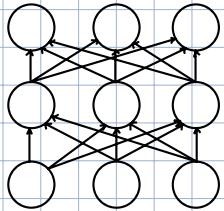


DEEP LEARNING Book - IAN GOODFELLOW (SOME POINTS FROM CHAPTER 1)

- QUINTESSENTIAL OF REPRESENTATION LEARNING ALGO → **AUTOENCODER**
- DEEP LEARNING → SOLVES CENTRAL PROBLEM OF REPRESENTATION LEARNING. → BUILDING COMPLEX CONCEPTS OUT OF SIMPLE CONCEPTS.
- ESSENTIAL OF DEEP LEARNING → **MULTILAYER PERCEPTRON (MLP)** OR **FEEDFORWARD DEEP NETWORK**
 - E.G. BREAKS DOWN COMPLICATED MAPPING INTO SERIES OF NESTED MAPPINGS.



- DEEP LEARNING HAD A LOT OF DIFFERENT NAMES. → CYBERNETICS (40s-60s) → CONNECTIONISM (80-90s) → DEEP LEARNING (OG)
- MORE USEFUL WITH MORE DATA AVAILABLE
- GROWN IN SIZE → COMPUTATIONAL INFRASTRUCTURE
- SOLVING INCREASINGLY COMPLEX PROBLEMS
- ONE OF DEEP LEARNING NAMES ⇒ **ARTIFICIAL NEURAL NETWORK (ANN)** → ENGINEERED SYSTEMS INSPIRED BY BRAIN.

- NEUROSCIENCE HAS GIVEN REASON TO HOPE THAT A SINGLE DL ALGO CAN SOLVE MANY DIFFERENT TASKS → IMPORTANT SOURCE OF INSPIRATION.
 - ↳ NEOCOGNITION → BASIS OF MODERN CONVOLUTIONAL NETWORK. (POWERFUL MODEL ARCHITECTURE FOR PROCESSING IMAGES)
- MOST NEURAL NETWORKS TODAY → BASED ON **RECTIFIED LINEAR UNIT**
- DON'T VIEW DEEP LEARNING AS AN ATTEMPT TO SIMULATE BRAIN.

CHAPTER 3

CHOOSING CLASSIFICATION ALGORITHM

- TAKES PRACTICE AND TIME
- "**NO FREE LUNCH THEOREM**" - DAVID H. WOLPERT - THERE ISN'T ONE CLASSIFIER THAT WORKS BEST FOR ALL SCENARIOS.
↳ CHOOSE A HANDFUL OF MODELS, THEN SELECT ONE THAT WORKS. → MODEL TRAINING VERY DEPENDANT ON DATA OVER TRAINING MODEL.
- **5 STEPS FOR SUPERVISED TRAINING:**
 1. SELECT FEATURES + COLLECT LABELED TRAINING DATA.
 2. CHOOSE PERFORMANCE METRIC
 3. CHOOSE LEARNING ALGORITHM + TRAINING MODEL
 4. EVALUATE PERFORMANCE OF MODEL.
 5. FINE-TUNE MODEL
- **DISCLAIMER** → FEATURE SELECTION + PRE PROCESSING + ETC.. → COVERED LATER IN BOOK

SCIKIT-LEARN - PERCEPTRON

- USER-FRIENDLY + OPTIMISED IMPLEMENTATION OF SEVERAL CLASSIFICATION ALGORITHMS.
- BUILDING PERCEPTRON WITH SCIKIT-LEARN USING IRIS DATASET.
- 'train-test-split' → SPLIT DATA INTO TEST + TRAINING DATA.
- MANY MACHINE LEARNING + OPTIMISATION ALGORITHMS → REQUIRE FEATURE SCALING FOR OPTIMAL PERFORMANCE → E.G. GRADIENT DESCENT.

LOGISTIC REGRESSION + CONDITIONAL PROBABILITIES

- LOGISTIC REGRESSION → **CLASSIFICATION MODEL** → PERFORMS WELL ON LINEARLY SEPARABLE CLASSES.
- **ONE OF THE MOST WIDELY USED ALGORITHMS FOR CLASSIFICATION**
- LINEAR MODEL FOR BINARY CLASSIFICATION.
- LOGISTIC REGRESSION FOR MULTIPLE CLASSES → **MULTINOMIAL LOGISTIC REGRESSION OR SOFTMAX REGRESSION** (LOOK INTO IT)
- **MAIN MECHANICS:**

- **ODDS** → ODDS IN FAVOUR OF PARTICULAR EVENT. $\Rightarrow \frac{P}{(1-p)}$, P = PROBABILITY OF POSITIVE EVENT.
POSITIVE EVENT LABEL $\rightarrow y=1 \Rightarrow p := P(y=1|x)$

↳ EVENT WE WANT TO PREDICT. DOESN'T NECESSARILY MEAN "Good".

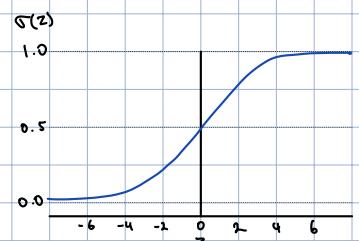
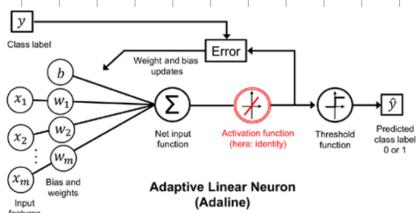
- **LOGIT FUNCTION** $\Rightarrow \text{logit}(p) = \log \frac{p}{(1-p)} \Rightarrow \log = \text{NATURAL LOG}$
↳ LOG OF ODDS
↳ TAKES INPUT VALUE IN THE RANGE 0 TO 1.

LOGIT FUNCTION MAPS PROBABILITY TO REAL-NUMBER RANGE
↳ CONSIDER INVERSE OF FUNCTION → MAP BACK TO $[0,1]$ RANGE.

$$\text{logit}(p) = w_1x_1 + \dots + w_mx_m + b = \sum_{j=1}^m w_j x_j + b = \mathbf{w}^T \mathbf{x} + b$$

- INVERSE OF LOGIT FUNCTION → **LOGISTIC SIGMOID FUNCTION** → **ABB. SIGMOID FUNCTION** → **S SHAPE**

$$\sigma(z) = \frac{1}{1+e^{-z}} \rightarrow z \text{ NET INPUT} \rightarrow z = \mathbf{w}^T \mathbf{x} + b$$



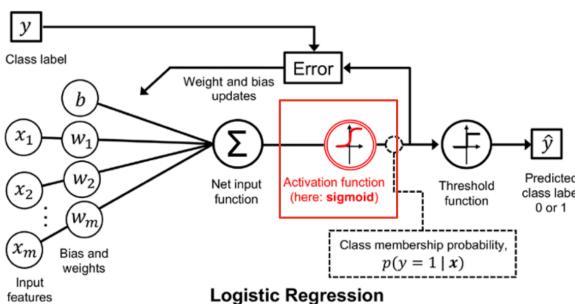
- $\sigma(z) \rightarrow 1$ IF $z \rightarrow \infty$ BECAUSE e^{-z} BECOMES VERY SMALL FOR LARGE VALUES.
- $\sigma(z) \rightarrow 0$ FOR $z \rightarrow -\infty$ AS A RESULT OF AN INCREASINGLY LARGE DENOMINATOR.
- CONCLUSION → SIGMOID TAKES REAL-NUMBER VALUES AS INPUT AND TRANSFORMS THEM INTO VALUES IN RANGE $[0,1]$, INTERCEPT $\sigma(0) = 0.5$

- E.G. $\rightarrow \sigma(z) = 0.8 = 80\%$ IS A IRIS-VERSICOLOR

$$p(y=0|x, w, b) = 1 - p(y=1|x, w, b) = 0.2 \rightarrow 20\% \quad \left| \begin{array}{l} 1 \\ 0 \end{array} \right. \begin{cases} 1 & \text{IF } z \geq 0.0 \\ 0 & \text{otherwise} \end{cases}$$

PREDICTED PROBABILITY CAN THEN SIMPLY BE CONVERTED INTO BINARY OUTCOME.

- LOGISTIC REGRESSION USED IN WEATHER FORECASTING. ALSO POPULAR IN MEDICINE.



MODEL WEIGHTS VIA LOGISTIC LOSS FUNCTION

- MINIMIZE LOSS TO LEARN PARAMETERS FOR ADALINE CLASSIFICATION MODEL.

$$L(w, b | x) = p(y|x; w, b) - \prod_{i=1}^n p(y^{(i)}|x^{(i)}; w, b) = \prod_{i=1}^n (\sigma(z^{(i)}))^{y^{(i)}} (1 - \sigma(z^{(i)}))^{1-y^{(i)}}$$

↳ IN PRACTICE → LOG-LIKELIHOOD

$$L(w, b | x) = \log L(y|x; w, b) = \sum_{i=1}^n [y^{(i)} \log(\sigma(z^{(i)})) + (1+y^{(i)}) \log(1-\sigma(z^{(i)}))]$$

LOGISTIC LOSS FUNCTION

$$L(\theta) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}; \theta) \rightarrow \text{LOG-LIKELIHOOD}$$

$$\log L(\theta) = \sum_{i=1}^n \log p(y^{(i)}|x^{(i)}; \theta)$$

1. APPLYING LOG FUNCTION REDUCES POTENTIAL NUMERICAL UNDERFLOW

- WORKING WITH PROBABILITY IN RANGE 0-1.

} $\log(p)$ MAKES VERY TINY PROBABILITIES MANAGEABLE + NUMERICALLY STABLE.

- ↳ SOMETIMES PROBABILITIES CAN GET REALLY CLOSE TO 0 OR 1.

→ IMPROVES RELIABILITY FOR MODEL TRAINING PROCESS

2. CONVERT PRODUCT OF FACTORS INTO SUMMATION OF FACTORS. → EASIER TO COMPUTE DERIVATIVES.

↳ e.g. $\log(a \cdot b \cdot c) = \log(a) + \log(b) + \log(c)$

↳ TAKING LOG OF PRODUCTS ALLOWS YOU TO WRITE THE LOSS AS A SUM, MAKING GRADIENT CALCULATION FOR WEIGHTS SIMPLER + STABLE.

FOUNDATION OF LOGISTIC REGRESSION → LIKELIHOOD FUNCTION.

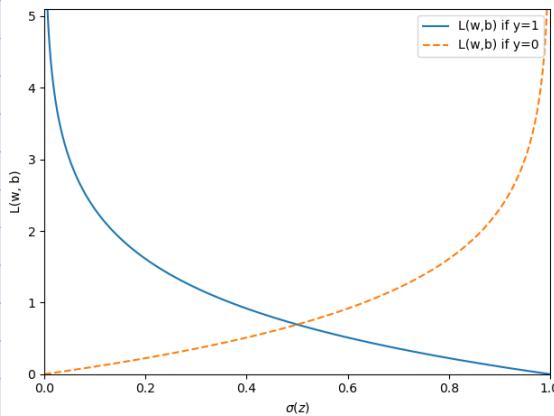
- ↳ GD ALGO → MAXIMISE THE LIKELIHOOD ⇒ DONE BY GRADIENT ASCENT.

↳ HOWEVER → CAN USE GRADIENT DESCENT WITH NEGATIVE LOG-LIKELIHOOD. → MINIMISING → CONVENIENT TO SWITCH TO NEGATIVE.

$$L(w, b) = \sum_{i=1}^n [-y^{(i)} \log(\sigma(z^{(i)})) + (1+y^{(i)}) \log(1-\sigma(z^{(i)}))]$$

$$L(\sigma(z), y; w, b) = -y \log(\sigma(z)) - (1-y) \log(1-\sigma(z))$$

$$L(\sigma(z), y; w, b) = -\log(\sigma(z)) \begin{cases} \text{if } y=1 \\ -\log(1-\sigma(z)) \text{ if } y=0 \end{cases}$$



MAIN POINT → WE PENAUSE WRONG PREDICTIONS WITH INCREASINGLY LARGER LOSS.

CONVERTING ADALINE IMPLEMENTATION INTO ALGORITHM FOR LOGISTIC REGRESSION

- ↳ COMPUTE LOSS OF CLASSIFYING PER EPOCH } SWITCH THIS TWO → BECAUSES LOGISTIC REGRESSION.

- ↳ SWAP LINEAR ACTIVATION WITH SIGMOID.

- ↳ FIT LOGISTIC REGRESSION MODEL → ONLY WORKS FOR BINARY CLASSIFICATION TASKS.

GRADIENT DESCENT LEARNING ALGORITHM FOR LOGISTIC REGRESSION

- ADALINE, WEIGHTS + BIASES DIDN'T CHANGE FOR LOGISTIC REGRESSION → GRADIENT DESCENT SIMILAR FOR LR.

- PARTIAL DERIVATIVE OF LOG-LIKELIHOOD WITH RESPECT TO j-th TERM:

$$\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w_j} \quad \text{WHERE } a(z) = \frac{1}{1+e^{-z}}$$

$$\frac{\partial L}{\partial a} = \frac{a-y}{a \cdot (1-a)}$$

$$\frac{\partial a}{\partial z} = \frac{e^z}{(1+e^{-z})^2} = a \cdot (1-a)$$

$$\frac{\partial z}{\partial w_j} = x_j$$

→ TAKE STEP OPPOSITE OF GRADIENT → FLIP $\frac{\partial L}{\partial w_j} = -(y-a)x_j$ + LEARNING RATE η

$$w_j := w_j + \eta(y-a)x_j$$

$$b := b + \eta(y-a)$$

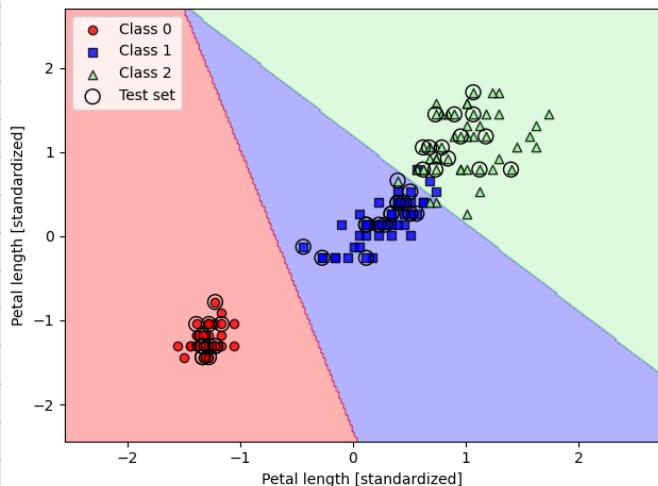
WEIGHT + BIAS = ADALINE WEIGHT + BIAS

$$\frac{\partial L}{\partial w_j} = (a-y)x_j$$

$$\frac{\partial L}{\partial w_j} = -(y-a)x_j$$

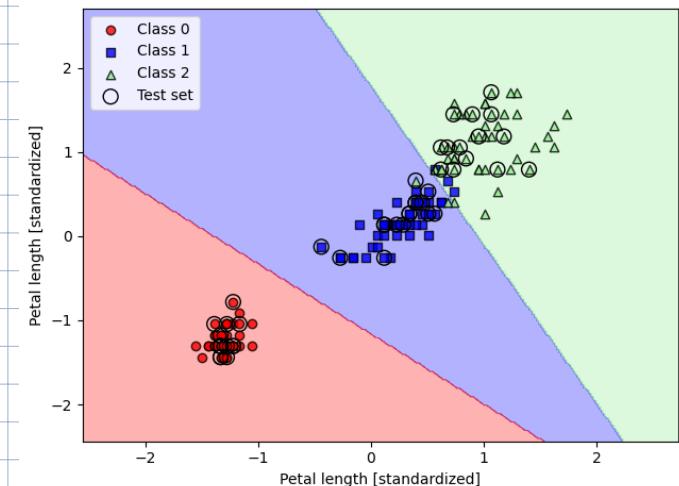
LOGISTIC REGRESSION WITH SCIKIT-LEARN

- SUPPORTS MULTI-CLASS LOGISTIC REGRESSION OFF THE SHELF



MULTI-CLASS = DLR

DEFAULT FOR SCIKIT-LEARN LOGISTIC REGRESSION

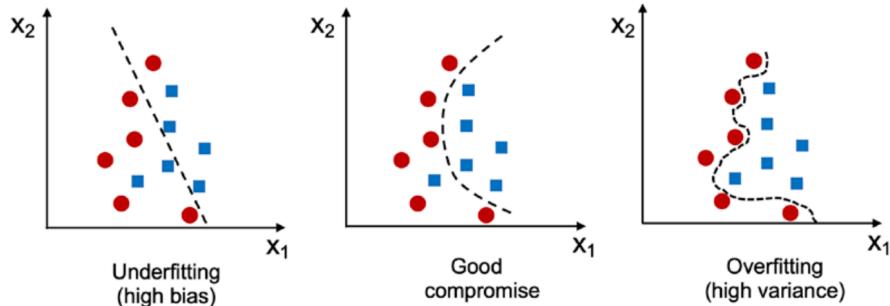


MULTI-CLASS = MULTINOMIAL

(CONVEX OPTIMISATION) → MINIMISING CONVEX LOSS FUNCTIONS (LOGISTIC REGRESSION) → RECOMMENDED TO USE MORE ADVANCED ALGO OVER SGD.

OVERFITTING VIA REGULARISATION

- COMMON PROBLEM FOR MACHINE LEARNING → DOESN'T GENERALISE WELL WITH UNSEEN DATA.
- OVERTFITTING → HIGH VARIANCE MODEL → CAUSED BY TOO MANY PARAMETERS (TOO COMPLEX)
- UNDERFITTING → HIGH BIAS → MODEL NOT COMPLEX ENOUGH



BIAS-VARIANCE TRADE-OFF

WHAT IS BIAS?

- MEASURES HOW FAR OFF THE PREDICTIONS ARE FROM THE CORRECT VALUES.
- HIGH-BIAS MODEL MAKES STRONG ASSUMPTIONS ABOUT DATA BUT UNDERFITS - CAN'T FIGURE OUT THE UNDERLYING PATTERN WELL.

WHAT IS VARIANCE?

- MEASURES CONSISTENCY OF MODEL PREDICTION
- HIGH-VARIANCE → MODEL IS SENSITIVE TO SPECIFIC TRAINING DATA (FITS NOISE + SMALL DATA FLUCTUATION)

THE TRADE-OFF

- GOTTA FIND THE SWEET SPOT BETWEEN THE TWO.
- ONE WAY TO FIND GOOD BIAS-VARIANCE TRADEOFF → TUNE COMPLEXITY OF MODEL VIA REGULARISATION.

REGULARISATION → VERY USEFUL FOR HANDLING COLLIGARITY, FILTERING OUT NOISE FROM DATA + EVENTUALLY PREVENTING OVERTFITTING.

↪ ADDITIONAL INFORMATION TO PENALISE EXTREME PARAMETER VALUES.

↪ MOST COMMON → L2 - REGULARISATION

REGULARISATION → FEATURE SCALING → STANDARDISATION → IMPORTANT.

↪ FOR THIS TO WORK → FEATURES MUST BE ON COMPARABLE SCALES.

$$\frac{\lambda}{2n} \|w\|^2 = \frac{\lambda}{2n} \sum_{j=1}^m w_j^2$$

λ = REGULARISATION PARAMETER

λ = SCALING FACTOR → CANCELS COMPUTING LOSS FUNCTION.

n = SAMPLE SIZE ADDED TO SCALE REGULARISATION TERM SIMILAR TO LOSS.

LOSS FUNCTION FOR LOGISTIC REGRESSION + REGULARISATION:

$$L(w, b) = \sum_{i=1}^n [-y^{(i)} \log(\sigma(z^{(i)})) + (1+y^{(i)}) \log(1-\sigma(z^{(i)}))] + \frac{\lambda}{2n} \|w\|^2$$

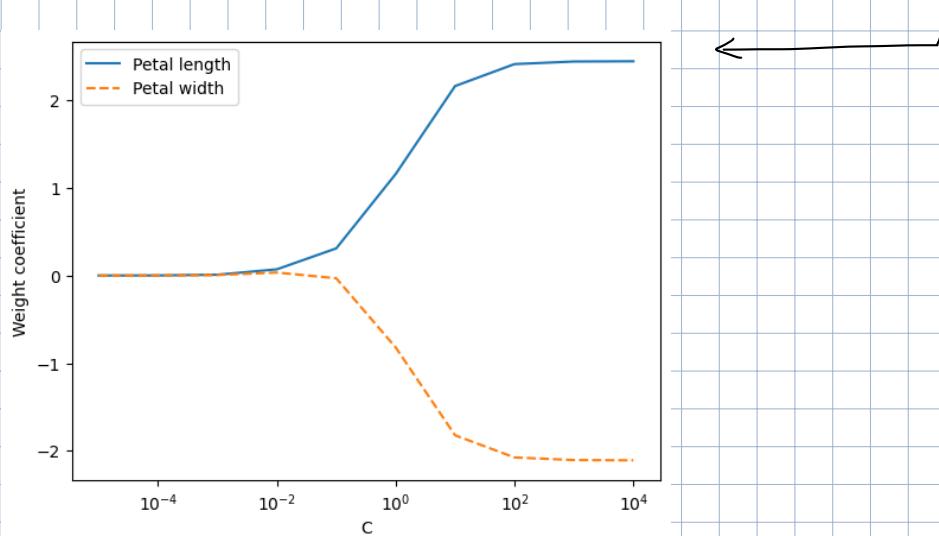
PARTIAL DERIVATIVE OF UNREGULARISED LOSS:

$$\frac{\partial L(w, b)}{\partial w_j} = \left(\frac{1}{n} \sum_{i=1}^n (\sigma(w^\top x^{(i)}) - y^{(i)}) x_j^{(i)} \right) \implies \frac{\partial L(w, b)}{\partial w_j} = \left(\frac{1}{n} \sum_{i=1}^n (\sigma(w^\top x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{n} w_j$$

CAN CONTROL HOW CLOSELY WE FIT THE DATA WITH $\lambda \rightarrow$ INCREASE λ , INCREASE REGULARISATION STRENGTH.

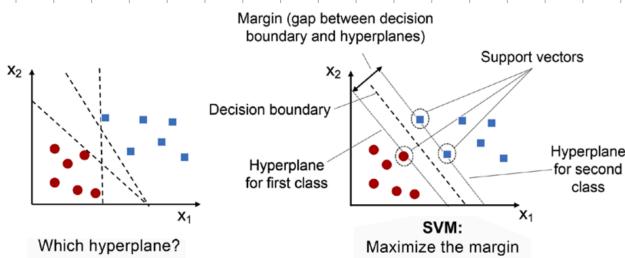
WHY DO WE NOT STRONGLY REGULARISE ALL MODELS IF IT REDUCES OVERFITTING?

\hookrightarrow IF REGULARISATION \rightarrow TOO HIGH \rightarrow WEIGHTS COEFFICIENTS APPROACH 0 \rightarrow THEN MODEL WILL UNDERFIT.



MAXIMUM MARGIN CLASSIFICATION WITH SUPPORT VECTOR MACHINES

- POWERFUL + WIDELY USED ALGORITHM \rightarrow SUPPORT VECTOR MACHINE (SVM) \rightarrow CONSIDERED EXTENSION OF PERCEPTRON.
- SVM - MAXIMISE MARGIN \rightarrow DISTANCE BETWEEN THE SEPARATING HYPERPLANE (DECISION BOUNDARY) AND TRAINING EXAMPLES CLOSEST TO HYPERPLANE



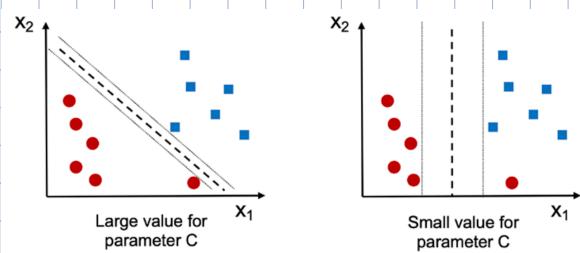
MAXIMUM MARGIN INTUITION

- LARGER MARGINS \rightarrow TEND TO HAVE A LOWER GENERALISATION ERROR.
- SVM LOOKS SIMPLE \rightarrow MATH IS COMPLEX \rightarrow NEED TO KNOW CONSTRAINED OPTIMISATION.

SUPPORT VECTORS

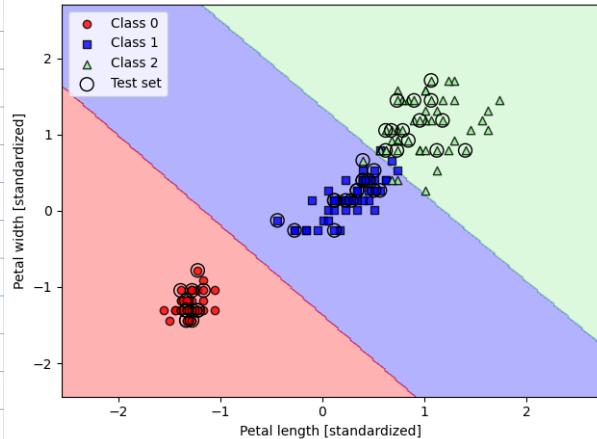
SLACK VARIABLE

- WOULD GO INTO MATHEMATICAL EXPLANATION \rightarrow LOOK INTO THIS + SVM.
- INTRODUCED \rightarrow VLADIMIR VAPNIK (1995) \rightarrow LED TO SOFT-MARGIN CLASSIFICATION
 - \hookrightarrow STRICT CONSTRAINTS \rightarrow FORCE SVM TO FIND A HARD BOUNDARY \rightarrow FAILS IF DATA ISN'T CLEANLY SEPARABLE.
 - \hookrightarrow RELAXING CONSTRAINTS \rightarrow ADDING SLACK VARIABLES \rightarrow SOME POINTS CAN BE MISCLASSIFIED \rightarrow BUT PENALTY FOR EACH MISTAKE.
 - \hookrightarrow SVM OPTIMISATION SEEKS BALANCE \rightarrow WIDE MARGIN (GOOD SEPARATION) + FEW MISCLASSIFICATION (LOW SLACK PENALTY)
- TO USE: INTRODUCE VARIABLE C \rightarrow HYPERPARAMETER FOR CONTROLLING PENALTY FOR MISCLASSIFICATION.



- LARGER C \Rightarrow LARGE ERROR PENALTIES
- SMALLER C \Rightarrow LESS STRICT ABOUT MISCLASSIFICATION ERROR
- C CONTROLS WIDTH OF MARGIN \rightarrow TUNES BIAS-VARIANCE TRADEOFF

THREE DECISION REGIONS OF SVM.



LOGISTIC REGRESSION VS. SVM

- LOGISTIC REGRESSION + SVM IN PRACTICAL CLASSIFICATION YIELD SIMILAR RESULTS.
- LOGISTIC REGRESSION → MAXIMISE CONDITIONAL LIKELIHOOD
- SVM → CARES ABOUT POINTS CLOSEST TO DECISION BOUNDARY (SUPPORT VECTORS)
- LOGISTIC REGRESSION → SIMPLER MODEL + EASY TO IMPLEMENT + EASILY UPDATABLE.