**Materials**
**[MLE Best practices](#)**

**Homework**

During the course you will be using Linux, thus if your OS is not Linux, install VirtualBox and use any Linux distribution (Ubuntu is the best choice for a Linux newbie).

1. Download any Linux distributive, use Ubuntu by default [https://ubuntu.com/desktop](https://ubuntu.com/desktop)
2. [How to install Linux on VirtualBox, watch up to 9th minute](#)
3. Install Docker and Docker Compose on Linux. How to run [docker without sudo](#)
4. Create a repository on your github.

A repository should contain separate folders for each further task. For each homework in addition to the code, you should provide a README file with the commands to perform the task and your description/comments.

*Useful materials & links for the course (optional, for further self-study of topics)**:**

| Topic | Links |
|---|---|
| Common | [Useful git repository](#) |
| Cloud solutions | Google cloud Cloud ML Engine<br>AWS SageMaker<br>Azure Machine Learning studio |
| Platforms | Use cases of existing architectures,<br>    ● Uber<br>    ● Netflix<br>    ● Airbnb<br>Ready to use solutions<br>    ● Domino<br>    ● Datalab<br>    ● MLFlow<br>    ● MetaFlow (AWS only)<br>EPAM accelerators<br>    ● Legion<br>Kubernetes-based:<br>    ● Kubeflow<br>    ● Polyaxon |
| ML development/research process | Defining iterations:<br>https://blog.insightdatascience.com/how-to-deliver-on-machine-learning-projects-c8d82ce642b0<br>Caveats:<br>https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf<br>Workflow managers<br>    ● Luigi |

| | |
|---|---|
| | <ul><li>Airflow</li><li>Dagster</li></ul>Dividing batch data flows from streaming data flows<br>Version control / model versions / collaboration tools<ul><li>DVC</li></ul>Experiment management:<ul><li>TRAINS (separate client-server, good-looking UI, better configuration then in MLFlow)</li></ul> |
| Optimizations | Scaling: Spark, distributed training/inference for Deep Learning<ul><li>Horovod</li><li>Ray</li></ul>Low-level: optimizing frameworks for inference, low-level memory manipulations/parallelism to optimize models training/inference<ul><li>Treelite</li></ul> |
| Domain specific knowledge | Storage solutions for large amount of dense matrices<ul><li>TileDB</li></ul>Text processing: metadata governance, layers (raw -> cleaned -> tokenized -> etc.)<br>CV: storage and basic manipulations (up/down sampling, cleaning, augmentations)<ul><li>PicPoc</li></ul> |
| AUX | Configuration management<ul><li>Hydra: configuration management/cli tool (github)</li></ul>Containers<br>REST API deployment option<br>Monitoring pipelines: visualizations, reproducibility<br>Tooling/reporting:<ul><li>Streamlit: Fast reporting with advanced data caching and interactive elements</li></ul> |