

Specialization -03

"Machine Learning Modeling pipelines in production"

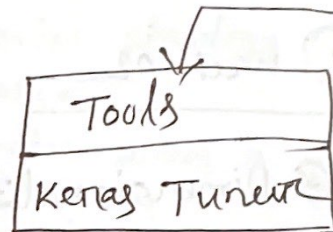
Week-01

④ Hyperparameter Tuning

- Neural architecture search (NAS) is a method for automating the design of ANN.

- Tune parameters

- Architecture choice
- activation functions
- weight initial strategy
- Learning rate, stop condition



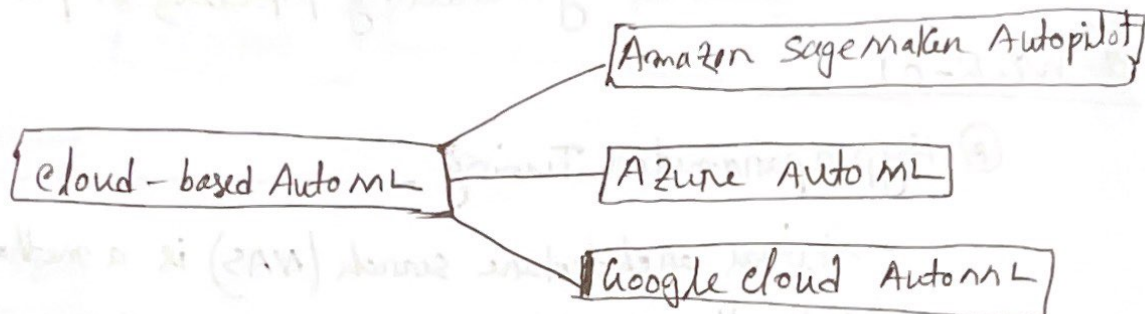
④ Auto ML → applies machine learning and search techniques to the process of create machine learning models and pipelines. It covers the complete pipeline from the raw dataset to the deployable.

Auto ML automates the entire ML workflow.

- NAS: a few search strategies

- Grid search
- Random search
- Bayesian optimization
- Evolutionary Algorithms
- Reinforcement Learning

* AutomL on the cloud



Week 02

* Dimensionality

- dimensionality reduction
 - Linear dimensionality reduction
 - principal component Analysis (PCA)
- classification: LDA
- Regression: PLS (partial Least Squares)
- Unsupervised: PCA
- others: SVD, ICA, NMF (Non-negative Matrix Factor)
 - ↳ decomposes non-square matrices.

* Quantization & Pruning (Model deployment on Edge devices)

-
- ```
graph TD; A[Quantization & Pruning] --> B[Shrinking model file size]; A --> C[post training quantization]; A --> D[Quantization Aware training (QAT)]; A --> E[ML Kit]; A --> F[Core ML]; A --> G[TF Lite]
```
- A central box labeled "Quantization & Pruning" has arrows pointing down to three boxes: "Shrinking model file size", "post training quantization", and "Quantization Aware training (QAT)". To the right of these boxes are three more boxes: "ML Kit", "Core ML", and "TF Lite", which are also connected to the central box by lines.
- Shrinking model file size
  - post training quantization
  - Quantization Aware training (QAT)
  - ML Kit
  - Core ML
  - TF Lite



- pruning

$$P_n = 1 - (1 - p)^n$$

## # Week - 03

### # High performance Modeling

- Distributed training
  - Data parallelism
  - Model parallelism
- tf. distribute. strategy
  - one device strategy
  - mirrored "
  - parameter server "

Tools → **GPipe** → distributed parallelism

- Knowledge Distillation (KD)

## # Week - 04

### # Model performance Analysis

### # model Debugging

- Benchmark models

### # Sensitivity Analysis & Adversarial Attacks

- partial dependence plots.

- PDP box and pyCE box

Tools

→ Tenson board

→ **TFma/tfma**

Model performance Analysis.

Tool

→ what-if

→ Adversarial attack is to confuse the model.

Adversarial Attack  
↓  
Tools

cleverhans

Foolbox

# ~~Ris~~ Residual Analysis

# Model Remediation

# Statistical process control

$$\mu = np_+$$

$$\sigma = \sqrt{\frac{p_+(1-p_+)}{n}}$$

# Sequential Analysis

$$p_+ + \sigma_+ \geq p_{\min} + 3\sigma_{\min}$$

# clustering

- Algorithm = OLINDA, MINAS, ECS Miner, GCS

\* google cloud AI continuous Evaluation

⊕ week 05

⊕ Explainable AI (XAI)

⊕ Responsible AI (RAI)

⊕ Interpretability

⊕ model specific or model Agnostic

→ Interpretable Models

• Linear • Logistic • Decision Tree • RuleFit • K-nearest Neighbors

• TF Lattice

## → Model Agnostic methods

- These methods separate explanations from the machine learning model.

- Partial Dependence plots (PDP)
- permutation Feature
- Shapley values
- Testing concept Activation
- LIME
- TCAV & LIME
- AI Explanations.