

CAI 4104/6108 – Machine Learning Engineering: ML Engineering (2)

Prof. Vincent Bindschaedler

Spring 2024

Administrivia HiPerGator

■ Teaching Allocation

- ◆ Please review the **Acceptable Use Policy**

- ☼ <https://it.ufl.edu/policies/acceptable-use/acceptable-use-policy/>

- ◆ Keep in mind: these are **shared resources**!

- ◆ All data will be **deleted** on **03 May 2024**

- ◆ Training: https://help.rc.ufl.edu/doc/New_user_training

- ◆ **Academic Integrity**: Make sure your files are not publicly readable (use chmod)



■ HiPerGator Help

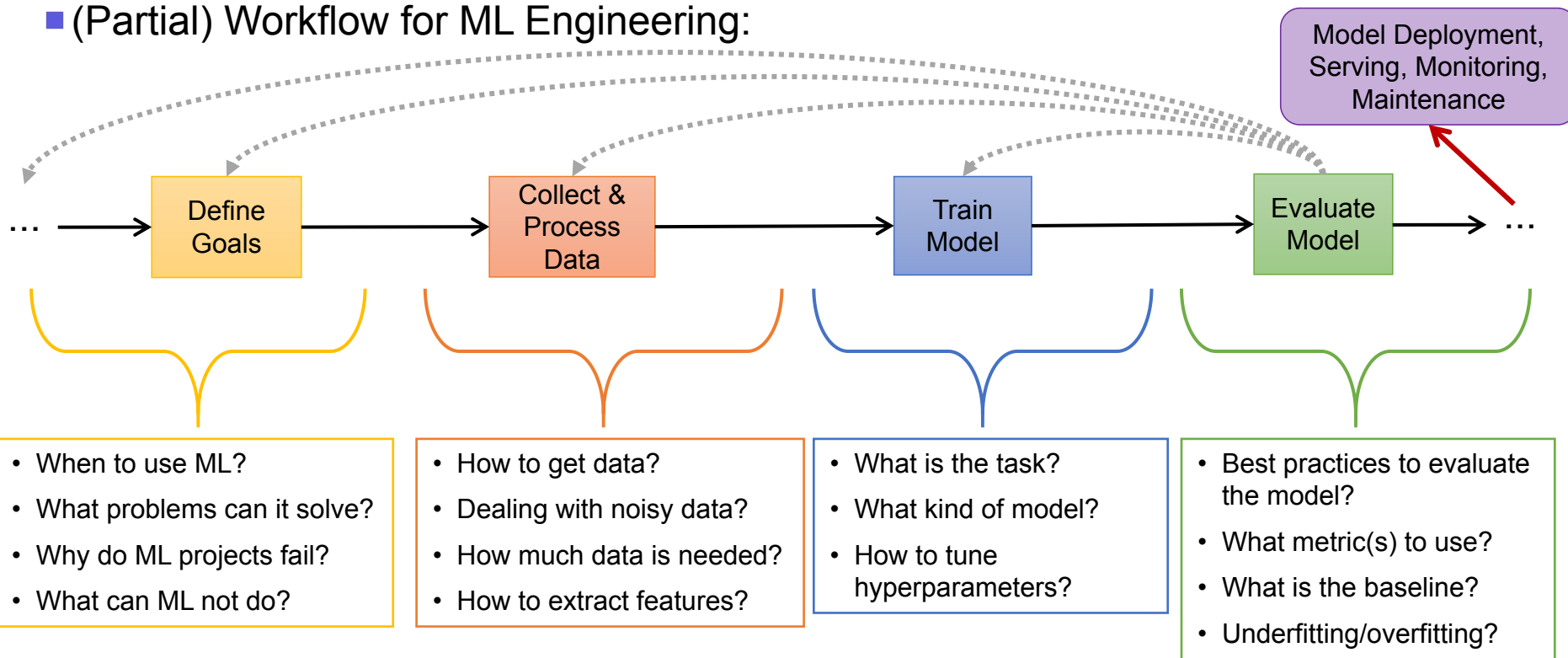
- ◆ Prerecorded training: https://help.rc.ufl.edu/doc/Prerecorded_Training

- ◆ Wiki: <https://help.rc.ufl.edu/>



Reminder: Machine Learning Engineering

■ (Partial) Workflow for ML Engineering:



Reminder: Many Reasons Not To Use ML

- Cannot get the (right) data or enough of it
- Problem does not require learning from data
- You can solve the problem in other ways
 - ◆ By developing new algorithms, using software development, etc.
- Cannot afford the cost of a mistake
- Judged **unethical/undesirable** to use ML
 - ◆ Should we use ML to determine prison/jail time?
- You need explanations, not just predictions
 - ◆ The problem requires **explaining some phenomenon** (e.g., physics)
 - ✿ ML is **not** magic: it doesn't have deep insights about the natural world!
 - ✿ It's (just) a set of techniques and tools to make accurate predictions from data
 - ◆ *Note: there is active research on interpretable/explainable ML; but many models are black boxes!*
- Many others...

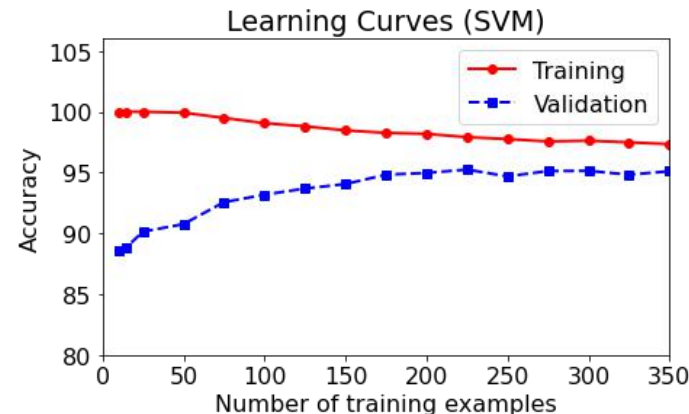
Reminder: Collecting Data

■ Many questions:

- ◆ Where is the data coming from? Is it reliable?
- ◆ Does it need to be anonymized? (What can happen if we don't?)
- ◆ Is the data noisy/incomplete? Is it missing values?
- ◆ Is the data biased in some way?
 - ✧ Is the data expected to be similar to the data available after the system is deployed?

■ How much data do we need?

- ◆ More than you expect
- ◆ Rules of thumb:
 - ✧ More examples/instances than attributes/features
 - Some models (e.g., SVM) still work with few examples but lots of features
 - ✧ Many examples (e.g., 30+ or 50+) for each class
- ◆ How to answer this in practice? => **learning curves**
 - ✧ You'll be able to try this in homework1



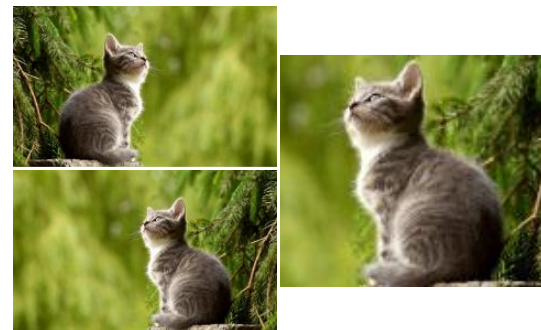
Reminder: Dealing With Unclean Data

- The data is **noisy**
 - ◆ For example: some attribute values/features are wrong
 - ✿ Some learning algorithms are more sensitive than others
- The data has missing features values
 - ◆ Mitigation strategies
 - ✿ **Algorithmically**: use a learning algorithm that can deal with this (e.g., some Decision Trees implementations)
 - ✿ **Removal**: discard examples with missing values
 - ✿ (If feature is categorical/nominal) **Treat missing as value**: replace missing values with a “special value” (e.g., -1)
 - ✿ **Imputation**: use a data imputation technique (e.g., resampling from marginals, replace with average)
- The data contains duplicates
- The data is imbalanced
 - ◆ Weighting classes (if allowed by learning algorithm); oversampling; undersampling
- The data is incomplete or not representative?
 - ◆ Find better data!

Too Little Data? Too Much Data?

■ Data Augmentation

- ◆ Idea: **augment** the existing training data by using the training data!
 - ✿ Can have a **major** impact on model performance
 - ✿ But it won't help if you truly don't have enough data!
- ◆ For example: for images: rotation, scaling, crop, flip, etc.
- ◆ Data augmentation techniques are usually specific to the data type



■ Data Sampling

- ◆ If you have a lot of data, it may not be practical (and necessary) to use all of it
- ◆ Different strategies for sampling
 - ✿ **Random sampling**: e.g., pick a uniformly random subset of instances, pick each instance independently with some probability p
 - ✿ **Stratified sampling**: divide dataset into groups (**strata**), then randomly sample from each stratum
 - Note: this can be used to reduce bias (i.e., make the dataset more representative)

Data Augmentation?

- Data Augmentation can be thought of a **regularization** technique

Single sample

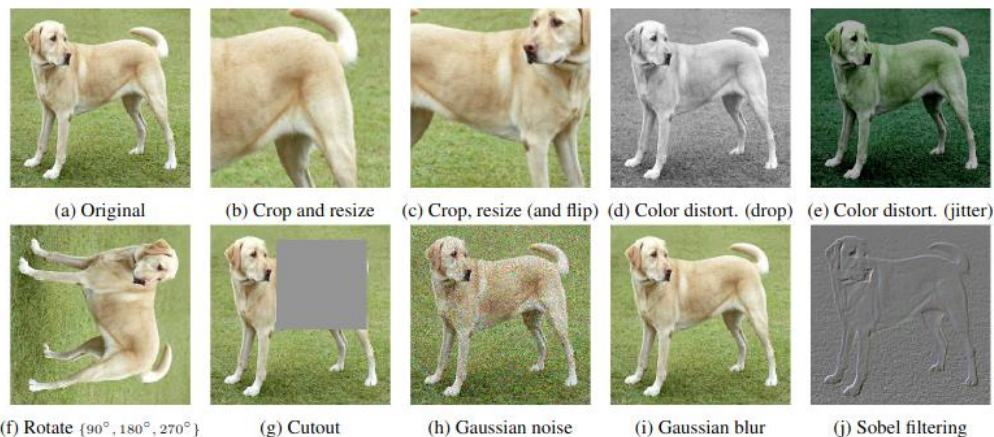
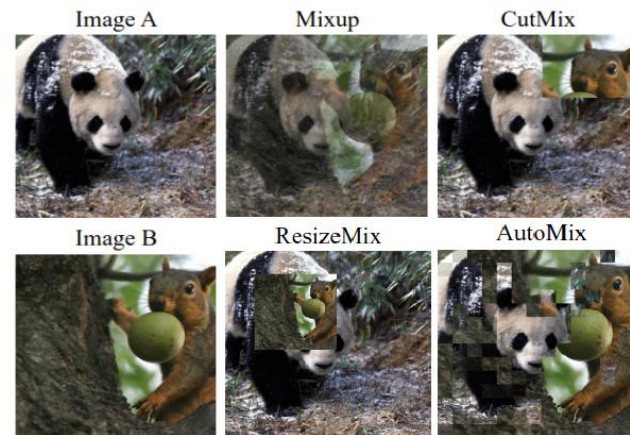


Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize)*, *color distortion*, and *Gaussian blur*. (Original image cc-by: Von.grzanka)

Source: Chen et al. "A simple framework for contrastive learning of visual representations." ICML 2020.

Multi-sample



Source: Li et al. "Openmixup: Open mixup toolbox and benchmark for visual representation learning." arXiv preprint arXiv:2209.04851.

■ Feature Engineering

- ◆ The process of transforming raw data into a dataset we can use
 - ✿ Time consuming process
 - ✿ Techniques are often specific to the data type you are dealing with (e.g., feature engineering for text)
- ◆ Goal: select a set of features that are highly predictive
 - ✿ But also not correlated among each other. **Why?**

■ Common Techniques

- ◆ For numerical features:
 - ✿ **Binning** (aka bucketing): aggregate range of numerical values into discrete bins
 - ✿ **Normalization**: remap values onto a range such as $[0,1]$ (**min-max normalization**) or $[-1,1]$
 - ✿ **Standardization** (aka z-score normalization): rescale feature values to follow a *standard normal distribution* (how? => subtract the mean and divide by the standard deviation)
- ◆ For categorical features:
 - ✿ **One-hot encoding**: turn a categorical feature with k possible values into a vector of k binary features
 - ✿ **Ordinal encoding**: if values are ordered/ranked, values can be encoded in order (e.g., using integers)

■ For numerical features:

- ◆ **Binning** (aka bucketing): aggregate range of numerical values into discrete bins
- ◆ **Normalization**: remap values onto a range such as $[0, 1]$ or $[-1, 1]$
- ◆ **Standardization** (aka z-score normalization): rescale feature values to follow a *standard normal distribution* (how? \Rightarrow subtract the mean and divide by the standard deviation)

■ For categorical features:

- ◆ **One-hot encoding**: turn a categorical feature with k possible values into a vector of k binary features with hamming weight 1
- ◆ **Ordinal encoding**: if values are ordered/ranked, values can be encoded in order (e.g., as integers)

■ Rules of thumb:

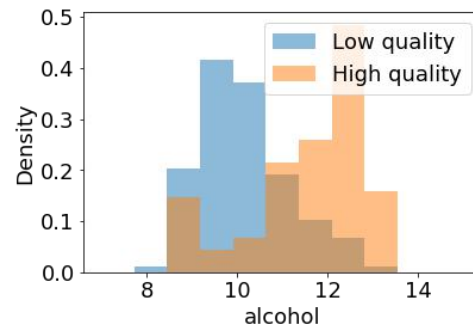
- ◆ If you have a mix of numerical and categorical features: use one-hot encoding / ordinal encoding
- ◆ Feature scaling (i.e., normalization & standardization) helps with most learning algorithms
 - ✱ For SVM (especially linear SVM), linear models (or if you are doing regularization): **you should rescale features!**

■ Preprocessing, data cleaning, and feature engineering

- ◆ Avoid manipulating the data in ad-hoc (i.e., non-reproducible ways) such as bash scripting, awk, manual editing, etc.
- ◆ For reproducibility, it is best to implement all the steps in a (Python) script
 - ✧ Ideally, you want to develop a **pipeline**

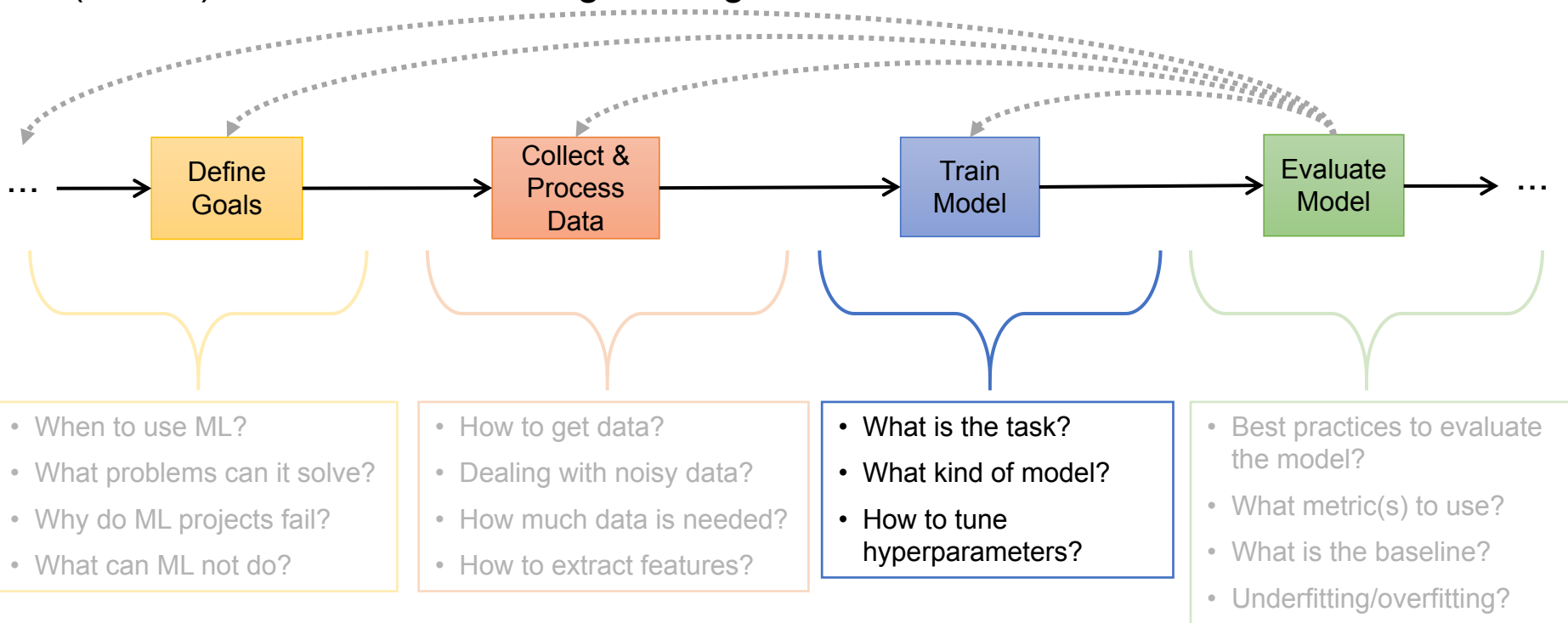
■ Visualizing data:

- ◆ Before getting too far into data cleaning or training a model, take a look at the (training) data!
 - ✧ Visualizing the data may yield insights into what features to use and what kind of model to use
 - ✧ You can also look at feature statistics: mean, standard deviation, min, max, mode, missing values, etc.
- ◆ Visualization examples:
 - ✧ Histograms,
 - ✧ Heatmaps, correlation plots,
 - ✧ Scatter plots, QQ-plots,
 - ✧ Class distribution plot per feature



Machine Learning Engineering

■ (Partial) Workflow for ML Engineering:



Types of Learning

■ Supervised Learning

- ◆ Learning from **labeled data** (i.e., each example or instance in the dataset has a corresponding label)
- ◆ Tasks: **classification** vs. **regression**

■ Unsupervised Learning

- ◆ Learning from **unlabeled data** (we must **discover patterns** in the data)
- ◆ Tasks: clustering (e.g., K-means), dimensionality reduction (e.g., t-SNE, PCA), etc.

■ Semi-supervised Learning

- ◆ Learning from **partially labeled data**

■ Reinforcement Learning

- ◆ There is an **agent** that can interact with its **environment** and perform **actions** and get **rewards**
- ◆ Learning a policy (i.e., a strategy for which actions to take) to get the most rewards over time

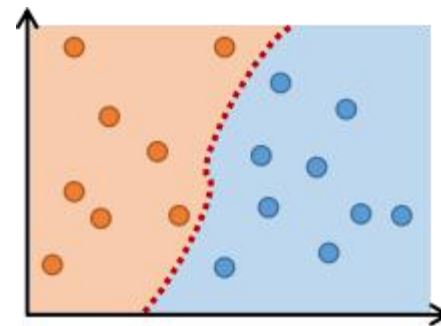
■ Transfer Learning

- ◆ Learning to **repurpose** an existing model for a new task

Supervised Learning

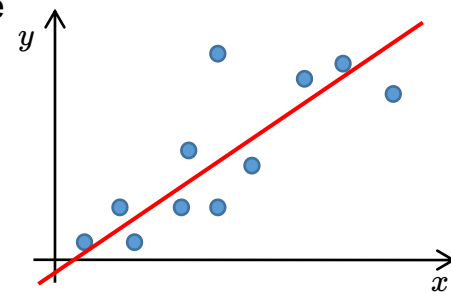
■ Classification

- ◆ Task: predict the corresponding **label**
- ◆ Different types:
 - ✧ **Binary classification**: there are only two classes (0,1 ; +,-, etc.)
 - ✧ **Multiclass**: more than two classes
 - ✧ **Multi-label**: each instance can belong to more than one class (e.g., label all objects in a photo)
 - ✧ **One-class**: there is only one class, we want to distinguish it from everything else



■ Regression

- ◆ Task: predict the corresponding **value** (typically a real number) or **target**
 - ✧ E.g.: you want to predict a person's future income based on their high school GPA



■ Others:

- ◆ Sequence-to-sequence, similarity learning/metric learning, learning to rank, etc.

Multiclass Classification

- **Multiclass classification** (aka *multinomial classification*)
 - ◆ There are $c > 2$ distinct classes: $y_i \in \{1, 2, \dots, c\}$ is the label of the i^{th} example
- Wait. How do we do this?
 - ◆ Recall the SVM formulation:
 - ✱ We want to find a **hyperplane** $w \cdot x - b = 0$, also we relabel the classes as +1 and -1. **Problem?**
- Some learning algorithms / models naturally support multiclass classification
 - ◆ E.g.: kNN, decision trees, neural networks
- For others, we can transform multiclass classification into a binary classification
 - ◆ **One-vs-rest** (OvR): Train c binary classifiers. f_i to classify class i versus not i
 - ✱ Predict: $y = \operatorname{argmax}_i f_i(x)$
 - ◆ **One-vs-one** (OvO): Train $c(c-1)/2$ binary classifiers. $f_{i,j}$ to classify class i versus class j
 - ✱ Predict: predict with all $c(c-1)/2$ and return the class that has the highest number of “votes”

Next Time

- Wednesday (1/24): Lecture
 - ◆ Topic: ML Engineering (3)

- Upcoming:
 - ◆ **Homework 1** will be out soon (due 2/2 by 11:59pm)