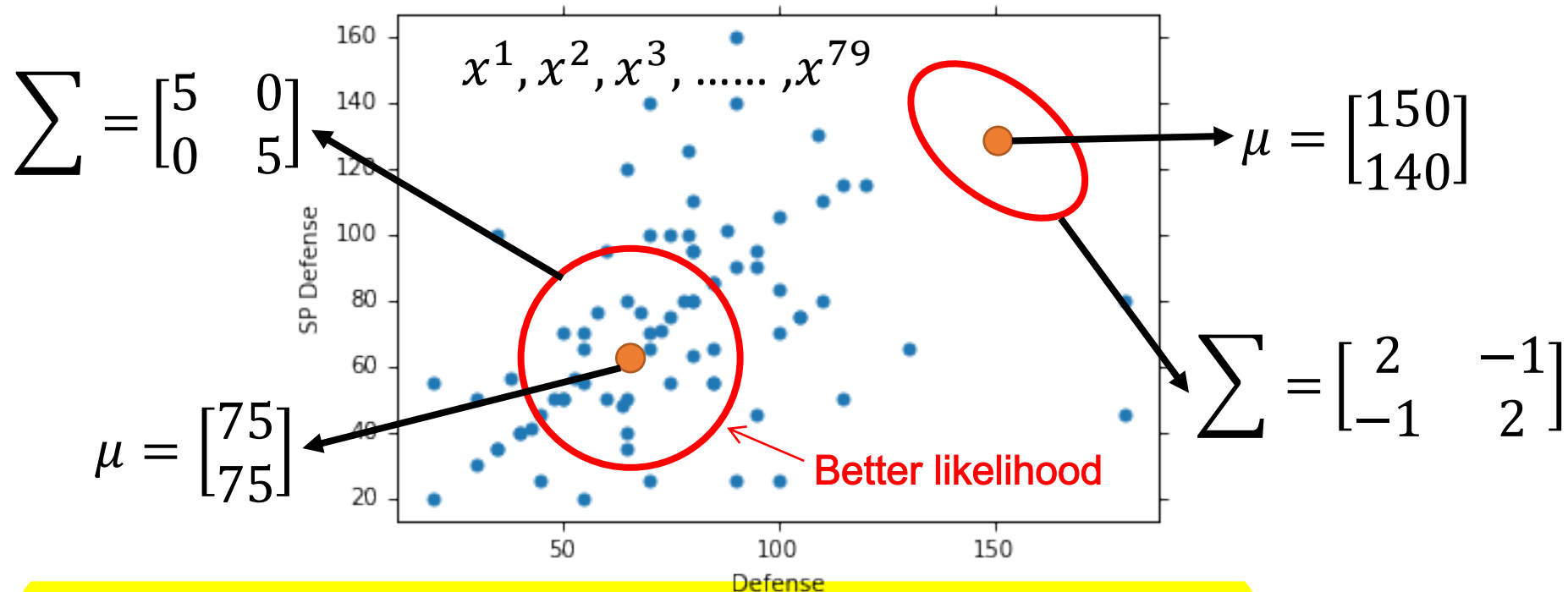
**Water
Type**

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Defense

Likelihood is kind of like the objective function of Gaussian distribution.

Maximum Likelihood $f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$



The Gaussian with any mean μ and covariance matrix Σ can generate these points. ➡ Different Likelihood

Likelihood of a Gaussian with mean μ and covariance matrix Σ

The bigger the better ➡ = the probability of the Gaussian samples $x^1, x^2, x^3, \dots, x^{79}$
 L stands for likelihood instead of loss function.

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

It looks like we don't need y to calculate likelihood, but we need to remember that these X come from the same y . So it's still a supervised learning.

Maximum Likelihood

Supervised GMM: classification
Unsupervised GMM: clustering

We have the “Water” type Pokémons: $x^1, x^2, x^3, \dots, x^{79}$

We assume $x^1, x^2, x^3, \dots, x^{79}$ generate from the Gaussian (μ^*, Σ^*) with the **maximum likelihood**

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

Closed-form solution for each class

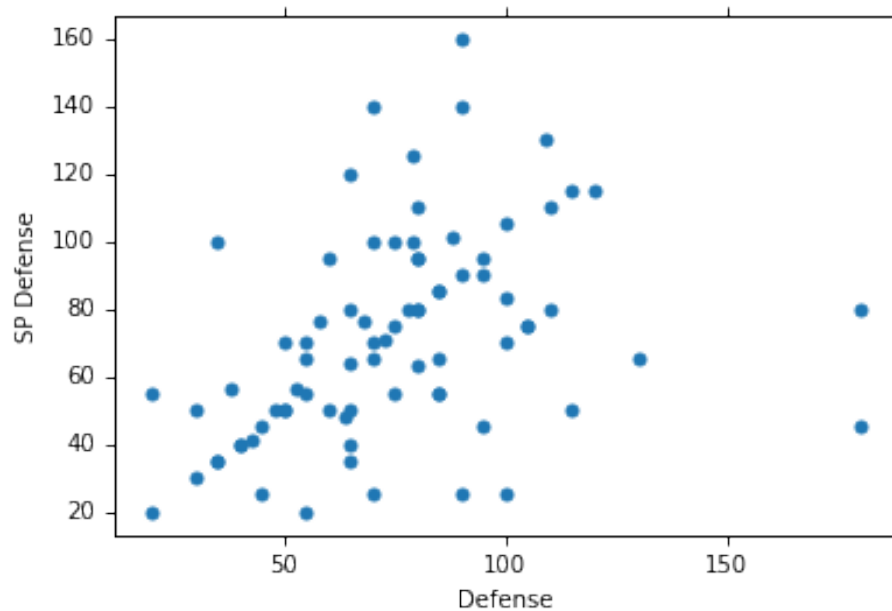
$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x^n$$

average

$$\Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x^n - \mu^*) (x^n - \mu^*)^T$$

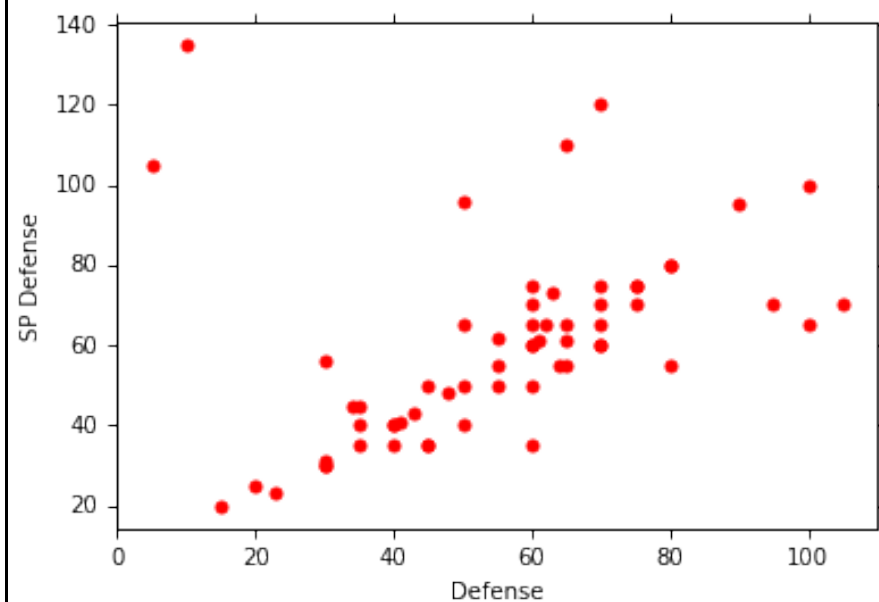
Maximum Likelihood

Class 1: Water



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

Class 2: Normal



$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

Now we can do classification 😊

$$f_{\mu^1, \Sigma^1}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1)\right\}$$
$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

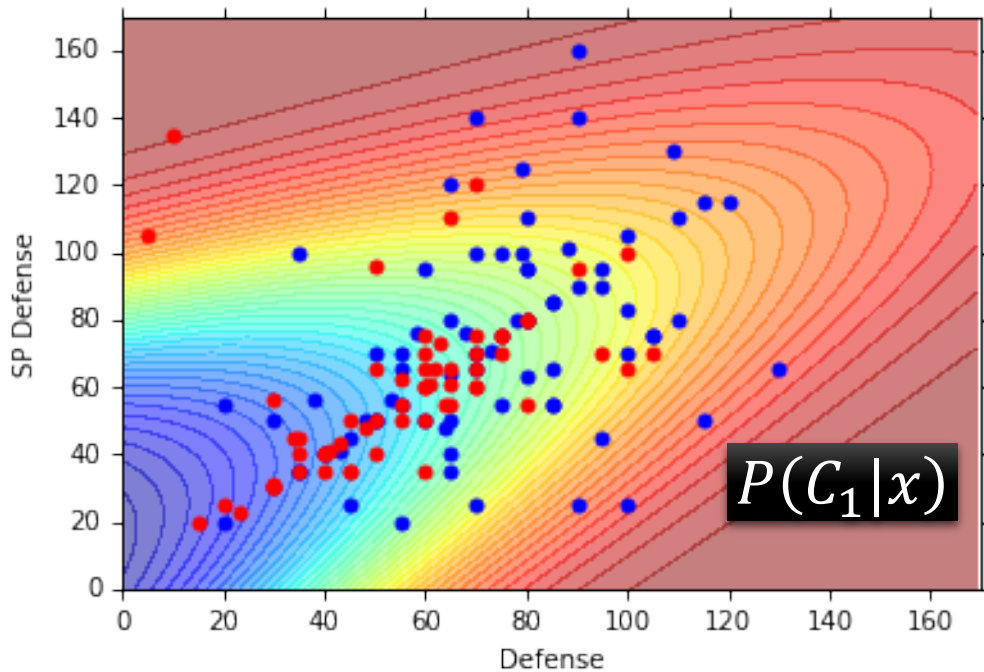
$P(C_1) = 79 / (79 + 61) = 0.56$

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$f_{\mu^2, \Sigma^2}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)\right\}$$
$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

$P(C_2) = 61 / (79 + 61) = 0.44$

If $P(C_1|x) > 0.5$ ➡ x belongs to class 1 (Water)



Blue points: C_1 (Water), Red points: C_2 (Normal)

How's the results?

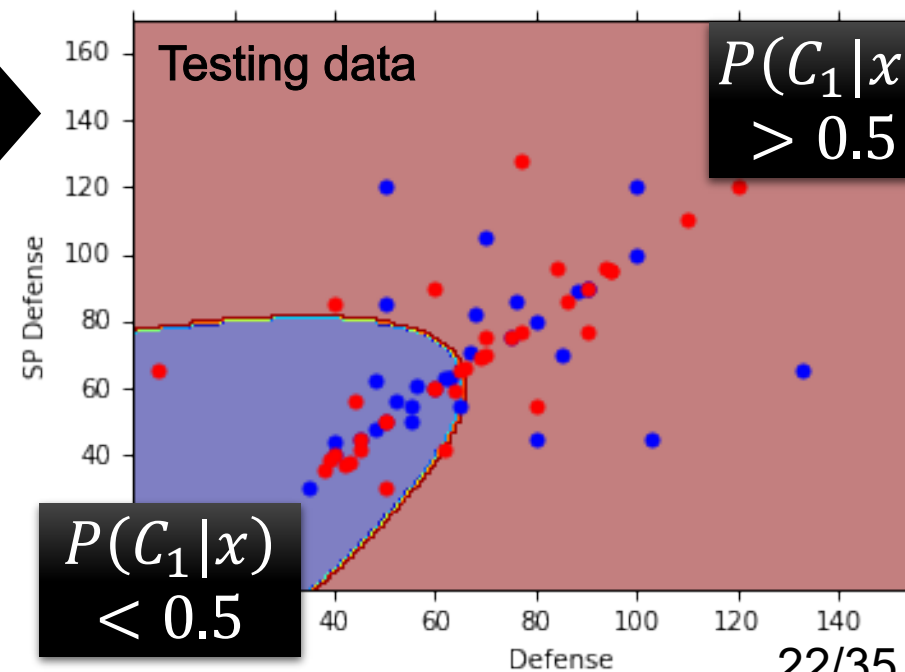
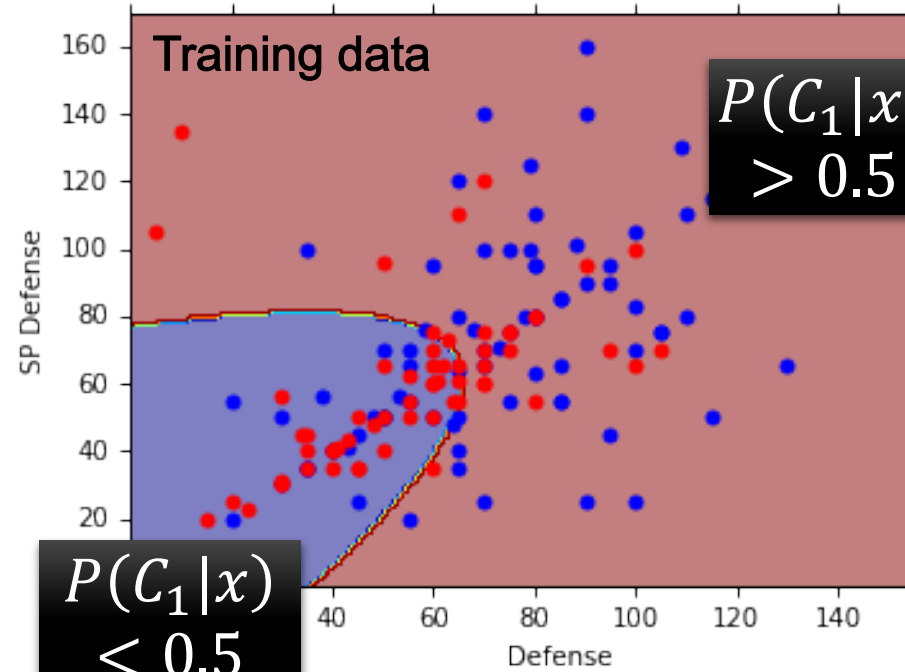
Testing data: 47% accuracy ☹️

All: hp, att, sp att,
de, sp de, speed (6 features)

μ^1, μ^2 : 6-dim vector

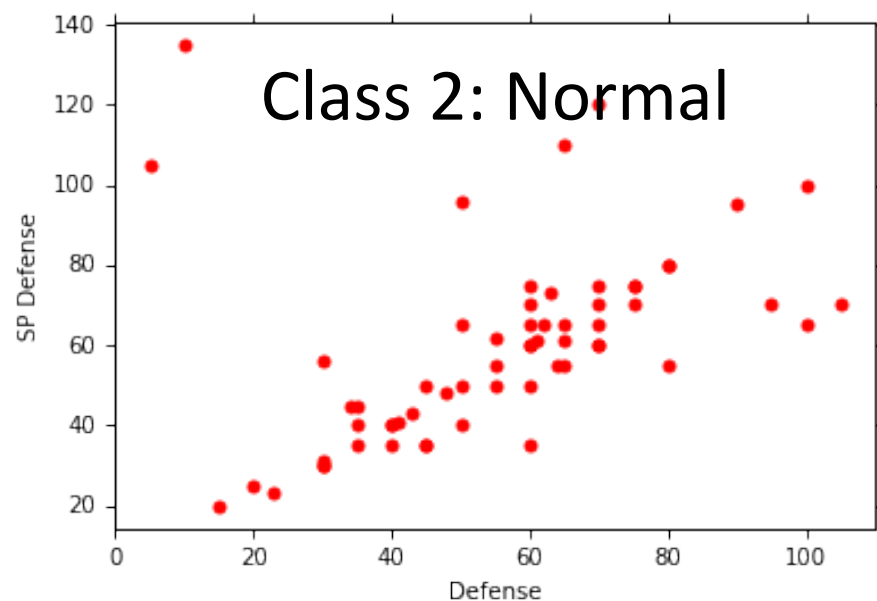
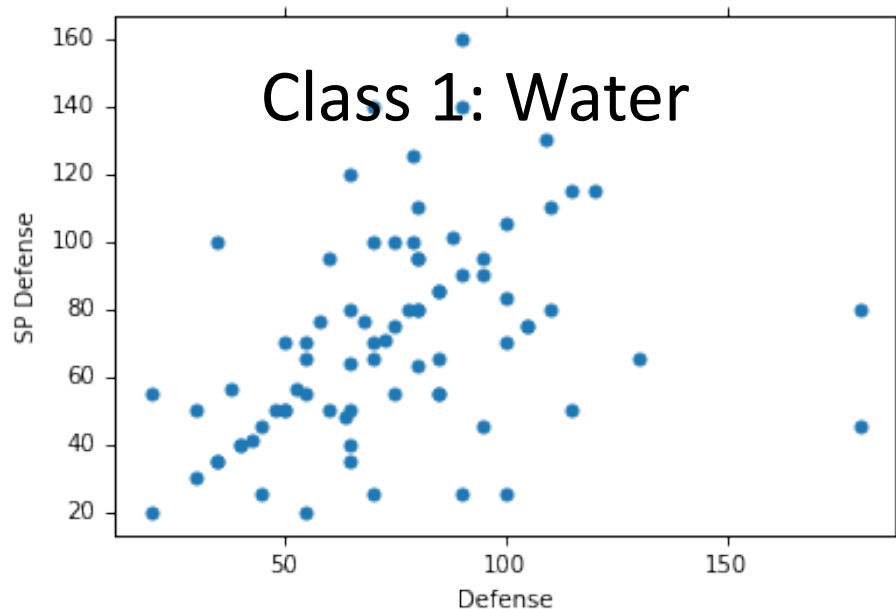
Σ^1, Σ^2 : 6 x 6 matrices

64% accuracy ...



Modifying Model

There are too many parameters in the model. \Rightarrow Overfitting



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

#parameters: n

n^2

The same Σ

We can share the covariance matrix. Less parameters

Modifying Model

Ref: Bishop,
chapter 4.2.2

- Maximum likelihood

$L(\mu^1, \Sigma^1)$
“Water” type Pokémons:

$x^1, x^2, x^3, \dots, x^{79}$

μ^1

$L(\mu^2, \Sigma^2)$
“Normal” type Pokémons:

$x^{80}, x^{81}, x^{82}, \dots, x^{140}$

μ^2

Σ

Find μ^1, μ^2, Σ maximizing the likelihood $L(\mu^1, \mu^2, \Sigma)$

$$L(\mu^1, \mu^2, \Sigma) = f_{\mu^1, \Sigma}(x^1) f_{\mu^1, \Sigma}(x^2) \cdots f_{\mu^1, \Sigma}(x^{79}) \\ \times f_{\mu^2, \Sigma}(x^{80}) f_{\mu^2, \Sigma}(x^{81}) \cdots f_{\mu^2, \Sigma}(x^{140})$$

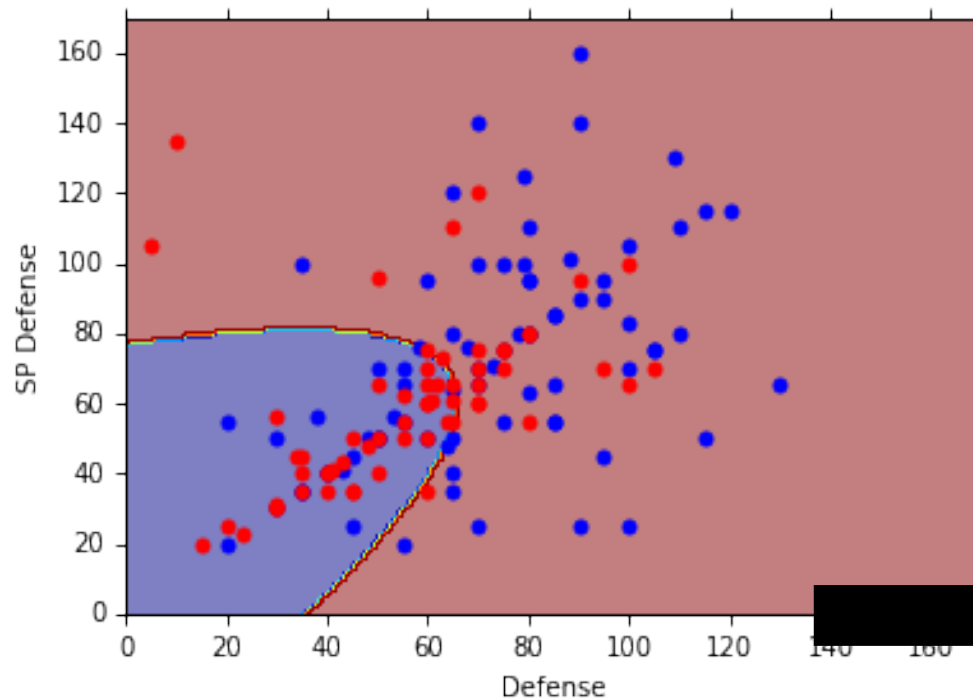
Compare to the original μ, Σ .

μ^1 and μ^2 is the same

$$\Sigma = \frac{79}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$$

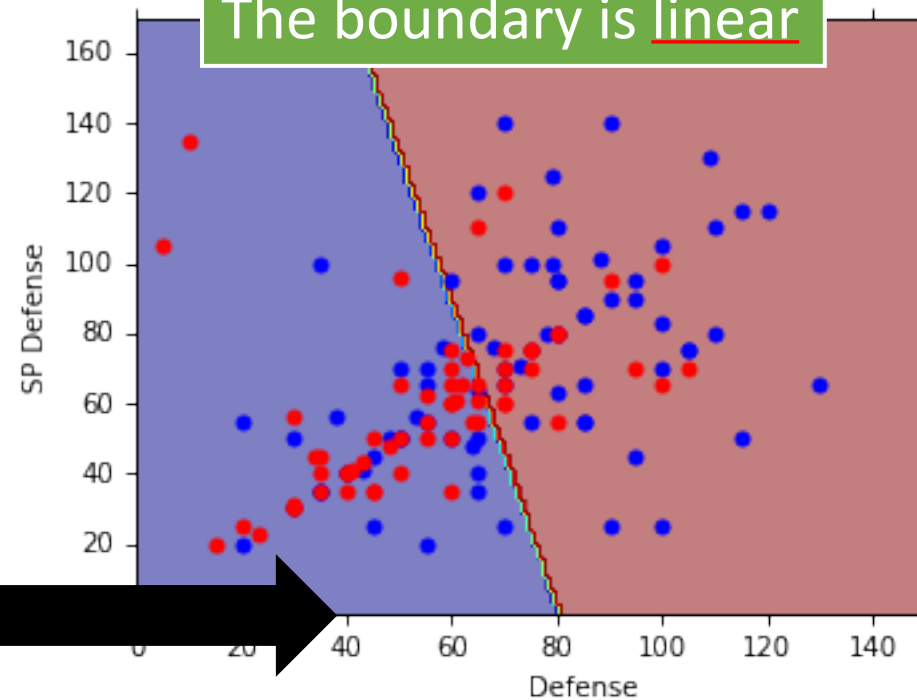
Modifying Model

Only use Defense and SP Defense:



$\therefore \equiv \sigma(wx+b)$

The boundary is linear



The same covariance matrix

All: hp, att, sp att, de, sp de, speed

54% accuracy \longrightarrow 73% accuracy

Three Steps

- Function Set (Model):

x 

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

If $P(C_1|x) > 0.5$, output: class 1
Otherwise, output: class 2

- Goodness of a function: **Distribution's likelihood**
 - The mean μ and covariance Σ that maximizing the likelihood (the probability of generating data)
- Find the best function: easy **MLE**

Probability Distribution

instead of Gaussian distribution

- You can always use the distribution you like 😊

We can assume that the distribution of each feature is independent.

$$P(x|C_1) = P(x_1|C_1) P(x_2|C_1) \cdots P(x_k|C_1) \cdots$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \\ \vdots \\ x_K \end{bmatrix}$$

1-D Gaussian

The covariance matrix will be diagonal.

ex: Is it a legendary Pokémon?

For binary features, you may assume they are from Bernoulli distributions.

It's very trivial that we can't choose Gaussian distribution for binary features.

If you assume all the dimensions are independent, then you are using Naive Bayes Classifier.

Posterior Probability

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

After some mathematical derivations, we can find that we can implement this in discriminative fashion. (We can then use gradient descent.)

GMM with shared covariance matrix

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

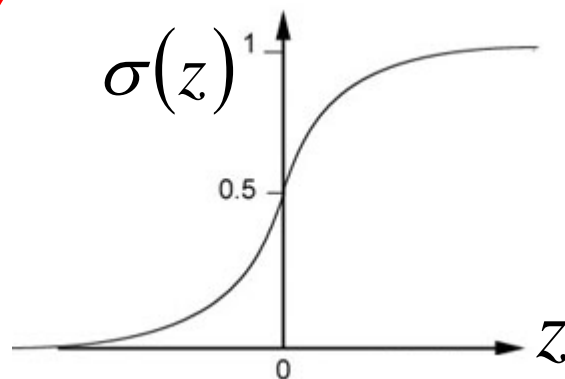
$$= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}}$$

$$= \frac{1}{1 + \exp(-z)}$$

$$= \sigma(\overset{wx+b}{\underset{\parallel}{z}})$$

Sigmoid function

$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$



Warning of Math

Prove that $z = wx + b$

Posterior Probability

$$P(C_1|x) = \sigma(z) \quad \text{sigmoid} \quad z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \rightarrow \frac{\frac{N_1}{N_1 + N_2}}{\frac{N_2}{N_1 + N_2}} = \frac{N_1}{N_2}$$

#class 1's occurrences

#class 2's occurrences

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$\ln \frac{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}}{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)] \right\}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)]$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} \left[\underbrace{(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1)}_{\textcircled{1}} - \underbrace{(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)}_{\textcircled{2}} \right]$$

$$\textcircled{1}: (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1)$$

$$= x^T (\Sigma^1)^{-1} x - \underbrace{x^T (\Sigma^1)^{-1} \mu^1 - (\mu^1)^T (\Sigma^1)^{-1} x}_{\text{blue underline}} + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$= x^T (\Sigma^1)^{-1} x - \underbrace{2(\mu^1)^T (\Sigma^1)^{-1} x}_{\text{blue underline}} + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$\textcircled{2}: (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)$$

$$= x^T (\Sigma^2)^{-1} x - 2(\mu^2)^T (\Sigma^2)^{-1} x + (\mu^2)^T (\Sigma^2)^{-1} \mu^2$$

$$z = \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

End of Warning

$$P(C_1|x) = \sigma(z)$$

$$z = \cancel{\ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}}} - \cancel{\frac{1}{2} x^T (\Sigma^1)^{-1} x} + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 + \cancel{\frac{1}{2} x^T (\Sigma^2)^{-1} x} - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

$$\Sigma_1 = \Sigma_2 = \Sigma \quad \text{Assume that the covariance matrices are the same.}$$

$$z = \underbrace{(\mu^1 - \mu^2)^T \Sigma^{-1} x}_{\mathbf{w}^T \text{ vector}} - \underbrace{\frac{1}{2} (\mu^1)^T \Sigma^{-1} \mu^1 + \frac{1}{2} (\mu^2)^T \Sigma^{-1} \mu^2}_{\mathbf{b} \text{ scalar}} + \ln \frac{N_1}{N_2}$$

Discriminative model
⇒ Use gradient descent

$$P(C_1|x) = \sigma(w \cdot x + b)$$

How about directly find \mathbf{w} and b ?

This is why the boundary is linear for shared covariance matrix.

In generative model, we estimate $N_1, N_2, \mu^1, \mu^2, \Sigma$

We may obtain better result with different \mathbf{w}, b . Then we have \mathbf{w} and b

Reference

- Bishop: Chapter 4.1 – 4.2
- Data: <https://www.kaggle.com/abcsds/pokemon>
- Useful posts:
 - <https://www.kaggle.com/nishantbhadauria/d/abcsds/pokemon/pokemon-speed-attack-hp-defense-analysis-by-type>
 - <https://www.kaggle.com/nikos90/d/abcsds/pokemon/mastering-pokebars/discussion>
 - <https://www.kaggle.com/ndrewgele/d/abcsds/pokemon/visualizing-pok-mon-stats-with-seaborn/discussion>