

~~Linear Regression Presentation~~

- Imagine you are trying to predict the  of a house based on its size....
- How do we find the best way to make this prediction?
- At its core:
 - LN helps us find relationship between an input & an output
 - It does this by fitting a straight line through data points!
 - This line is defined by a function $f(x)$
 - $f(x)$ represents our prediction model.

LARGER!

Neural Network

Input Layer

- Neurons in input layer directly output the input values, no activation function

Hidden Layers

- Must use non-linear activation function for a useful ANN
- Famous: Sigmoid, Tanh, ReLU

Output Layer

- Depends on what our ANN is used for:
 - Regression Problems \rightarrow Linear Function
 - Binary Classification \rightarrow Sigmoid Function $\rightarrow \frac{e^x}{1 + e^x}$
 - Multi-class classification \rightarrow Softmax Function $\rightarrow p(y=i|z) = \frac{e^z}{\sum_{j=1}^n e^z_j}$

Weights

- In a fully connected feed forward neural network, each neuron in one layer is connected to every neuron in the next layer. These connections are represented by weights, which determine the strength of the influence of one neuron on another.

Forward Pass

A sample (x, y)

- Present x at input layer
- Compute $z^{(l)}$ & $a^{(l)}$ for all layers
- Compute output of Loss function

1 Sum or Avg loss of all samples = total loss

TRAINING

- TRAIN PARAMETERS = WEIGHTS & BIASES

Process

- o Initialize W & b randomly for all layers
- o Choose a α
- o Repeat UNTIL CONVERGENCE or Stopping Condition
 - UPDATE PARAMETERS = $\frac{\partial J}{\partial \theta}$ OF COST FUNCTION WRT EACH W & b

HOW TO FIND?

$$W_{\text{new}} = W_{\text{old}} - \alpha \nabla J \quad b_{\text{new}} = b_{\text{old}} - \alpha \nabla J$$

BACKPROPAGATION

- BACKWARD PASS
- CALCULATE $\frac{\partial J}{\partial \theta}$ FROM OUTPUT LAYER TOWARDS INPUT LAYER
- COMBINE THEM USING LEARN RATE TO GET HOW LOSS CHANGES WRT TO EACH WEIGHT

PREVENT OVERFITTING

- o REGULARIZATION METHODS
 - WEIGHT DECAY = Squash HIGH MAGNITUDE W TOWARDS ZERO
 - DROPOUT REGULARIZATION = Leave some NEURONS UNUSED IN THE TRAINING
 - Use FULL NETWORK FOR PREDICTION THROUGH
 - EARLY STOP = WHEN TEST ERROR STOPS IMPROVING, STOP TRAINING
 - DATA AUGMENTATION = Distorts DATA in a HUMAN RECOGNISABLE WAY. SHIFT, ROTATE.

Convolutional Neural Networks

- Design with Fully Connected Neural Network MANDATORY Imposes us
 - ↳ TO FLATTEN I.E. IMAGE
 - ↳ LEADS TO LOSS OF SPATIAL INFORMATION
 - ↳ AND WE NEED A HUGE AMOUNT OF PARAMETERS \rightarrow MAY LEAD TO OVERFIT
- CNN WORKS BETTER FOR THIS

Process

- CHECK FOR CERTAIN FEATURES IN DIFFERENT IMAGE PATCHES
- CREATE HIERARCHY OF FEATURES
 - ① DETECT LOW LEVEL FEATURES
 - ② DETECT MID LEVEL FEATURES
 - ③ DETECT HIGH LEVEL FEATURES

ENHANCED PROCESS

- 1 CONVOLUTION LAYERS WHICH EXTRACT FEATURES FROM LOW-LEVEL TO HIGH-LEVEL
- 1 THEN FLATTEN OUTPUT & APPLY FULLY CONNECTED NN TO CLASSIFY
- A SO NETWORK LEARN NON-LINEAR COMBS OF HIGH LEVEL FEATURES

CONVENTIONAL LAYERS

- THE CORE OF CNN ARE FILTERING LAYERS & POOLING LAYERS
- WHERE FILTERING LAYERS:
 - o APPLY FILTERS TO INPUT DATA PERFORMING FEATURE EXTRACTION
- WHILE POOLING LAYERS:
 - o PERFORM DOWNSAMPLING OF THE DATA
 - MOST COMMON METHOD IS MAX-POOLING
 - CAN EMPHASIS DOMINANT FEATURE

Filters, Padding, Stride

- A FILTER IS A SMALLER MATRIX THAN INPUT MATRIX THAT SLIDES ON IT BY A CERTAIN STRIDE (LT OF STEPS). CREATING A CONVOLUTED FEATURE.
 - o MULTPLY OVERLAPPING VALUES & SUM THEM

THANK TO THIS MULTIPLE LAYER WE ARE THEN ABLE TO HAVE A BETTER CLASSIFICATION!

ENSEMBLE METHODS

- IN MACHINE LEARNING NO SINGLE MODEL CAN ALWAYS PERFORM BEST ON EVERY PROBLEM.
Some models overfit, others underfit. What if we could COMBINE MULTIPLE MODELS TO BALANCE OUT THEIR WEAKNESSES & AMPLIFY THEIR STRENGTHS?
- Now with Ensemble Methods we can!

two types

HOMOGENEOUS = ≠ TYPES OF CLASSIFIERS TRAINED ON THE SAME DATA

HETEROGENEOUS = SAME TYPES OF CLASSIFIERS TRAINED ON DIFFERENT DATA

Aim = ENHANCE DIVERSITY!

WHY ENSEMBLE?

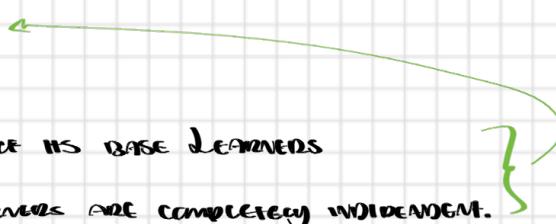
- CAN PERFORM BETTER THAN ANY OF ITS BASE LEARNERS
- Under condition that base learners are completely independent.

AIM OF ≠ GROUPS

- ONE GROUP AIMS TO REDUCE VARIANCE ANOTHER TO REDUCE BIAS

DIVERSE PREDICTORS

BOOTSTRAP



Denote Variance

BAGGING → Bootstrap Aggregation

- MAKE BOOTSTRAPPED SUBSET OF THE TRAINING DATA
- TRAIN A SEPARATE MODEL ON EACH SUBSET
- TRAIN CAN BE DONE IN PARALLEL

PREDICTION: COMBINATION OF PREDICTIONS OF ALL MODELS! AVG

CLASSIFICATION: USE VOTING!! HARD VOTING

SOFT & HARD VOTING

HARD: EACH CLASSIFIER MAKES A PREDICTION & MAJORITY VOTE WINS

SOFT: REQUIRES ESTIMATES $p(K|x)$, FINAL PRED IS CLASS WITH HIGHEST PROB

EVALUATION: OUT-OF-BAG EVALUATION

BAGGING = BOOTSTRAPPED DATA = SAMPLING WITH REPLACEMENT

THE NON SAMPLED POINTS WILL BE USED FOR VALIDATION

CAN BE REDUCED BY A VALIDATION SET!

ERROR = AVG ERROR BETWEEN MODELS THAT HAVE THOSE AS OOB

Random Forest:

- ESSENTIALLY A BAGGING DECISION TREE METHOD WITH MODIFIED SPLIT CRIT.
- MAKE BOOTSTRAPPED DATA SETS.
- TRAIN A FULL TREE ON THIS DATA SETS
- BEFORE FINDING EACH SPLIT:
 - CHOOSE RANDOMLY K FEATURES
 - CONSIDER ONLY THOSE FOR SPLIT
- OUT = AGGREGATION OF ALL TREES
- READING ME BECAUSE NOT MANY HYPERPLANES:
 - m = # BOOTSTRAPPED DATASET
 - K = # OF RANDOMLY SELECTED FEATURES.

* PICKS BEST CURRENT

Extra Trees - Extremely Randomized Trees

- FIT EACH DT ON WHOLE DATASET
- USE RANDOMLY SAMPLE K FEATURES
- RANDOMLY SELECTS A CUT-POINT FOR AN REGIONES OF PICKS BEST

Reduce Bias

BOOSTING:

SIGNIFICANTLY OUTPERFORMS BAGGING!

SEQUENCE OF WEAK LEARNERS COMBINED INTO A STRONGER ONE!

SEQUENTIAL !!!!!

How it works:

Draw a weak learner \rightarrow Get training error \rightarrow train next based on error

Prediction: Combination of all weak learners trained

Ada Boosting \rightarrow most commonly used for classification!

Process:

Equal weights for all training samples $w_i^{(1)} = \frac{1}{N}$

For $c=1, \dots, n$ iterations (trees):

- Train a model $h_c(x)$

- Calculate total error, on training set $e_c = \sum_{i=1}^N w_i^{(c)} h_c(x_i) / y_i$

- Based on calculate amount of say: $\alpha_c = \ln\left(\frac{1-e_c}{e_c}\right)$

- Update weights:

• Increase misclassification:

$$w_i^{(c+1)} = w_i^{(c)} e^{\alpha_c}$$

- Correct classifications:

Unchanged or decrease

$$w_i^{(c+1)} = w_i^{(c)}$$

model has to focus more on this!!

WHERE PREDICTION WAS DIFFERENT

- Normalize weights to prevent instability:

$$w_i^{(c+1)} = \frac{w_i^{(c)}}{\sum_{j=1}^n w_j^{(c)}}$$

Decision Trees

- Supervised data
- Great interpretability due to its IF-ELSE structure
- Each split focuses on a single feature
- Process:

Ask question recursively to create a tree

- trees can be used both for regression & classification

MEAN

PREDICTION: actual label

- Most commonly a split is binary but can also be multi-way

Training

- A balanced tree is a good tree

- To choose the node splitting condition:

- Try all possible splits for current node

- Use Heuristic to measure how good each split is
- Pick the best split!

EVALUATING A SPLIT

REGRESSION: Calculate mean \bar{y}_{xz} & resulting action y_{xz}

$$RSS = \sum_{i:x_i \in R_1} (y_i - \bar{y}_{xz})^2 + \sum_{i:x_i \in R_2} (y_i - \bar{y}_{xz})^2$$

Pick the split with the minimum RSS!

CLASSIFICATION: Concept of **Impurity**

A measure of how mixed data is w.r.t class members

Combine the impurity using their weighted avg

$$D_L\phi(L) + D_R\phi(R)$$

SPLIT THAT MINIMIZES THIS IS THE BEST SPLIT!

$$\phi(L) \quad \phi(R)$$

$$D_L = \frac{N_L}{N_u} \quad D_R = \frac{N_R}{N_u}$$

Impurity Function

- CAN BE DEFINED AS A FUNCTION OF POSTERIOR PROBS.

$$\phi(p_0, p_1, \dots, p_K) \quad \text{where estimate p.e. } p_k = \hat{p}(k|t) = \frac{N_k}{N_t}$$

- Some of the most used impurity functions:

Gini Impurity: $\phi(t_i) = 1 - \sum_{k=1}^K p_k^2 \quad \Rightarrow \text{Gini Index}$

Entropy: $\phi(t_i) = - \sum_{k=1}^K p_k \log(p_k)$

PCA

IN DATA SCIENCE DATASETS ARE ALWAY HIGH DIMENSIONAL

SO WE WISH TO REPRESENT DATA IN A USEFUL WAY

WE TRY TO SEE IF DATA CAN BE REDUCED TO LOWER DIMENSIONS!

- UNSUPERVISED METHOD FOR SUCH THINGS
- PCA FINDS A PATTERN THAT IS:

WHAT DIRECTIONS IN OUR DATASET SPREAD THE DATAPOINTS THE MOST?

SEEK THE DIRECTION OF MAXIMAL VARIANCE

Process

- TWO NEW BASIS VECTORS!

- BUILD AN ORTHOGONAL BASIS WHERE NEW BASIS VECTORS EXPLAIN THE DIRECTION OF GREATEST VARIANCE IN DATA.

- PROJECT TO LOWER DIMENSIONS!

- FIRST K PRINCIPAL COMPONENTS SPAN A K-DIMENSIONAL SUBSPACE

- IT IS THE BEST K-DIM VIEW OF DATA

PRINCIPAL COMPONENTS

FEATURE EXAMPLE

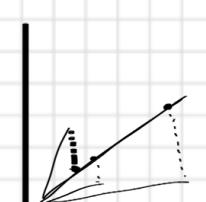
- Consider set of p features x_1, \dots, x_p
- PCA = GIVES NEW SET OF THE p FEATURES, PRINCIPAL COMPONENTS, EACH BEING A LINEAR COMB OF THE ORIGINAL p FEATURES!
- Every principal component is going to capture a \neq amount of variance

IMAGINE THAT WE HAVE n DATA POINTS, EACH WITH p FEATURES.

$$\text{DATA POINT } i \longrightarrow x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}, \text{ WE CAN DEFINE } \bar{x} = \sum_{i=1}^n x_i / n$$

GOAL: FIND DIRECTION OF MAXIMAL VARIANCE

LINEAR ALGEBRA PROBLEM

- In order to solve this we need to transform it into a Linear Algebra Problem
 - IMAGINE THAT x_i = ONE DATA POINT
vector of features
 - e_i = some vector
- } WHEN WE DOT TWO VECTORS, ONE GETS PROJECTED onto the other
- WE CAN e_i , a col vector, OUR FIRST PRINCIPAL COMPONENT
 - IF e_1 POINTS IN THE DIRECTION OF MAXIMAL VARIANCE
 our goal is then to maximize $\text{Var}(e_1^T x)$
- 
- Maximization turns into $\text{Var}(e_1^T x) = e_1^T S e_1$ Data Covariance Matrix
- $$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$
- Ensures e_1 is a unit vector
- We can maximize this to constraint $e_1^T e_1 = 1$ by Lagrange multipliers:
- $$F(e_1, \lambda) = e_1^T S e_1 - \lambda (e_1^T e_1 - 1) = 0$$

NOTE

SO WE, MAX WRT e_1

PRINCIPAL COMPONENT IS AN EIGENVECTOR OF THE DATA COVARIANCE MATRIX, REPRESENTING DIRECTION OF MAXIMAL VARIANCE IN DATASET

$$\frac{\partial F}{\partial e_1} = 2S e_1 - 2\lambda e_1$$

WHICH MEANS THAT e_1 IS THE EIGENVECTOR CORRESPONDING TO λ_1 EIGENVALUE

$$S e_1 - \lambda_1 e_1 = 0$$

$$S e_1 = \lambda_1 e_1$$

AND MAX WRT λ_1
WE REACH ∞ ...

$$\begin{aligned} F(e_1, \lambda_1) &= e_1^T S e_1 - \lambda_1 (e_1^T e_1 - 1) \\ &= \cancel{e_1^T} \lambda_1 e_1 - \cancel{\lambda_1} e_1^T e_1 + \lambda_1 \\ &= \lambda_1 \end{aligned}$$

WHAT NEEDS TO CHOOSE λ_1 TO BE LARGEST EIGENVALUE

TO FIND THE SECOND & OTHER COMPONENTS WE WILL ADD 0 FOR EACH NEW CONSTRAINT WE IMPOSE.

EIGENDECOMPOSITION

- In PCA we consider EIGENDECOMPOSITION OF S given by: $S = A \Lambda A^T$

- WHERE Λ IS THE DIAGONAL MATRIX WITH EIGENVALUES $\lambda_1 > \dots > \lambda_p > 0$

- AND WHERE $A = p \times p$ ORTHOGONAL MATRIX WHERE COLUMNS = CORRESPONDING EIGENVECTORS

ALL OF THIS TO SAY THAT PCA CHOOSES EIGENDecomposition THAT ORDERS THE EIGENVECTORS ACCORDING TO DECREASING EIGENVALUES

SO A PRINCIPAL COMPONENT CAN BE REPRESENTED BY:

$$e_k^T x = e_{k1}x_1 + e_{k2}x_2 + \dots + e_{kp}x_p$$

COEFFICIENTS ALSO CALLED LOADINGS

PCA CREATES UNCORRELATED FEATURES THAT ARE LINEAR COMBINATIONS OF EXISTING FEATURES.

RECONSTRUCTION OF X

BY USING ALL PC'S WE CAN RECONSTRUCT EXACTLY A DATA POINT x

BY USING ONLY THE FIRST m PRINCIPAL COMPONENTS WE GET AN APPROXIMATION

DATA CENTERING

HERE'S WE MINIMIZE APPROXIMATION ERROR

MEANS MOVING COORDINATE SYSTEM TO THE MEAN \bar{x}

CHOOSING # OF PC

- THERE ARE MANY HEURISTICS ON HOW TO CHOOSE A GOOD # OF THEM
- TREAT THEM AS TUNING PARAMETERS
- WE CAN USE IF HELP US WITH A SCREE PLOT EXPLAINING:
 - HOW MUCH VARIANCE IS EXPLAINED BY EACH PC
 - CUMULATIVE IS MORE INTUITIVE

STANDARDIZATION

- WE NEED TO STANDARDIZE OTHERWISE FEATURES WITH LARGER VARIANCES WOULD DOMINATE THE PRINCIPAL COMPONENTS BASED ON ITS VARIANCE.
- SO BY CENTERING & STANDARDIZING WE GET:

VARIANCES WITH MEAN 0 & VARIANCE 1

BAYES THEOREM

- THE Bayes Theorem ESTIMATES THE POSTERIOR PROB OF CLASSES GIVEN OBSERVED FEATURES

$$P(C_k | X) = \frac{P(X|C_k)P(C_k)}{P(X)}$$

- WITH THE AIM OF MAKING AS FEW MISCLASSIFICATION AS POSSIBLE
- IN ORDER TO DO THIS WE NEED TO DEFINE A DECISION RULE:

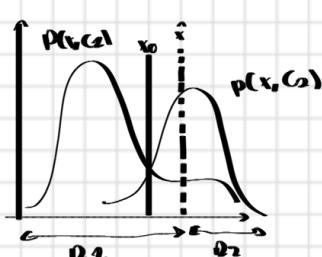
↳ DEFINES WHEN DO WE CLASSIFY A POINT AS C_k ?

↳ ANY DECISION RULE DIVIDES THE FEATURE SPACE INTO DECISION REGIONS:

$$R_k \quad k \in 1, \dots, K$$

- We want that the aim = AS FEW MISCLASSIFICATION AS POSSIBLE

$$\text{minimize} [p(\text{mistake})] = p(X \in R_1, C_2) + p(X \in R_2, C_1)$$



$$= \int_{R_1} p(X, C_2) dX + \int_{R_2} p(X, C_1) dX$$

$$\int_{R_1} p(x) p(C_2 | x) dx + \int_{R_2} p(x) p(C_1 | x) dx$$

LIMIT BY ASSIGNING CLASS WITH HIGHEST $p(C_k | x)$

WE KNOW THAT:

$$p(X, C_k) = P(X) p(C_k | X)$$

- Just developed classifier is called **Bayes classifier**

- WE CLASSIFY x WITH HIGHEST PROB:

$$d(x) = \arg \max P(Y=g|x)$$

- NOT ALL MISTAKES ARE EQUAL THOUGH!**

- DEATH WITH CANCER:

* FALSE POSITIVE causes stress to patient

* FALSE NEGATIVE may lead to death

GOAL = LIMIT CERTAIN MISTAKES

- Loss & Loss matrix:

$$\begin{matrix} \text{PRED} \\ \text{true } c \end{matrix} \begin{bmatrix} C & N \\ 0 & 1000 \\ N & 0 \end{bmatrix}$$

WE ASSIGN A LOSS = $L_{kj} > 0$ WHICH DEPENDS ON THE TRUE CLASS

- New way to express our goal

- Minimize THE EXPECTED LOSS!

$$E[L] = \sum_j \int_{R_j} L_{kj} p(X, C_k) dX$$

$$E[L] = \sum_j \int_{R_j} L_{kj} p(x) p(C_k | x) dX$$

} ASSIGN A NEW X TO CLASS j THAT MINIMIZES

$$\sum_k L_{kj} p(C_k | X)$$

SVM

- It is Supervised Learning Method

- Can be used both for:

Classification - Regression

- Intuition:

RESEARCHES THE BEST DECISION BOUNDARY FOR TRAINING DATAPoints

- IF C_k are linearly separable \rightarrow Best decision boundary is LINE

MEANING A HYPERPLANE WITH THE **MAXIMUM MARGIN**

(WIDEST STREET)

FROM THE CLOSEST SAMPLES OF BOTH CLASSES

- DECIDE LINEAR CLASSIFIER:

$$y(x) = w \cdot x + b \longrightarrow y(x) = 0$$

DECISION BOUNDARY

IF $y(x_i) \geq 0$ THEN SAMPLE $i=1$

IF $y(x_i) < 0$ THEN SAMPLE $i=-1$

- THE RELEVANT SAMPLES IN ORDER TO MAKE THE WIDEST STREET

• ARE THE ONLY SAMPLES THAT CONTRIBUTE TO DECISION BOUNDARY!

How to find Support Vectors?

- SVM FINDS THEM DURING TRAINING
- USES THEM TO CLASSIFY NEW SAMPLES!

Smart Vectors!!

HARD Margin

- USED WHEN TWO CLASSES ARE LINEARLY SEPARABLE

INSTEAD OF $y(x) = w \cdot x + b$:

$$y(x) = \sum_{i=1}^n d_i y_i (x_i \cdot x) + b$$

TRAINING SAMPLE
UNSEEN VALUE

WHEN $d_i > 0$
Support Vector!

DOT PRODUCT → every to sample x

- AFTER TRAINING WE CAN KEEP ONLY SUPPORT VECTORS AND THEIR ASSOCIATED d_i 'S FOR PREDICTIONS

DISCARD REST OF THE TRAINING SET!

- RULE: NO TO SAMPLES ARE ALLOWED ON THE WRONG SIDE, NOT EVEN ON THE MARGIN!

- TO MAXIMIZE WIDTH f OF PREDICTED

① DRAW MODEL, FIND d_i FOR TRAINING SAMPLE

$$y(x) = \sum_{i=1}^n d_i y_i (x_i \cdot x) + b$$

② KEEP ONLY SUPPORT VECTORS & THEIR d_i

WHERE $d_i > 0$

③ USE SUPPORT VECTORS & x_i TO CLASSIFY NE SAMPLES

VS

SOFT Margin

WHAT HAPPENS IF CLASS CONDITIONAL DISTRIBUTION OVERLAP?



- IT IS IMPOSSIBLE TO FIT A HARD MARGIN SVM FOR THIS
- WE RELAX THE CONSTRAINT
- SOME TO SAMPLE CAN FALL ON WRONG SIDE
- WE PENALIZE IT!
- OBJECTIVE:

MAXIMIZE STREET'S WIDTH WHILE MINIMIZING THE PENALTY

- END UP WITH THE FOLLOWING:

$$y(x) = \sum_{i=1}^n d_i y_i (x_i \cdot x) + b \quad (0 \leq d_i \leq C)$$

BECOMES TO A REGULARIZATION PARAMETER

- C MAKES A TRADEOFF BETWEEN PENALTY AND MARGIN MAXIMIZATION!

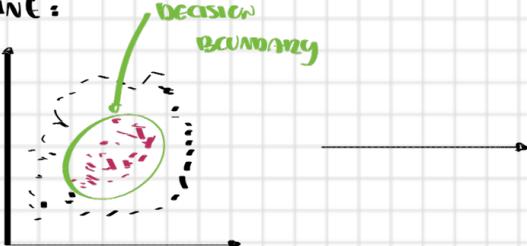
What makes SVM interesting?

- IT'S THE DOT PRODUCT
- LET'S THINK ABOUT NON-LINEAR DECISION BOUNDARY

IMPORTANT

SOMETIMES WE CAN MAP FEATURES FROM THE ORIGINAL DIMENSION SPACE TO A HIGHER DIMENSION WHERE THE CLASSIFICATION PROBLEM IS LINEAR

- IMAGINE:



A LINEAR DECISION BOUNDARY IN THE HIGH DIMENSIONAL FEATURE SPACE CAN REPRESENT A NON-LINEAR BOUNDARY IN THE ORIGINAL SPACE

- TRANSFORMATIONS CAN ADD MANY NEW NON-LINEAR FEATURES TO HOPEFULLY MAKE A BETTER CLASSIFIER!

↳ IT IS WHAT GIVES US THE POSSIBILITY TO MAP FROM A LOWER TO HIGHER DIM.

Hence, PREDICTION MODEL BECOMES

$$y(x) = \sum_{i=1}^n w_i y_i (\phi(x_i) \cdot \phi(x)) + b$$

- WE UNFORTUNATELY KNOW THAT:

HIGHER DIMENSION = HIGHER COMPUTATIONAL POWER REQUIRED

KERNEL TRICK!

- DEFINES A KERNEL FUNCTION THAT COMPUTES THE DOT PRODUCT OF THE TRANSFORMED INPUT WITH A TIME COMPLEXITY BASED ON THE ORIGINAL SPACE

- So now Prediction model:

$$y(x) = \sum_{i=1}^n w_i y_i K(x_i, x) + b$$

- THE BEST THING IS THAT:

We can transform feature vectors even to infinite dims & still not bother about needing more computational power

- WE COULDN'T NOT KNOW $\phi(x)$

WE ONLY NEED $K(x_i, x)$

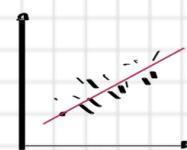
LINEAR REGRESSION

- IMAGINE YOU ARE TRYING TO PREDICT THE \$ OF A 🏠 BASED ON ITS SIZE....

- HOW DO WE FIND THE BEST WAY TO MAKE THIS PREDICTION?

- At its core:

- LN HELPS US FIND RELATIONSHIP BETWEEN AN INPUT & AN OUTPUT
- IT DOES THIS BY FITTING A STRAIGHT LINE THROUGH DATA POINTS!
- THIS LINE IS DEFINED BY A FUNCTION $f(x)$
- $f(x)$ REPRESENTS OUR PREDICTION MODEL.



↳ WE ASSUME A FUNCTIONAL RELATIONSHIP BETWEEN X & Y $\rightarrow y = f(x) + \epsilon$

↳ IN LINEAR REGRESSION WE TAKE $f(x)$ TO BE LINEAR $\rightarrow y = \beta_0 + \beta_1 x + \epsilon$

↳ MORE GENERALLY WITH p FEATURES & AN INTERCEPT TERM THE MODEL $\rightarrow y = \mathbf{x}\beta + \epsilon$

DESIGN MATRIX!

DUMMY VARIABLES

VARIABLES ARE NOT ALWAYS CONTINUOUS:

- YES/NO 0 or 1

WE THEN INTRODUCE DUMMY VARIABLES

WHERE WE INTRODUCE ONE DUMMY VAR FOR EACH LEVEL A OF FEATURE X

$$\mathbb{1}_{\{x=A\}} = \begin{cases} 1 & x=A \\ 0 & x \neq A \end{cases}$$

ONE HOT ENCODING

Complex Functional Relationships

- Variables can also interact with each other
- To handle this we need to account for interaction terms defined as x_1x_2 i.e.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \epsilon \quad \text{NON-LINEAR!!}$$

trick:

introduce transformation of feature were we get

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

NOTE:

DECISION IS CONVEX
IF PLOTTED AGAINST
ORIGINAL FEATURES
BUT LINEAR IN NEW
ONES.

Estimate Coefficients β

- Ordinary least squares regression finds β that minimizes the RSS

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y})^2$$

- More generally we prefer to fit our model to maximize the model likelihood

- MEANS: which model params would be most likely to give us the data?

- Start by defining a likelihood function

$$p_y(y_i | x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \hat{y})^2}{2\sigma^2}}$$

GIVEN THAT WE ASSUME RESIDUES TO BE N DISTRIBUTED WE CAN USE PDF OF IT

PDF FOR EACH ERROR

The likelihood for the entire dataset is then:

$$\prod_{i=1}^n p_y(y_i | x) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \hat{y})^2}{2\sigma^2}}$$

$$L(\beta, \sigma^2) = -\frac{1}{2} \left(n \log \sigma^2 + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \hat{y})^2 \right)$$

Prediction

- Once model is ready:

Given a set of features $x_i = (x_{i1}, \dots, x_{ip})$

We can predict \hat{y}_i

- We predict by saying that:

\hat{y}_i = value of regression line at x_i

Confidence & Prediction Intervals

- We can interpret coefficients using confidence intervals
- Prediction interval > Confidence Interval

Don't say, word to overtime

SCRUTINIZING Model!

- Now that we have our model we should ask ourselves:

DOES THE MODEL FIT WELL? \longrightarrow EFFECTIVELY SCRUTINIZE MODEL

- i.e. test whether a specific β is 0:

$$\frac{\beta_j - 0}{\hat{SE}(\beta_j)}$$

EVALUATE
SIGNIFICANCE
OF PREDICTOR
VARIABLE!

i.e. # of post boxes

- See if several β 's are β_0 :

CREATING a SIMPLER COPY OF THE MODEL

$$M_2: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \epsilon$$

$$M_0: Y = \beta_0 + \beta_1 x_1 + \epsilon$$

use F-TEST:

measures improvement in RSS in the two models

$$F = \frac{(RSS_{M_0} - RSS_{M_2}) / q}{RSS_{M_2} / (n - p - 1)}$$

of params dropped

total # of observations

High: Significant Reduction in M_2 's RSS

Low: It is not worth to have such complex model

Machine Learning Project!

In the Project we explore three different:

- Classification Methods
- To classify clothing Images!

The three methods we used are:

- Decision Trees
- FFNN
- KNN and Bagged KNN \longrightarrow IMPLEMENTED WITH LIBRARY SUPPORT

Exploratory Analysis

- At the beginning we did some exploratory analysis.
- Plotting the features distributions to immediately gain some valuable informations.
 - We were able to see that each class was **equally represented**
 - We knew that model would have had an equal representation during learning!
- Given the high dimensionality of the data we scaled the features & applied PCA in order to reduce as much as possible the dimensions!
- A scree plot helps us understand better how many PCs we needed in order

To capture at least 90% of the total variance

- We need at least 112 PCs

Decision Trees

- In order to evaluate the splits in a tree there is the need of a impurity function

As our implementation choice, our model works based on: **Information Gain**

So we choose Entropy over Gini Index

- Trains the model scaled & reduces features.
- Information Gain criterion was used mainly to better handle **CATEGORICAL** outputs
 - Entropy provides a more **intuitive** measure of disorder!
- We adopted three stopping criterions:

MAXIMUM DEPTH: CAPPED AT MAX_DEPTH VAR

MINIMUM SAMPLES SPLIT: ONLY SPLIT IF CONTAIN AT LEAST A MIN # OF SAMPLES [5]

HOMOGENEOUS LABELS: DECLARE LEAF IFF ALL SAMPLES HAVE SAME LABEL

- By a GRID SEARCH we were able to find best values for our HyperParams:

MIN_SAMPLES_SPLIT = 2

MAX_DEPTH = 10

PCA = 40 components → 80% variance explain

- Ended up with 97% F1 score for **SAME** confusion matrices

Feed Forward Neural Networks

Key elements of our implementation are:

Early Stopping: to prevent overfitting, at most 20 epochs with no loss improve?

Dropout Regularization: only on Keras model

Model Evaluation Metrics: Accuracy, per class loss change, confusion matrix, precision, recall & F1

- AFTER NUMEROUS EXPERIMENTS WITH THE MODEL WE DECIDED TO USE A SINGLE HIDDEN LAYER WITH 64 NEURONS & $\text{ELU} = \text{SIGMOID}$ FUNCTION
- WEIGHTS INITIALIZED WITH **GUARDED UNIFORM DISTRIBUTION**
- BIASES INITIALIZED TO 0
- η SET TO 0.0001
- EPOCHS TO 1000
- WE SAW THAT MODEL CONSTANTLY IMPROVED DURING SUCH HIGH # OF EPOCHS
- GIVEN HIGH # OF EPOCHS → EARLY STOPPING SET TO [20] **AFTER** [20] **> CRASH**
- RESULTS OF LIBRARY & SCRATCH WERE VERY SIMILAR 83%-86% ACC
- BOTH MODELS STRUGGLED IN CLASSIFYING SHIRTS, SOLVABLE BY **Data Augmentation**
- 'Scratch' → 95% acc., a bit of overfitting
- By comparing Confusion matrix can see the similarity!

Scratch
model uses ReLU activation
Keras
model uses SIGMOID activation

while its neurons
for the 'scratch'

KNN

- I initially implemented a version of KNN trained on non-reduced features
 - USED AS BASELINE & AS COMPARISON MODEL!
 - THE MODEL CHOSE WAS MADE WITH $K=7$ DUE TO HIGH ERROR IN RECOGNITION TO CLASS 4, T-SHIRT,
 - TOOK THE MODEL 301 SECONDS TO COMPUTE FARMS FOR CROSS VALIDATION.
- THEN STARTED WORKING WITH PCA-REDUCED FEATURES
 - SIGNIFICANTLY FASTER, TOOK ONLY 3.6s TO COMPUTE FARMS
 - W/BETWEEN A RANGE OF $K[4,8]$ TRAINED MODELS, CHOICE WITH $K=7$ REPORTED HIGHER ACCURACY & BETTER METRICS OF PRECISION, RECALL, F1 FOR MOST CLASSES

Classifying class 4

WHY WE PICKED IT

BAGGED KNN

USING IT TO REDUCE VARIANCE HENCE OVERFITTING

- MADE SURE THAT TRADEOFF BETWEEN VARIANCE & BIAS WAS TAKEN INTO ACCOUNT
 - SEEMED THAT A GOOD TRADEOFF WAS GIVEN BY MODELS WITH $K=6$ OR $K=5$
- WE GOT GREATER ACC FOR LOWER K'S AS EXPECTED
BAGGING & VARIANCE = LOWER K'S HAVE LOWER VARIANCE

I THEN BEGAN ANALYZING THE MATRICES THAT MODELS WITH $K \leq 5$ ARE TO ANALYZE

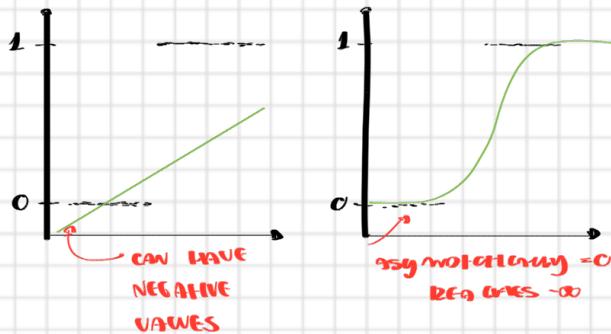
TESTING!

SO TESTING NORMAL KNN PCA-ADDED GIVES US 82.5% ACC!

WHITE OUT OF ALL MODELS WITH $K \leq 5$, $K=3$ IS THE BEST WITH ACC 83.3%

LOGISTIC REGRESSION

- IN MANY APPLICATIONS WE ARE INTERESTED IN PREDICTING WHAT CLASS A DATAPPOINT BELONGS TO:
- PREDICT BINARY OUTCOMES \rightarrow SOLVE A CLASSIFICATION PROBLEM
- LET'S IMAGINE WE WANT TO CLASSIFY A CERTAIN POINT AS 0 OR 1 BASED ON SOME FEATURES:



ASYMPTOTICALLY = 1
FOR ∞

$$P(x) = \frac{e^{Bx+bx}}{1+e^{Bx+bx}}$$

DEFINITION $P(x)$
W/BETWEEN
 $0-1$

- SOMETIMES WE LIKE TO THINK ABOUT LIKELIHOODS NOT IN TERM OF $P(x)$ BUT AS ODDS:

$$\text{ODDS} = \frac{P(x)}{1-P(x)} = e^{Bx+bx}$$

WHAT ARE THE ODDS THAT THIS HAPPENS RATHER THAN THE OTHER?

THEY ARE POSITIVE REAL IF $\text{ODDS}(A) > \text{ODDS}(B)$ IFF $P(A) > P(B)$

So if we take the log $\rightarrow \log(\text{odds}) = \beta_0 + \beta_1 x_1$ A one unit change
is a change of β_1 w/ logit

Multiclass Features!

- Generalizing to general feature:

- Logistic Regression models conditional prob of $y=1$ given x_1, \dots, x_p

$$P(y_i=1 | x_i=x) = \frac{e^{x\beta}}{1+e^{x\beta}}$$

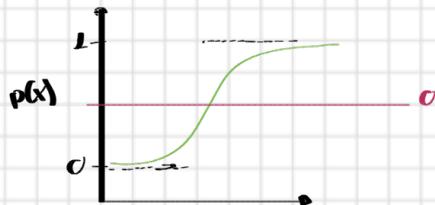
$$\text{odds} = \frac{P(y_i=1 | x_i=x)}{1-P(y_i=1 | x_i=x)} = e^{x\beta}$$

$$\log(\text{odds}) = x\beta$$

Prediction

- to make \hat{y} to the class with highest probability

- take into account a threshold $\rightarrow 0.5$ default



$$\hat{y}_i = \begin{cases} 1, & P(y_i=1 | x_i) > 0.5 \\ 0, & P(y_i=1 | x_i) \leq 0.5 \end{cases}$$

Evaluation

$$\begin{matrix} \hat{y}=1 & \hat{y}=0 \\ y=1 & \begin{bmatrix} 50 & 10 \\ 10 & 30 \end{bmatrix} \\ y=0 & \end{matrix}$$

Accuracy, Precision, Recall

Test Error Rate - prob of $\hat{y}_i \neq y_i$:

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

Estimating Coefficients

- Parameters are estimated by maximum likelihood
- Linear Function \rightarrow Binomial Linearization

$$L(\beta, y) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} 1-p(x_j)$$

$$\log(L(\beta, y)) = \sum_{i=1}^n (y_i \log(p(x_i)) + (1-y_i) \log(1-p(x_i)))$$

Multiclass Classification

- Softmax Regression

$y_i \in \{1, \dots, k\}$ Use one-Hot-Encoding $y_i = (0, 1, \dots, 0)$

clearly follows a multinomial distribution with (p_{i1}, \dots, p_{ik})

So now we have $k-1$ log-odds

$$P(y=k | x=x) = \frac{e^{x\beta_k}}{\sum_{c=1}^k e^{x\beta_c}}$$

then chooses \hat{y} that belongs to class with highest prob

$$\hat{y} = \text{argmax } \hat{P}(y=k | X=x)$$

Tradeoff Bias & Variance

- When we talk about prediction models,

AIM = Predict UNSEEN datapoints **well!!**

- Recall LINEAR CLASSIFIER:

$$y = f(x) + \epsilon$$

IMAGINE THAT WE HAVE TWO MODELS TO FIT THE DATA

$$\hat{f}_1(x) \quad \& \quad \hat{f}_2(x)$$

AND WE ARE NOW ASKED TO PREDICT 1 UNSEEN DATA POINT (x_0, y_0)

How do we evaluate which model is the best?

+NOTE *

LATER WE ARE GOING TO BE ASKED TO PREDICT n UNSEEN (x_i, y_i) $i=1 \dots n$

- With access to TRAINING DATA we want:

$$E[\text{test MSE}] = E[(y_0 - \hat{f}(x_0))^2] \longrightarrow \text{as small as possible!}$$

- IMAGINE TRAINING A LEARNING \hat{f} & PREDICTING A POINT (x_0, y_0)

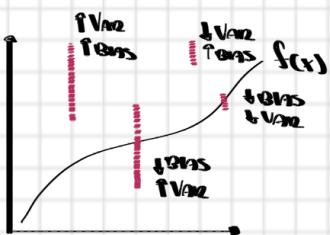
$$E[\text{test MSE}] = E(\hat{f}(x_0) - E\hat{f}(x_0))^2 + (E\hat{f}(x_0) - f(x_0))^2 + \text{Var}(\epsilon)$$

Variance
Swallow of data

Bias

Error introduced by approximation
& real life problem

RECALL:



= Errors made by \hat{f} instances of f

GIVEN THAT ALL THREE ARE NON-NEGATIVE
IF ANY IS LARGE, MSE IS LARGE
 \hat{f} LOW BIAS!!!
& LOW VARIANCE!!!

Decompose Formula:

$$E[\text{test MSE}] = E(\hat{f}(x_0) - E\hat{f}(x_0))^2 + (E\hat{f}(x_0) - f(x_0))^2 + \text{Var}(\epsilon)$$

REDUCIBLE ERROR

IRREDUCIBLE ERROR

Lower it by using estimate \hat{f} with both:

Low Bias
Low Variance

Completely

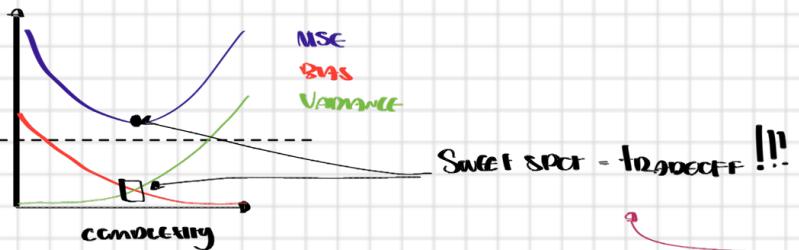
INTERPOLATES DATA
FITS A STRAIGHT LINE

} THINK ABOUT THESE TERMS OF VARIANCE

Has low bias, has high variance

Has high bias, low variance

DIFFICULT TO
FIND ONE
THAT HAS
BOTH LOW!!!



Error Formula

- to be a great model fit, we wanna do well with all points not just one!

We want $E[\text{Error}_0] = \text{Good}$ across all possible values of x_0 so...

So we end up with:

$$\text{MSE} = E[\text{Var}(f(x_0))] + E[\text{Bias}(f(x_0))]^2 + \text{Var}(\epsilon)$$

Clustering

- Clustering analysis is all about **discovering grouping in the data**

- unsupervised methods

- Given 'unlabelled data' \rightarrow Group data into some K of clusters

- It is vital the choice of K

- We seek to partition data into K clusters C_1, \dots, C_K

Assuming that each data point can only belong to a single cluster

- We say that datapoints in the same cluster are **similar**

- 2 methods covered:

K-means clustering
Hierarchical clustering

K-means Clustering

- First choose K

- Also then partitions data into K clusters C_1, \dots, C_K

WHERE: An datapoint belongs to exactly one cluster & all are non-empty.

We seek a partition that minimizes **within-cluster variation**

$$\underset{C_1, \dots, C_K}{\text{minimize}} \sum_{k=1}^K w(C_k) \quad \text{SIMILARITY} = \text{DISTANCE}$$

INSTEAD
USE CENTROIDS!

$$w(C_k) = \frac{1}{|C_k|} \sum_{i \in C_k} \|x_i - v_k\|^2$$

$$w_{C_k}(v_k) = \sum_{i \in C_k} \|x_i - v_k\|^2.$$

IF WE KNOW WHICH $i \in C_k$
THEN v_k = CLUSTER MEAN

$$v_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

Steps

- ① Randomly initialize v_1, \dots, v_K

- ② Assign all data points to closest v_k

- ③ Calculate mean & move centroid

} REPEAT UNTIL CENTRES STOP MOVING

THE ALGO MAY
NOT FIND AN
OPTIMAL PARTITION
OF DATA

CAREFUL

- Choice of K heavily impact which data points are clustered together
- Anything that impacts distances can impact clustering

High Dimensionality
Outliers
Scale of Variables

SILHOUETTE

- We can use Silhouette score to evaluate a partition

We calculate $a(x_i)$ & $b(x_i)$

$$\left\{ \begin{array}{l} a(x_i) = \frac{1}{|C_{\text{all}}|-1} \sum_{x_j \in C_{\text{all}}, j \neq i} d(x_i, x_j) \quad \text{AVG DISTANCE BETWEEN } x_i \text{ & ALL OTHER PANTS IN THE PARTITION} \\ b(x_i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{x_l \in C_j} d(x_i, x_l) \quad \text{AVG DISTANCE WITHIN } x_i \text{ & ALL PANTS IN CLOSEST CLUSTER} \\ S(x_i) = \frac{b(x_i) - a(x_i)}{\max(a(x_i), b(x_i))} \quad -1 \leq S(x_i) \leq 1 \quad \text{HIGH IS BETTER} \end{array} \right.$$

$\forall x_i \in C_k$

HIERARCHICAL METHODS

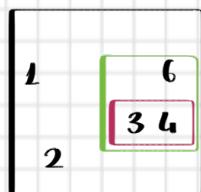
- AGGREGATIVE VS DISITIVE STRATEGIES

- Agglomerative = From $K \rightarrow K-1$ cluster
- Disitive = From $K \rightarrow K+1$ clusters

- It is all based on DISSIMILARITY

- . THE HIGHER THE BETTER

AGGREGATIVE

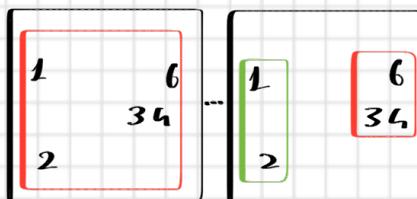


Start with every obs being a separate cluster

Repetitively merge two cluster

We merge the ones with smallest dissimilarity

DISITIVE



Works the opposite direction

Start with all obs in one cluster

Recursively split one cluster

MEASURE DISSIMILARITY BETWEEN CLUSTERS

- DISSIMILARITY BASED ON DISTANCE

SINGLE-LINK CLUSTERING

$$D(C_i, C_j) = \min \{d(x_i, y_j) | x_i \in C_i, y_j \in C_j\}$$

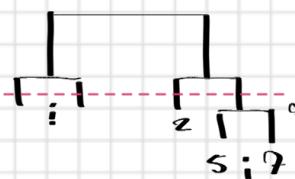
↳ shortest d
within any pair of point

Complete-link clustering $D(C_i, C_j) = \max \{d(x,y) | x \in C_i, y \in C_j\}$ → longest & minm any pair of pts

Group-average clustering $D(C_i, C_j) = \text{mean} \{d(x,y) | x \in C_i, y \in C_j\}$ → avg d within an pair of pairs

DENDROGRAMS

- Hierarchical clustering goes vs dendograms



Moving Bottom Up Clusters FUSE

While moving Bottom Down Clusters Get DIVIDED

WE CAN USE HORIZONTAL CUTS TO PARTITION DATA

All nodes below intersection are grouped together

Performance Metric

- Up until now we've understood that the common goal in prediction models is **MINIMIZING THE MISSCLASSIFICATION RATE**
- When minimizing Loss, we emphasized the fact that not all errors are equal
- Given that \neq task we require \neq performance metrics:

 - WE HAVE MANY OF THEM!
 - Metrics are usually evaluated on a test set!

CONFUSION MATRIX, PRECISION, RECALL

- At the \heartsuit of many performance matrix there is the confusion matrix.

	$\hat{y}=1$	$\hat{y}=0$
$y=1$	50	10
$y=0$	10	30

Accuracy = proportion of correct classification

$$(50+30)/100 = 0.8 \text{ acc}$$

Missclassification rate = proportion of missclassification

$$(10+10)/100 = 0.2 \text{ mr}$$

- Whenever we deal with binary matrix:

WE CAN EASILY WRITE THE FOLLOWING SCHEME

	$\hat{y}=1$	$\hat{y}=0$
$y=1$	TP	FN
$y=0$	FP	TN

We can use matrix entries to construct a range of informative measures!

Precision = % of true positives out of predicted pos.

$$TP / (TP+FP) \rightarrow 50 / (50+10)$$

How precise is our model?

Recall = % of true positives out of positive cases

$$TP / (TP+FN) \rightarrow 50 / (50+10)$$

How good is it at detecting cases?

- Tasks have \neq aims:

HIGH PRECISION

(SICKNESS)

i.e. it is not important to classify an C as C, but when do we need to **RIGHT!!**

i.e. it is important that we classify an C as C, some missclassifications are ok! (Cancer)

HIGH RECALL

We tend to train a classifier that has both high!

- It is very important to layout here:

Our choices can impact the confusion matrix! i.e. (Logistic Regression threshold)

F₁-Score

- We said that we want both Precision & Recall to be high

We can introduce F₁-score

- It is the Harmonic mean between precision & recall!

- Given that harmonic mean punishes small low numbers more than arithmetic does
to have high F₁-score → need high precision & recall

$$F_1 = 2 \frac{P \cdot R}{P+R}$$

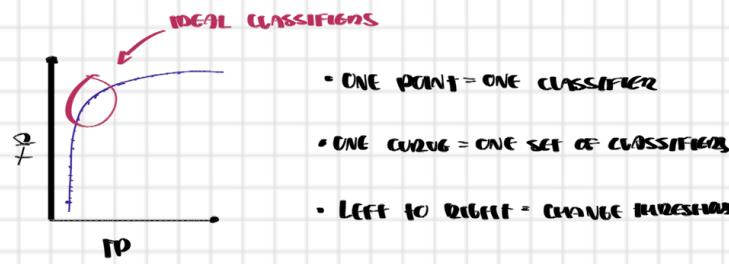
- For K classes there are K sets of precision & recall so also F₁-scores

We can get macro-F₁ score → avg across classes

$$\frac{1}{K} \sum_{k=1}^K F_1 k$$

ROC Curve

- Used for binary classification
- Plots $\frac{TP \text{ rate}}{\text{FP rate}}$ against $\frac{FP}{TP + FN}$



Classifier Good?

- Is model good?
- Do I want a simpler one?
- Best-performing model can still be bad!

Gradient Descent

- It is an optimization algorithm
- means min or max an objective function wrt some parameters.

We deal with:

Analytical methods & Numerical methods

Find a closed form sol
exact solutions

Approximated solutions
will be faster

Gradient Descent

- Most commonly used to find the $\min(f(x))$

Iterative

- Set a small Hyperparameter δ

- Randomly initialize parameters $\theta^{(0)}$
- Repeat until convergence \rightarrow TAKE A STEP AT A TIME DOWN THE GRADIENT

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla f$$

We stop once we've reached the minimum!



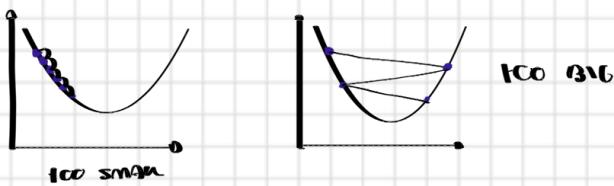
NOTE

∇f points in the direction of greatest ascent. To move down the gradient, we use its negative.

$$\begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_n} \end{bmatrix}$$

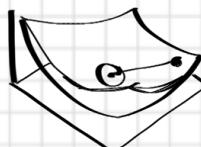
Vector of partial derivatives of a function w.r.t its params

SENSITIVITY



CONVEX & DIFFERENTIABLE

- Basic GD methods find the optimum minimum IF function is:
DIFFERENTIABLE & CONVEX



LINEAR REGRESSION APPLICATION

- Use GD to train machine learning models
- WHERE $f(\theta)$ is the cost/loss function f : θ = vector of model parameters
- TAKE i.e. Linear Regression WITH JUST ONE PARAMETER

$$\hat{y} = \beta_0 + \beta_1 x \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$\text{OPTIMIZATION: } \min_{\beta} J(\beta)$$

$$J(\beta) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \rightarrow \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i - y_i)^2$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$\frac{\partial J(\beta)}{\partial \beta_0} = \frac{2}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i - y_i) = \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \xrightarrow{\text{A}} \begin{bmatrix} \text{A} \\ \text{B} \end{bmatrix} = \nabla f$$

$$\frac{\partial J(\beta)}{\partial \beta_1} = \frac{2}{n} \sum_{i=1}^n x_i (\beta_0 + \beta_1 x_i - y_i) = \frac{2}{n} \sum_{i=1}^n x_i (\hat{y}_i - y_i) \xrightarrow{\text{B}} \begin{bmatrix} \text{A} \\ \text{B} \end{bmatrix} = \nabla f$$

$$\beta_0^{(t+1)} = \beta_0^{(t)} - \alpha \cdot \text{A}$$

$$\beta_1^{(t+1)} = \beta_1^{(t)} - \alpha \cdot \text{B}$$

UPDATE SIMULTANEOUSLY!!

DIFFERENT TYPES

BATCH GD \rightarrow An entire dataset is used to compute GD

Stochastic GD \rightarrow Considers just 1 datapoint at the time \rightarrow Good to escape local minima

Mini-Batch GD \rightarrow Computes ∇f on small random sets (mini-batches) \rightarrow May escape local minima

LDA & QDA

- Up until Generative Models we've only seen Discriminative ones
 - Both model do the same thing $\rightarrow P(Y|X)$
 - Just discriminative model them directly from the data i.e. Logistic Regression
 - While Generative ones take a step back & model the joint distribution of features

$$P(y|x) = \frac{P(x,y)}{P(x)}$$

Joint Distribution

- The joint distribution of features X & class Y can be specified in two parts:

$$P(x,y) = \underbrace{P(y=k)}_{\text{CLASS PRIOR}} \underbrace{P(x|y=k)}_{\text{CLASS CONDITIONAL}} \rightarrow \text{describes } Y \text{ in the distribution of features}$$

- What do they do?
 - They give us the posterior class distribution
 - But also give estimates of Joint distribution, can be used for simulating data
- Modeling the full joint distribution gives posterior probabilities

LDA

- Generative model that gives linear decision boundaries in classification problems

- Model the joint distribution $P(x,y)$
- End up with a discriminant function $g_x(x)$

Used to classify the correct class where among all values of $g_x(x)$:

We choose class with highest

$$y = \arg \max g_x(x)$$

Assumptions:

Class conditionals are Gaussian with:

- Class-specific mean
- Common Variance

We can build a Bayes Classifier Based on this

We can derive an optimal solution under the assumption that we know the true dist of data

- Class conditionals are Gaussian
- We know all parameters (mean, variance)
- We know the class priors.

$$P(x|Y=\text{BLACK}) = N(2,1)$$

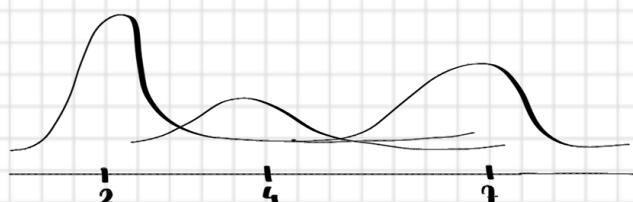
$$\mu_{\text{BLACK}} = 0.6$$

$$P(x|Y=\text{RED}) = N(4,1)$$

$$\mu_{\text{RED}} = 0.1$$

$$P(x|Y=\text{BLUE}) = N(7,1)$$

$$\mu_{\text{BLUE}} = 0.3$$



Now if we look at the log of said distribution:

$$\begin{aligned}\log p(x|y) &= \log p(y=k) + \log p(x|y=k) \\ &= \log \pi_k + \log \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu_k)^2}{\sigma^2}} \right\}\end{aligned}$$

For which by further simplification

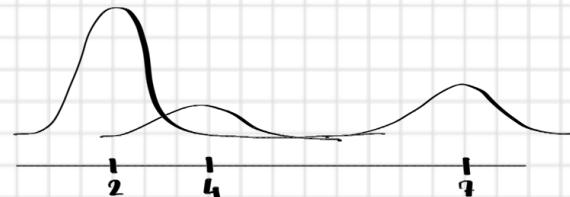
$$g_k(x) = \frac{\pi_k}{\sigma^2} x - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

QDA

- Given that the assumption of a common variance in LDA may be too strict
- In QDA:
 - class conditionals are Gaussian with both class specific mean & class specific variance

Class specific mean & class specific variance

$$\begin{aligned}P(x|y=\text{Black}) &= N(2, 0.25) & \pi_k = 0.6 \\ P(x|y=\text{Red}) &= N(4, 1) & \pi_k = 0.4 \\ P(x|y=\text{Blue}) &= N(7, 0.81) & \pi_k = 0.3\end{aligned}$$



QDA results in HIGHER POSTERIOR

more flexibility of quadratic boundaries

Plug-in Classifiers

- Can Standard \rightarrow a lot of assumptions!!!
- All the knowledge of the params are approximations

Options: $\hat{\pi}_k = \frac{n_k}{n}$

mean: $\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$

Var(lm): $\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^k \sum_{y_i=k} (x_i - \hat{\mu}_k)^2$

Var(ana): $\hat{\sigma}^2 = \frac{1}{n_k-1} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$

Multiple Features

- With multiple features we need multivariate distributions for class conditionals

We talk about:

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$$

$$EX = \begin{bmatrix} E x_1 \\ \vdots \\ E x_p \end{bmatrix}$$

$$\text{Var } X = E[(X - EX)(X - EX)^T]$$

✓ Covariance matrix

- Looking at factors involving x the PDF is

$$p(x, \mu, \Sigma) = \frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)} \rightarrow d\theta \cdot \frac{1}{2} \|x-\mu\|^2$$

EFFECTIVELY BEING A CONCAVE CURVE

LDA & QDA multiple features

LDA \rightarrow INDIVIDUAL MEANS, EQUAL COVARIANCE MATRICES

QDA \rightarrow INDIVIDUAL MEANS, CLASS SPECIFIC COVARIANCE MATRICES

FOR LDA $g_k(x) = 2 \log \pi_k - (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)$

AND THE DECISION BOUNDARIES BETWEEN CLASS j AND k ARE AS FOLLOWS:

$$g_j(x) = g_k(x)$$

