

NN Final 2018 & 2019

1. The machine learning Recipe

1. Train & evaluate network on training data
• If (underfitting) # High Statistical Bias
 • Consider a larger network (more hidden nodes / layers)
 • Train Longer
 # works sometimes
• Check for bias again & keep changing until good bias is reached.

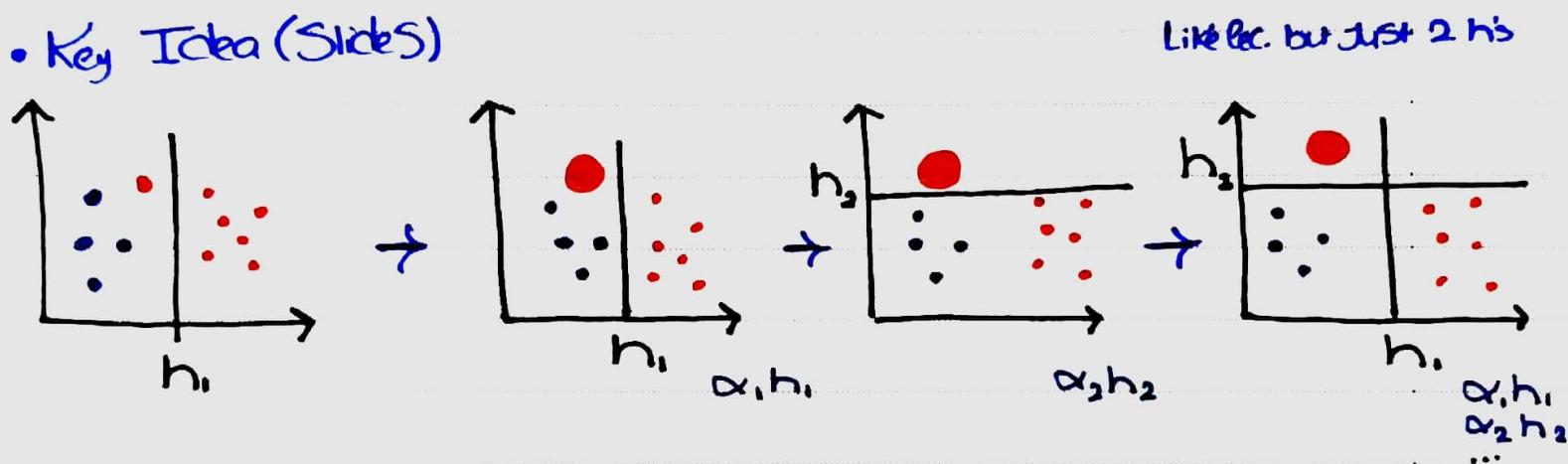
2. Evaluate on validation data

1. Train & evaluate network on training data
• If (overfitting) # High Statistical Variance
 • Consider more data
 • Regularization
• Check for bias again, then for variance & so on till good bias & variance are reached.
- Search for a better NN architecture
 # works sometimes

2. The Key Idea of AdaBoost + Diagram

- Sequentially Combine many weak classifiers in order to realize a strong classifier that approximates Bayes Classifier (Optimal Decision).

- Key Idea (Slides)



→ each classifier emphasizes the mistakes of the previous one and in the end a weighted combination of them decides the output.

<Bonus>

$$1. w_n = \frac{1}{M} \quad 1 \leq n \leq M$$

2. For $1 \leq t \leq T$:

$$3. h_t = \text{train}(\text{weights}, \text{data})$$

$$4. \text{err}_t = \sum_{m=1}^n w_m \cdot I(h_t(x_m) \neq c_m) / \sum_{m=1}^n w_m$$

$$5. \alpha_t = \log \left(\frac{1 - \text{err}_t}{\text{err}_t} \right)$$

$$6. \text{For } m=1 \text{ to } m=M: \quad w_m = w_m \cdot e^{\alpha_t I(h_t(x_m) \neq c_m)}$$

$$7. \quad w_m = w_m / \sum_{m=1}^n w_m \quad 1 \leq m \leq M$$

$$8. \text{Output}(x) = \arg \max_K \sum_{t=1}^T \alpha_t I(h_t(x) = K)$$

</Bonus>

3 - why can't AdaBoost work for 2 classes

- α_t , which signifies the influence of each classifier should be zero (no influence) when the classifier is as bad as random guessing & +ve otherwise.

→ For K classes, the accuracy by random guessing is $\frac{1}{K}$ $(\text{err} = 1 - \frac{1}{K} = \frac{K-1}{K})$

Hence,

$$\begin{aligned}\alpha_t &= \log\left(\frac{1-\text{err}_t}{\text{err}_t}\right) = \left(\frac{1-(K-1)/K}{(K-1)/K}\right) \\ &= \log\left(\frac{1}{K-1}\right)\end{aligned}$$

$$= -\log(K-1) \quad \cdot \text{This only evaluates to zero if } K=2$$

Thus, we can generalize the formula to be

$$\alpha_t = \log\left(\frac{1-\text{err}_t}{\text{err}_t}\right) + \log(K-1)$$

• and now $\alpha_t = 0$ for random guessing and $\alpha_t > 0$ otherwise.

Q4) Wrapper VS. Filter Method



Wrapper Method

- takes classifier into account while seeking the optimal feature set
- accurate but slow & lacks generality
- e.g., Sequential Forward Selection

Filter Method

- does not take specific classifier into account & rather relies on info content, correlation, class means distance & other statistics
- fast & more general but tends to select large subsets
- e.g., a priori

Q5) AI VS. ML VS. DL

AI: Any technique that enables computers to mimic human intelligence using logic, if then rules, decision trees, ML and DL algos

ML: Subset of AI that includes abstruse statistical techniques that enable machines to improve at tasks with experience. Includes DL

DL: SubSet of ML that includes algorithms that enable Software to train Itself on tasks like Speech and Image recognition by exposing multi-layered Neural Networks to vast amounts of data.

<Bonus>

- Note that deep learning algorithms can directly train on raw data without prior feature extraction (can be end-to-end) but require more data volume & computational power.

Classical AI: AI excluding machine & deep learning
Used when no learning data are needed (simple applications).

</Bonus>

6- Next Step after training an NN and getting 54% training accuracy and 51% validation accuracy.

→ The neural network is severely underfit (huge bias) so we should consider:

- larger network (layers / neurons)
- longer training time

If neither that nor changing the architecture helps we should consider substantial changes in terms of the features / choice of model.

- Note that with such bias, it was pointless to check performance on validation set in the first place.

7- Linear Perceptron Update Rule

$$\underline{w}_{\text{new}} = \underline{w}_{\text{old}} + \alpha (d_m - y(m)) \underline{u}(m)$$

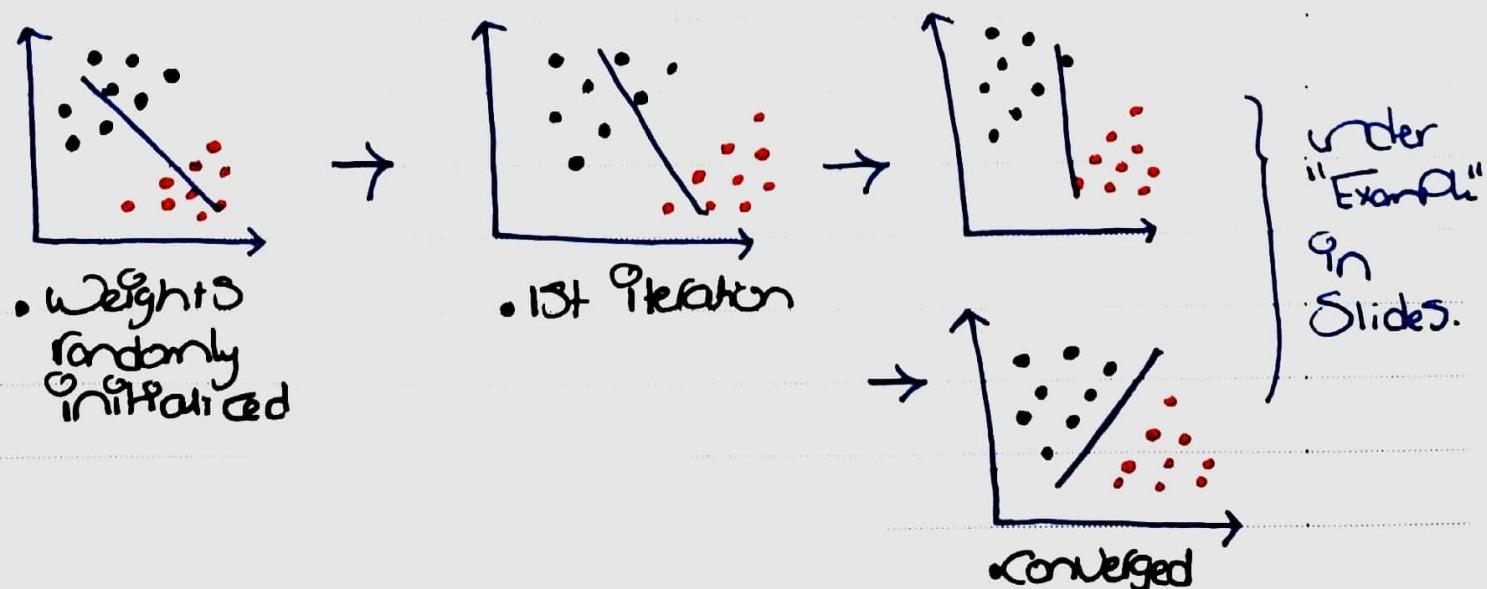
↓ desired output
 ↓ actual output
 ← augmented input

- The goal is to find $\underline{w} = (w_0, w_1, \dots, w_n)$ such that

$$w_0 + w_1 x_1 + \dots + w_n x_n = 0$$

Perfectly separates the two classes.

- Each weight update changes the hyperplane's shift/orientation until it classifies all points correctly



- In Case $d(m) = y(m)$

$$\underline{w}_{new} = \underline{w}_{old}$$

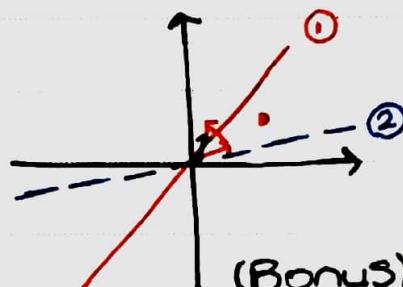
(Correctly classified, no update)

- $d(m) = 1, y(m) = 0 \quad (d(m) - y(m)) = 1$

• Should be above

• actually below

$$\underline{w}_{new} = \underline{w}_{old} + \alpha \underline{y}(m)$$



(Bonus)

$\underline{w}^T \underline{u}$
we get a better
hyperplane

- $d(m) = 0, y(m) = 1$

$$\underline{w}_{new} = \underline{w}_{old} - \alpha \underline{y}(m)$$

<Bonus>

1. Initialize weights w_0, w_1, \dots, w_n randomly

2. While training accuracy < 100% :

3. For $m=1$ to $m=N$:

4. Present $\underline{u}(m)$ and corresponding $d(m)$

5. $y(m) = \underline{w}^T \underline{u}(m)$

6. $\underline{w}_{new} = \underline{w}_{old} + \alpha (d(m) - y(m)) \underline{u}(m)$

</Bonus>

- I doubt if the derivation for "takes the right step each update" is needed but you can write that.

8. Naïve Estimator

1. Associate with each training Pattern a bin of width Δ

2. Associate with each bin the Probability

$$P(X=x) = \frac{n_x}{n} \quad \begin{matrix} \leftarrow & \text{No. of examples} \\ & \text{falling into bin} \end{matrix}$$

and add any overlaps.

- Histogram analysis applies the same formula but uniformly partitions the range of the dataset to produce the bins first.

⇒ Naïve estimator guarantees that the bins are centered at training pattern and hence results in less bias.

- Both correspond to a discontinuous density estimate although the underlying one is smooth.

9. Kernel Density Estimate

$$1. \hat{P}(x) = \frac{1}{Mh^n} \cdot \sum_{m=1}^M g\left(\frac{x - x_m}{h}\right) \quad . Q_n(x) = \frac{1}{n} g(x)$$

- Places a bump function $g(x)$ at each point x_m
- Need $\int_{-\infty}^{\infty} g(x) dx = 1$
-
- Than sums them all

- Typical Choice for $g(\underline{x})$:

$$g(\underline{x}) = \frac{e^{-\underline{x}^T \underline{x}}}{(2\pi)^{\frac{n}{2}}}, \quad \mathcal{Q}_n(\underline{x}) = \frac{e^{-\frac{\underline{x}^T \underline{x}}{n}}}{(2\pi)^{\frac{n}{2}} h^n}$$

• The BW
h acts
as the
variance.

10. Why is ReLU better than Sigmoid

$$\begin{aligned} \text{ReLU}(\underline{x}) &= \begin{cases} \underline{x} & \underline{x} \geq 0 \\ 0 & \underline{x} < 0 \end{cases} \\ &= \max(0, \underline{x}) \\ \left(\frac{d}{d\underline{x}} \right) &= \begin{cases} 1 & \underline{x} > 0 \\ 0 & \underline{x} \leq 0 \end{cases} \end{aligned}$$

$$\begin{aligned} \sigma(\underline{x}) &= \frac{1}{1 + e^{-\underline{x}}} \\ \left(\frac{d}{d\underline{x}} \right) &= \begin{cases} 1 & \underline{x} > 0 \\ 0 & \underline{x} \leq 0 \end{cases} \end{aligned}$$

- Sigmoid suffers from slow learning because it has weak gradients (that even vanish for large values) meanwhile ReLU has steady derivatives and thus exhibits much faster learning
- Besides, ReLU is easier to implement (involves no exponentiation) and empirically works well (esp. learning from images).

II. The Least Square Classifier

→ Goal is to find w_0, w_1, \dots, w_n such that
 $y = \underline{w}^T \underline{u} = w_0 + w_1 x_1 + \dots + w_n x_n$ is a hyperplane
that Predicts the Continuous target variable y so that
Classification can be performed by setting on it a threshold
(e.g., $y > 0 \rightarrow C1$ and $y < 0 \rightarrow C2$).

- This is done by minimizing the loss function

$$\begin{aligned} J &= \frac{1}{N} \sum_{m=1}^N \underbrace{(y_m - d_m)}_{\underline{u}^T \underline{w}}^2 \\ &= \frac{1}{N} (\underline{U} \underline{w} - \underline{d})^T (\underline{U} \underline{w} - \underline{d}) \\ &\quad \downarrow \begin{pmatrix} u^{(1)} \\ u^{(2)} \\ \vdots \end{pmatrix} \downarrow \begin{pmatrix} d^{(1)} \\ d^{(2)} \\ \vdots \end{pmatrix} \end{aligned}$$

→ By $\frac{\partial J}{\partial w} = 0$ we get

$$\underline{w} = (\underline{U}^T \underline{U})^{-1} (\underline{U} \underline{d})$$

which is the \underline{w} for the unique hyperplane that
minimizes the loss.

12. Main Issues of GMMs

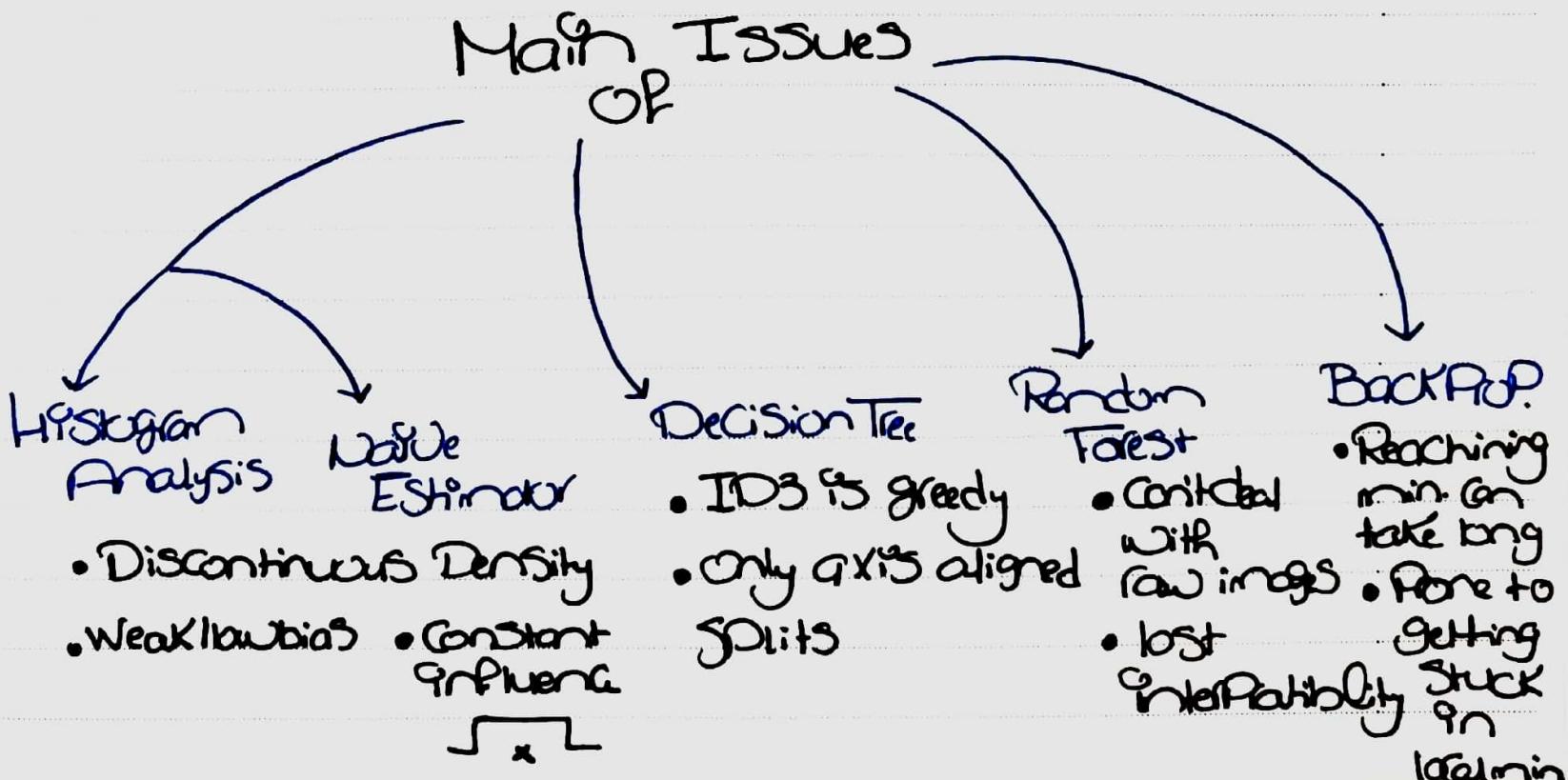
1. Initialization

- Because the underlying function to be optimized is non-convex, GMMs are sensitive to initialization.
- Bad Initialization → Bad estimate

#Can use K-means
to help initialize.

2. Need to decide #Components

- Can either try different values and take best according to Validation Set or use Information theory.



13. TyRESOP Features to be removed

- Irrelevant Features: Don't help discriminate between the two classes (e.g., height for male vs. female)
- Correlated Features: Vary very closely with each. The presence of one of them should be enough. (e.g., area and Perimeter of O)

14. Regularization & why its used.

→ Technique used to Prevent Overfitting which often occurs when we have a large network with many parameters.

- Can be achieved by

* Adding a regularization term to the loss

$$J = \sum_{m=1}^M J_m + \frac{\lambda}{2M} \|W\|_2^2$$

← vector of all network weights

* Dropout regularization where nodes in a layer may drop out each iteration according to a set probability.

Both methods in principle simulate the effect of training a smaller network.

15. Time Series Prediction



De Seasonalize:

1. For each season

$$a(\text{year}) = \frac{1}{12} \sum_{\text{window}} x(t)$$

For each month

$$z(i) = \frac{x(i)}{a(\text{year})}$$

For each month

$$u(i) = \frac{\sum_{\text{years}} z(i)}{\# \text{years}} \quad \left(\frac{\# \text{years}}{\sum_{j=1}^{\# \text{years}} z_j(i)} \right)$$

$$x(t)|_{\text{des}} = \frac{x(t)}{u(\text{month}(t))}$$

Predict Trend: $\hat{x}(t)|_{\text{des}}$

Recover Seasonality:

$$\tilde{x}(t) = \hat{x}(t)|_{\text{des}} \cdot u(\text{month}(t))$$

16.

1. For each class C_i assume & find Parameters (e.g. Gaussian) or estimate the conditional density $P(\underline{x}|C_i)$
2. Estimate $P(C_i)$ from dataset / real world
3. Given \underline{x} , Predict the class K where
$$K = \arg \max_i P(\underline{x}|C_i)P(C_i)$$

This is Bayes Classifier which comes from Bayes Rule

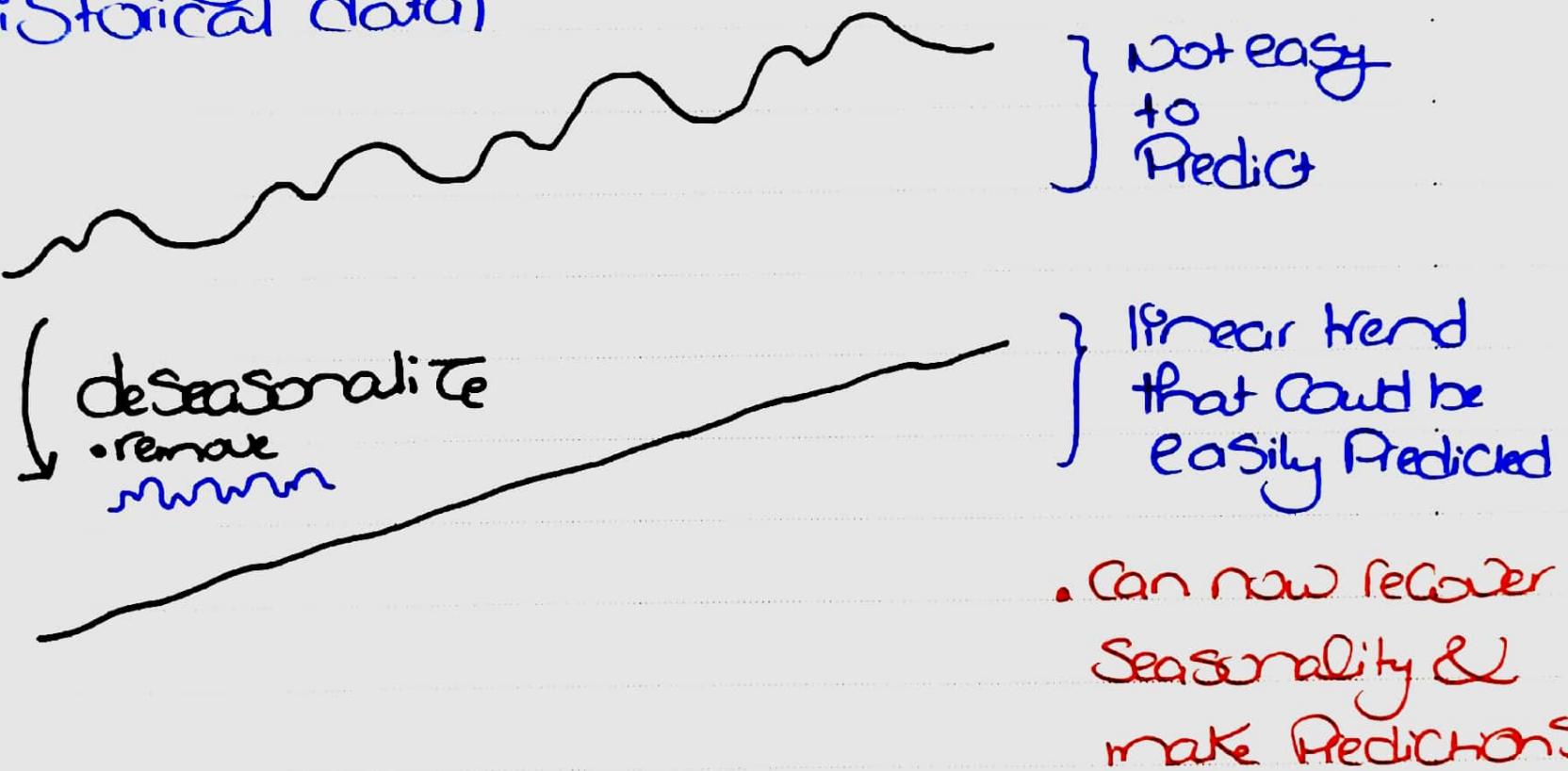
$$P(C_i | \underline{x}) = \frac{P(\underline{x}|C_i)P(C_i)}{P(\underline{x})}$$

→ Given a new sample \underline{x} we compute the Probability of each class being the corresponding one then taking the highest

→ We can ignore the marginal $P(\underline{x})$ since it's common to all evaluations. This corresponds to the Bayes classifier.

17. Explain the Importance of deSeasonalization using diagrams

- It allows us to Perform Time Series Forecasting (Predicting future values based on time-stamped historical data)



18. BackPropagation Steps & Disadvantages

1. Randomly Initialize all weights & biases
(e.g., in $[-1, 1]$)

2. For $m=1$ to $m=M$:

3. Forward Pass

Feed $u(m)$ & find Layer Outputs

$$l=1, \dots, L \quad \begin{cases} X_i^{(l)} = \sum_{j=1}^{N(l-1)} w_{ij}^{(l)} Y_j^{(l-1)} + w_{i0}^{(l)} \\ Y_i^{(l)} = P(X_i^{(l)}) \end{cases} \quad l \in \{1, \dots, L\}$$

4. Backward Pass

Use d_m & Forward Pass Outputs

to Compute $\partial J_m / \partial w_{ij}^{(l)}$ for all layers i.e.

$$l=L, \dots, 1 \quad \begin{cases} \cdot \partial J_m / \partial X_I^{(l)} \\ \cdot \partial J_m / \partial w_{ij}^{(l)} \end{cases} \quad \begin{matrix} I \in N(l) \\ \forall i, j \end{matrix}$$

5. Update Weights

$$w_{ij}^{(l)} = w_{ij}^{(l)} - \alpha \frac{\partial J_m}{\partial w_{ij}^{(l)}} \quad \forall i, j, l$$

6. Compute total loss and Stop once Converged

Disadvantages:

- Can take long to reach the minimum (esp. when close)
 - α is low \rightarrow too slow
 - α is large \rightarrow oscillations / not converging
- Decaying α is a solution

• Prone to getting stuck in local minima
 * loss is highly non-convex
 \rightarrow Solution can be to repeat training many times @ different weight initializations



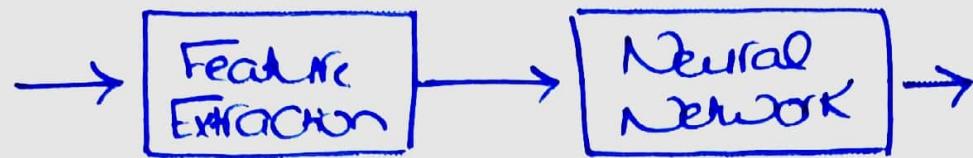
19. Three Methods to update weights

$$w_{\text{new}} = w_{\text{old}} - \alpha \sum_{n=1}^N \frac{\partial J_n}{\partial w_{\text{old}}}$$

- Full batch Gradient Descent
 - \rightarrow Slow but consistent
 - * Consider whole dataset for a single update

- Stochastic GD
 - \rightarrow Fast but prone to oscillations / bouncing around min.
 - \rightarrow Doesn't make use of all parallelism possible
 - * Consider single example each step
- STEP POWER OF 2
- N
- Mini-batch
- \rightarrow Faster & more consistent
- * Consider fraction of dataset N in each step

20.



- CNNs Precede the feedforward layers with feature extraction layers that extract features in light of minimizing the loss

Two types of layers

Convolutional
of K filters each $P \times P$

→ Objective is to learn features. Each feature has as much channels as input so each pair is convolved then added (along with bias) to produce a feature map

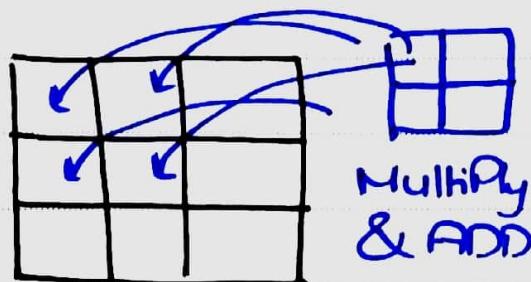
Pooling layer with $L \times L$ window

→ Objective is to summarize the feature map by avg-ing or max-ing over a window (no learning)

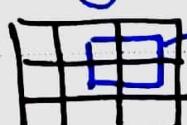
→ When window is as large as feature map it converts d-channel feature map to d-vector

Suitable Input for NN

• Convolution



• Pooling



Global Pooling

Less Memory Footprint

Parameter Sharing

- Same filter applied on different feature map pixels

(not that each has its own weights)

Sparcity of Connections

- Output depends on a limited no. of inputs (below filter).
(not that it depends on all of them)

21. Three Limitations of Deep Learning

1. Unintended results from fitness function
(on minimize loss in expected ways)

2. Sensitive to Standard Adversarial Attacks
( +  → wrong classification)

3. OverSensitivity to changes in context
(data set has limited no of them, in real world
its Combinatorially large)

4. Large Amount of labeled data (expensive)

23. K-Fold Cross Validation

1. Partition the data set into K folds
 2. Train the model K times. In each consider $K-1$ folds only then validate on the last one.
 3. Average the K resulting performance estimates
- It's used to arrive at (Statistically Significant*) Performance estimate of the model. In this sense it makes the validation set only optional
 - Can hence be used to tune hyperparameters

My thanks <3