

The Last Few Bits..

CS771: Introduction to Machine Learning

Piyush Rai

Debugging ML Algorithms



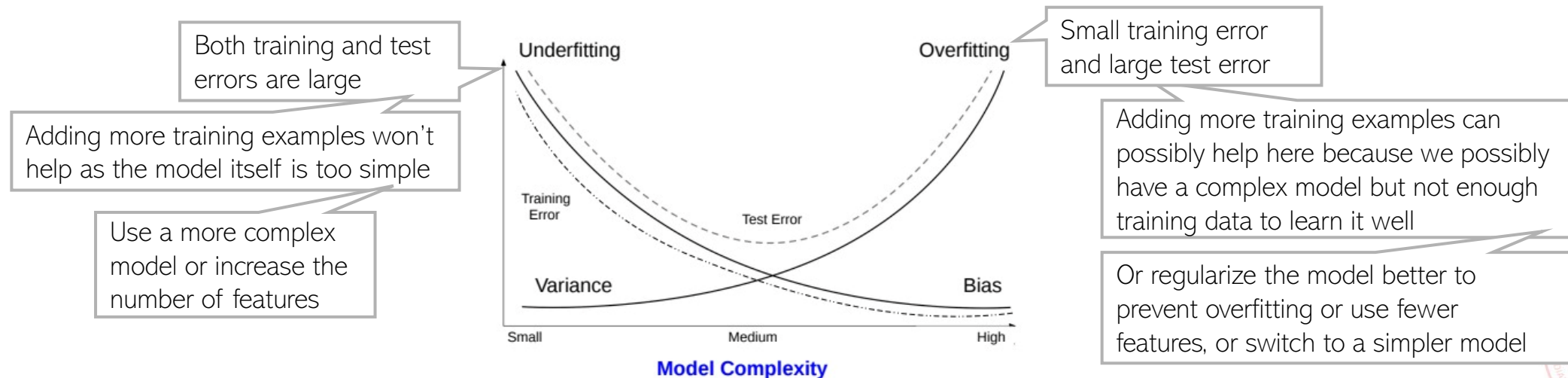
What is going wrong?

- What to do when our model (say logistic regression) isn't doing well on test data
 - Use more training examples?
 - Use a smaller number of features?
 - Introduce new features (can be combinations of existing features)?
 - Try tuning the regularization parameter?
 - Run (the iterative) optimizer longer, i.e., for more iterations
 - Change the optimization algorithm (e.g., GD to SGD or Newton..) or the learning rate?
 - Give up and switch to a different model (e.g., SVM or deep neural net)?



High-Bias or High-Variance?

- The bad performance (low accuracy on test data) of a model could be due either
 - **High Bias**: Too simple model; doesn't even do well on training data
 - **High Variance**: Even small changes in training data lead to high fluctuation in model's performance
- High bias means **underfitting**, high variance means **overfitting**
- Looking at the training and test error can tell which of the two is the case

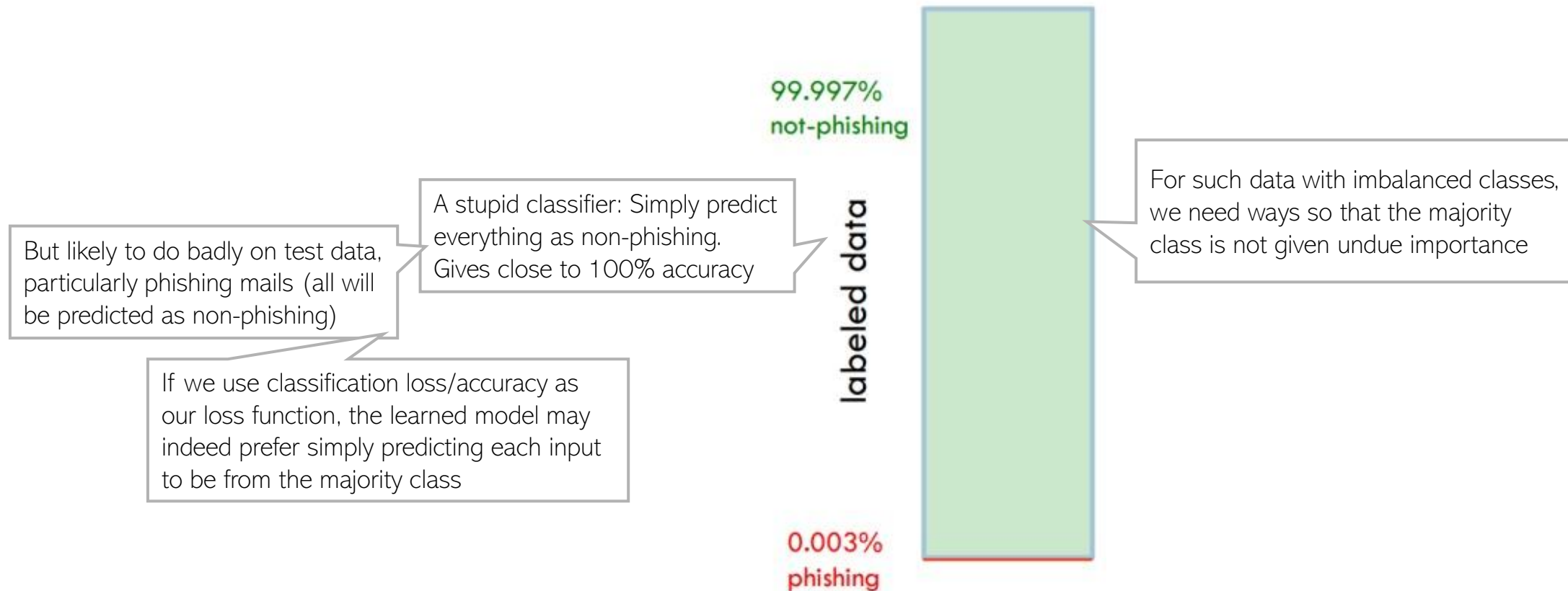


Learning from Imbalanced Data



Learning when classes are imbalanced

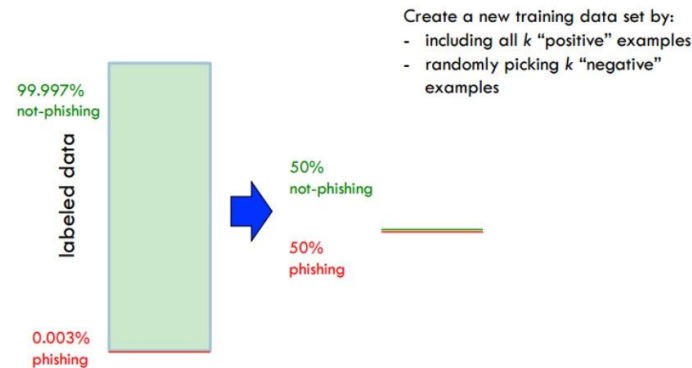
- When classes are imbalanced, even a “stupid” classifier can give high accuracy but looking at accuracy alone may be misleading



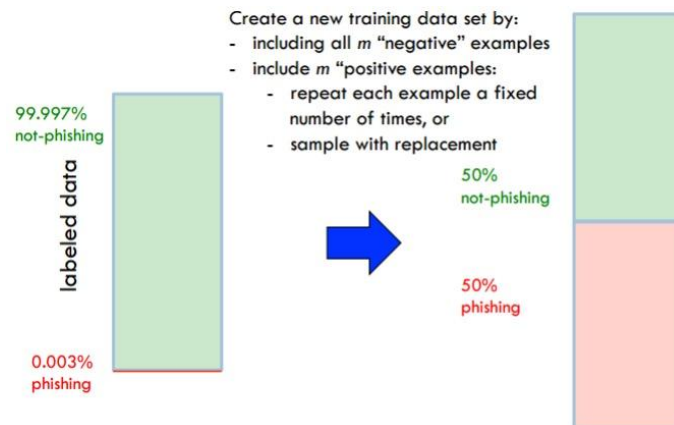
Solution 1: Balancing the training data

7

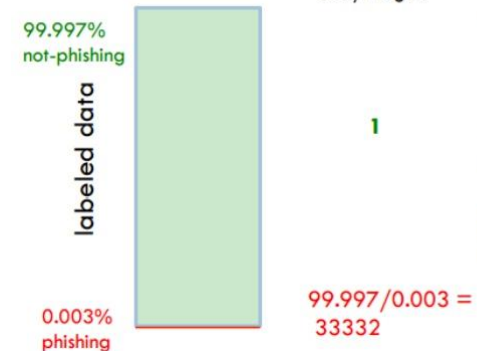
- Can balanced the training data by
 - Under-sampling the majority class examples



- Over-sampling the minority class examples



Equivalent to



Weighted loss function with much larger importance given to loss function terms of positive examples than negative examples

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \beta_{y_i} \ell(\mathbf{x}_i, y_i, \mathbf{w})$$

where $\beta_{+1} \gg \beta_{-1}$

Add costs/weights to the training set

"negative" examples get weight 1

"positive" examples get a much larger weight

change learning algorithm to optimize weighted training error



Solution 2: Changing the loss function

- Don't use loss functions that define loss or accuracy on per-example basis

This loss function is a simple sum of losses on individual training examples.
Not ideal for imbalanced classes

$$L(\mathbf{w}) = \sum_{i=1}^N \ell(x_i, y_i, \mathbf{w})$$

- Instead, use loss function that use **example pairs** (one positive and one negative)
- Assuming our model to be defined by some function $f(\mathbf{x})$ (e.g., $\mathbf{w}^T \mathbf{x}$), define a loss

An input with positive label

An input with negative label

$$\ell(f(\mathbf{x}_n^+), f(\mathbf{x}_m^-)) = \begin{cases} 0, & \text{if } f(\mathbf{x}_n^+) > f(\mathbf{x}_m^-) \\ 1, & \text{otherwise} \end{cases}$$

Now we don't care about per-example accuracy but care about whether the positive examples get a higher score than the negative examples (i.e., we are only preserving their relative rank)

Such loss functions can known as **"pairwise loss functions"**

$$\sum_{n=1}^{N_+} \sum_{m=1}^{N_-} \ell(f(\mathbf{x}_n^+), f(\mathbf{x}_m^-)) + \lambda R(f)$$

Usual regularizer on f



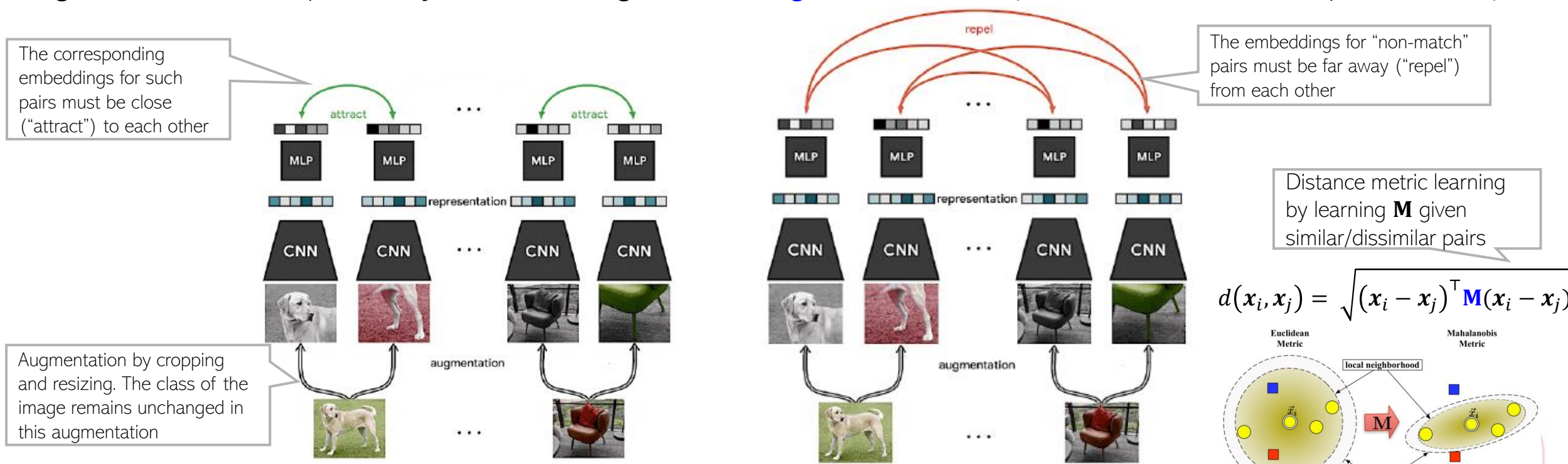
Contrastive Learning



Learning Good Features by Comparison

Or “triplets” (e.g., “cat” is more similar to “dog” than to a “table”)

- Can learn good features by comparing/“contrasting” similar and dissimilar object pairs
- Such pairs can be provided by to the algorithm (as supervision), or the algorithm can generate such pairs by itself using “data augmentation” (as shown in example below)



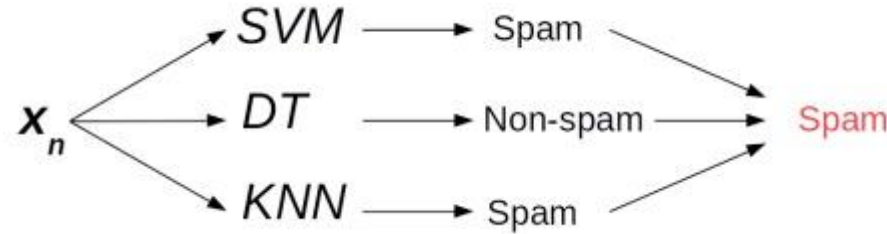
- Such “contrastive learning” of features is also related to “distance metric learning” algos

Ensemble Methods



Some Simple Ensembles

- **Voting** or **Averaging** of predictions of multiple models trained on the same data

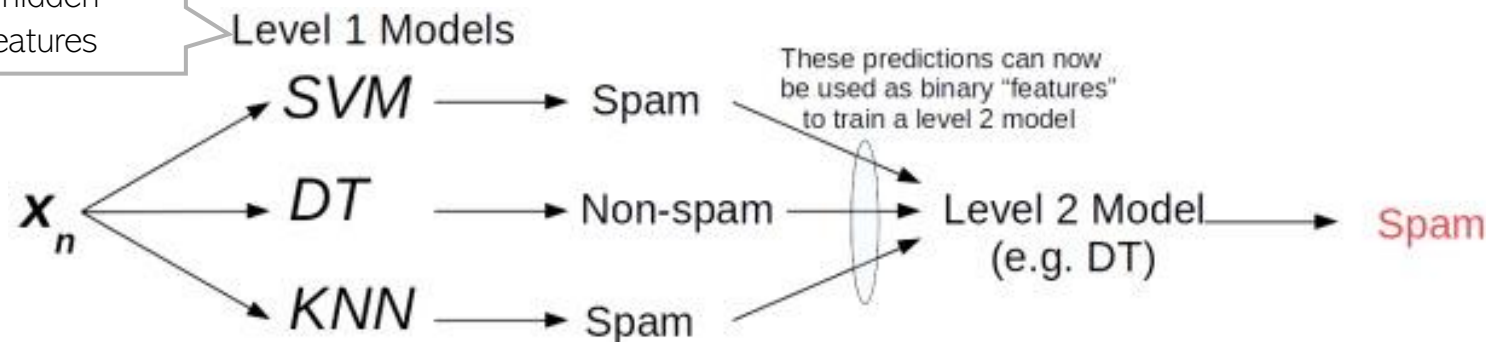


- **“Stacking”**: Use predictions of multiple already trained models as “features” to train a new model and use the new model to make predictions on test data

Stacking sort of has a flavor of deep learning where hidden layers extract good features

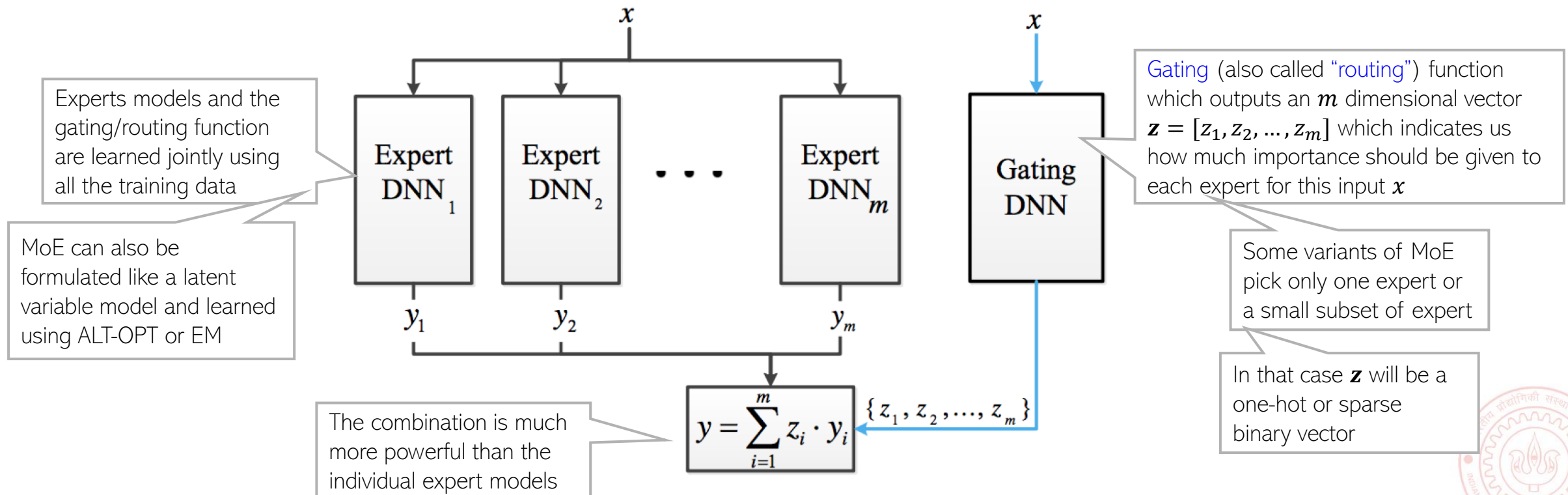
Here, that role is being played by these other already trained models

In stacking, level 1 and level 2 models are trained independently (first level 1 and then level 2)



Mixture of Experts (MoE) based Ensemble

- Mixture of Experts (MoE) is a very general idea
- We assume m “simple” models, usually of the same type, e.g., m linear SVMs or m logistic regression models, or m deep neural nets (usually all with same architecture)

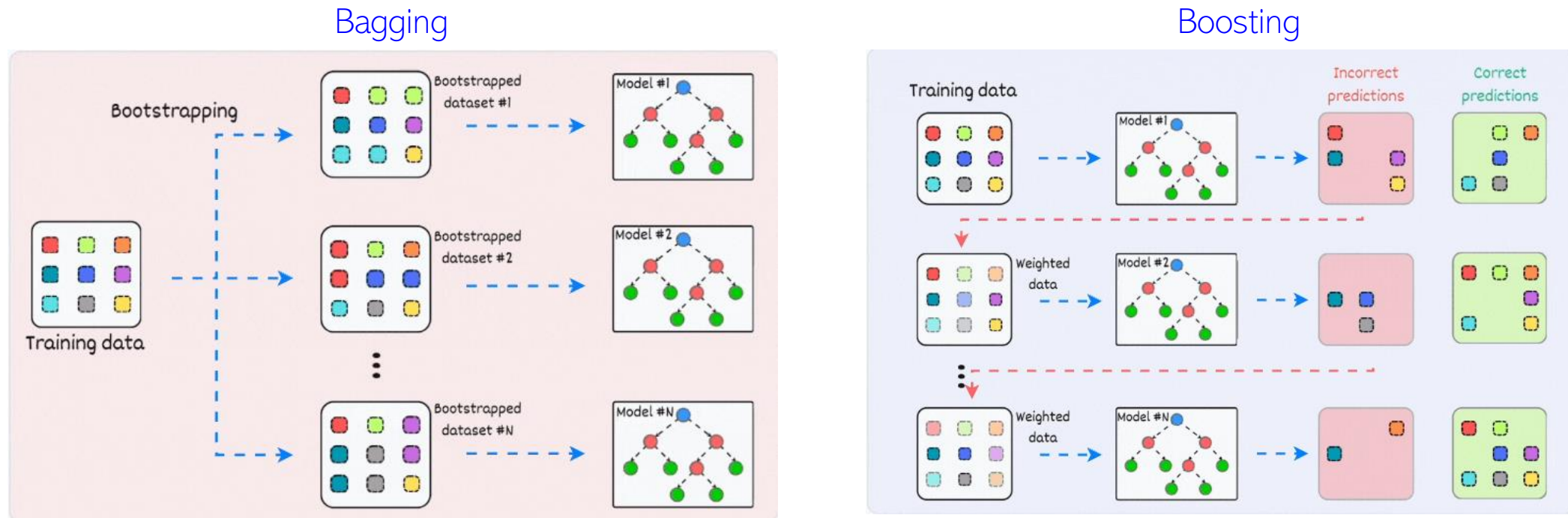


- MoE is very popular in classical ML as well as “modern” deep learning



Ensembles using Bagging and Boosting

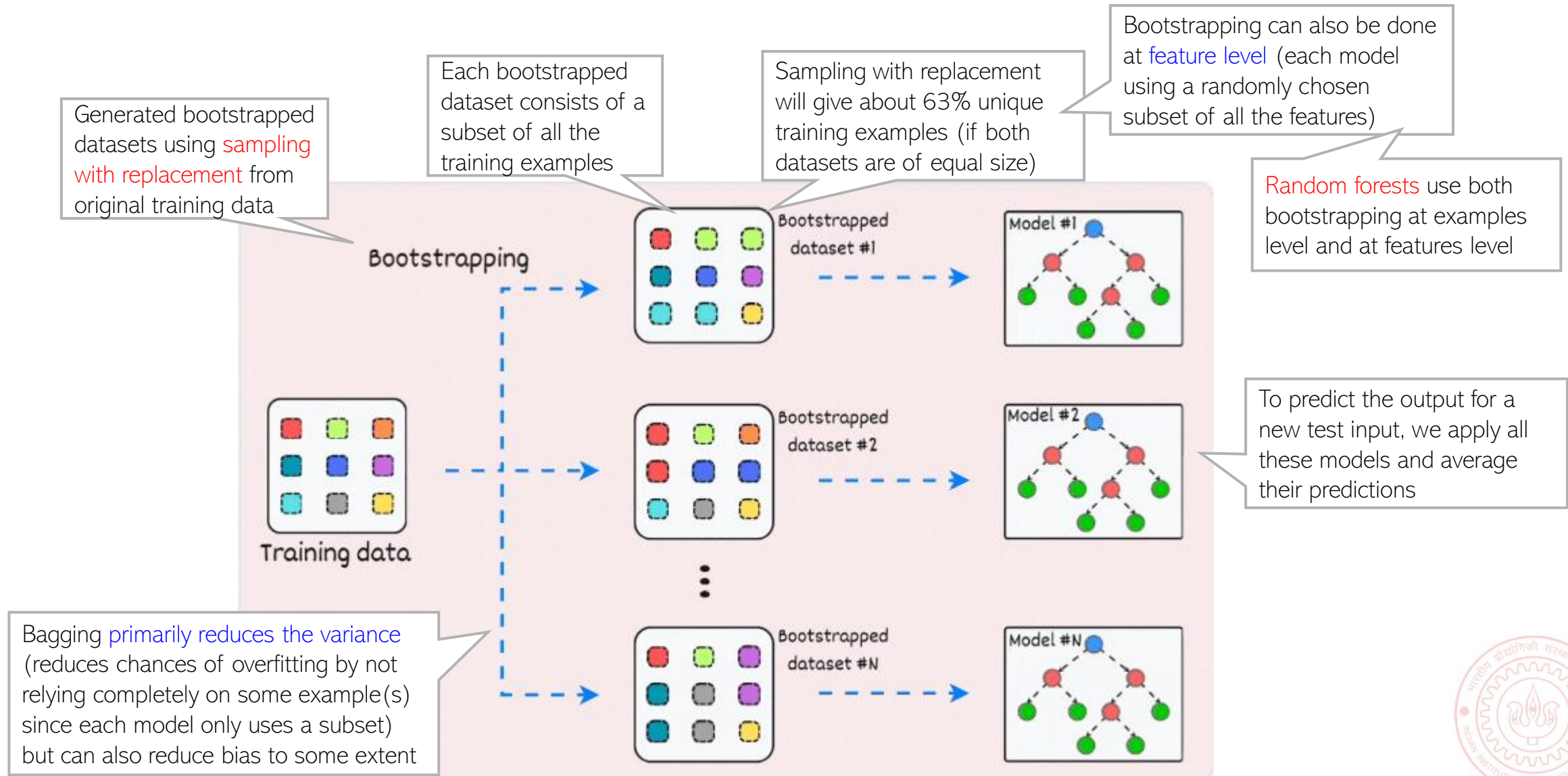
- Both use a single training set \mathcal{D} to learn an ensemble consisting of several models
- Both construct N datasets from the original training set \mathcal{D} and learn N models



- Bagging can do this in parallel for all the N models
- Boosting requires a sequential approach for N rounds



Bagging (Bootstrap Aggregation)



Boosting

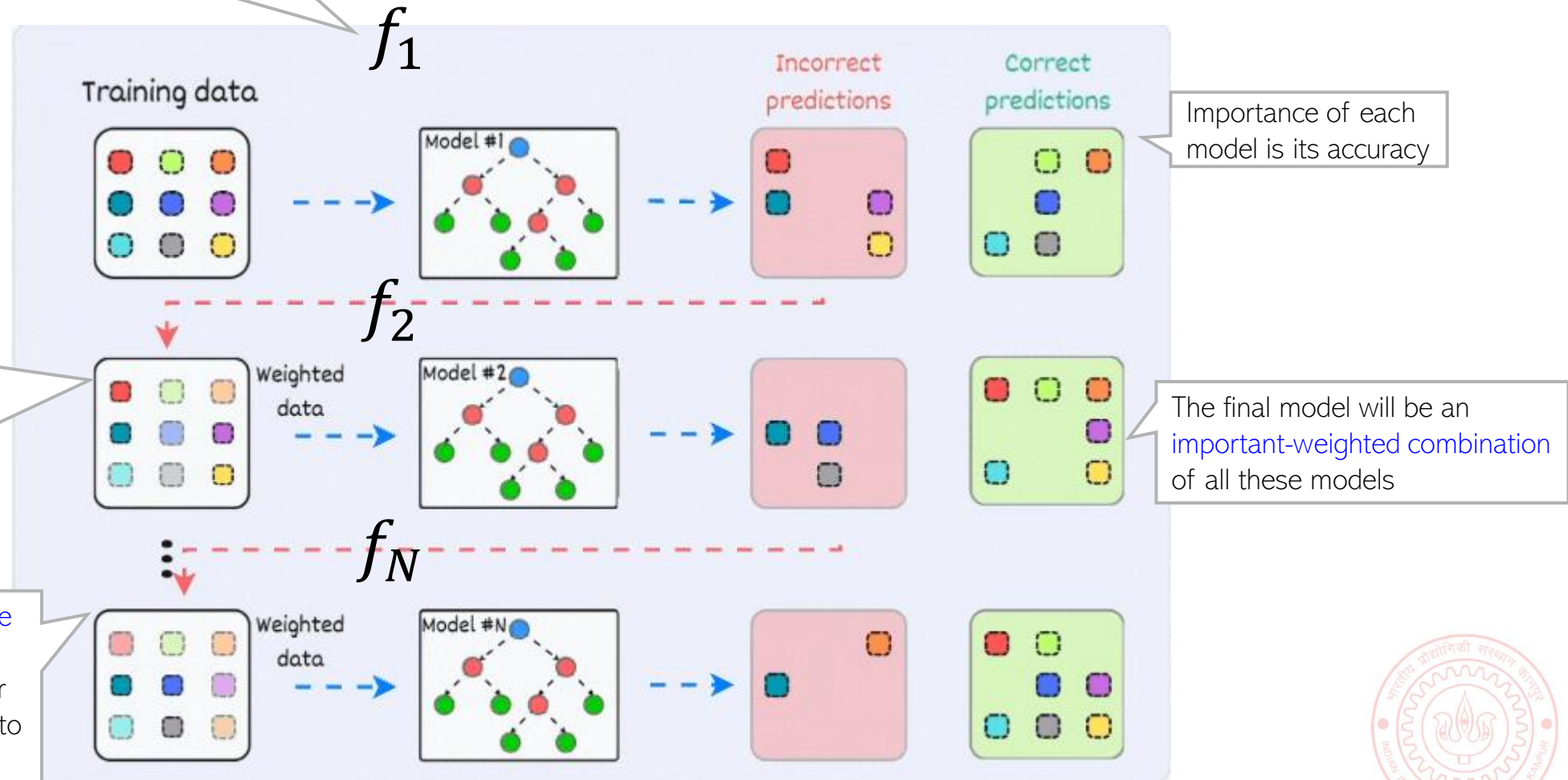
Boosting assumes that the individual models are simple/weak models that can be easily learned

Boosting trains them sequentially and combines them to get a “boosted” powerful model

Note that here we have two types of importances: importance of each training example and importance of each model

“Weighted data” means that we are increasing the importance of examples that were mis-predicted in the previous round and decrease it for examples that were correctly predicted

Boosting primarily reduces the bias by making the weak (underfitted) models stronger but can also reduce variance to some extent



A Boosting Algo: AdaBoost (Adaptive Boosting)

- In many ML problems, we can assign importance weight to each example, e.g., by weighing each term in the loss functions, i.e., $\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \beta_i \ell(\mathbf{x}_i, y_i, \mathbf{w})$

Importance of the training example (\mathbf{x}_i, y_i)

We might know this beforehand or estimate it during training

- AdaBoost is based on optimizing such a loss function

Initially assume equal importance for all training examples

- Initialize the ensemble as $\mathcal{E} = \{\}$ and $\boldsymbol{\beta}$ as $\boldsymbol{\beta}^{(0)} = [\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}]$

- For round $t = 1, 2, \dots, T$

- $\mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N \beta_i^{(t-1)} \ell(\mathbf{x}_i, y_i, \mathbf{w})$ and add it to ensemble $\mathcal{E} = \{\mathcal{E} \cup \mathbf{w}^{(t)}\}$

- Define the **total loss** of $\mathbf{w}^{(t)}$ as $L(\mathbf{w}^{(t)}) = \sum_{i=1}^N \beta_i^{(t-1)} \ell(\mathbf{x}_i, y_i, \mathbf{w}^{(t)})$

Or the importance weighted total error

- Compute the **"importance"** of $\mathbf{w}^{(t)}$ for the \mathcal{E} as $\alpha_t = f(L(\mathbf{w}^{(t)}))$

f is some function such that α_t is high if total loss $L(\mathbf{w}^{(t)})$ is low, and vice-versa

- Increase/decrease importance β_i of each training instance (\mathbf{x}_i, y_i) for next round as

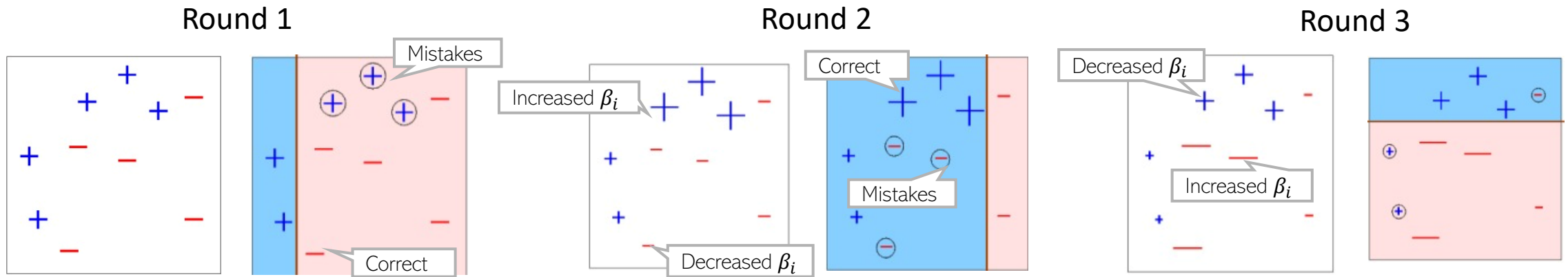
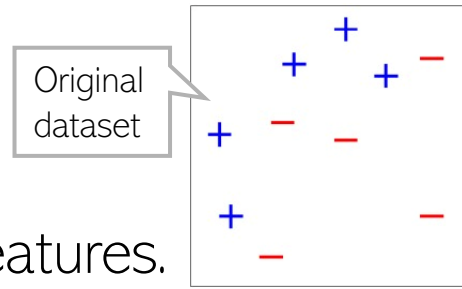
$$\beta_i^{(t)} \propto \begin{cases} \beta_i^{(t-1)} \times \exp(\alpha_t \ell(\mathbf{x}_i, y_i, \mathbf{w}^{(t)})) & \text{(Increase if } \mathbf{w}^{(t)} \text{ mispredicted } (\mathbf{x}_i, y_i)) \\ \beta_i^{(t-1)} \times \exp(-\alpha_t \ell(\mathbf{x}_i, y_i, \mathbf{w}^{(t)})) & \text{(Decrease if } \mathbf{w}^{(t)} \text{ correctly predicted on } (\mathbf{x}_i, y_i)) \end{cases}$$

- Final model is $\hat{\mathbf{w}} = \sum_{t=1}^T \alpha_t \mathbf{w}^{(t)}$ importance-weighted average of all $\mathbf{w}^{(t)}$'s

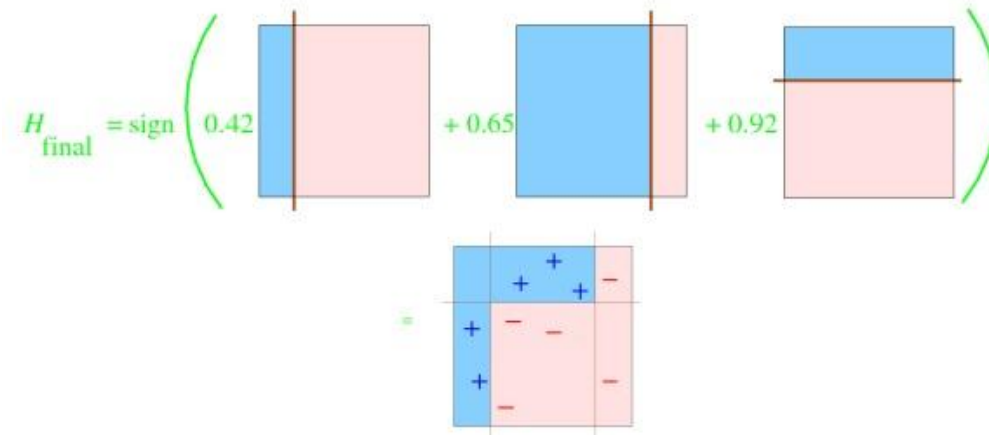


AdaBoost: An Illustration

- Suppose we have a binary classification problems with each input having 2 features.
- Suppose we have a weak model like a simple DT (decision stump)
- Illustration of AdaBoost using a decision stump if run for 3 rounds



- The ensemble represents the overall model
- We got a nonlinear model from 3 simple linear models
- Note that the ensemble was constructed sequentially



Gradient Boosting

- Consider learning a function $f(x)$ by minimizing a squared loss $\frac{1}{2}(y - f(x))^2$
- Gradient boosting is a sequential way to construct such $f(x)$
- For simplicity, assume we start with $f_0(x) = \frac{1}{N} \sum_{i=1}^N y_i$
- Given previously learned model $f_m(x)$, let's assume the following “improvement” to it

$$f_{m+1}(x) = f_m(x) + h(x)$$

“Residual” which, if added to $f_m(x)$, will make the new prediction $f_{m+1}(x)$ closer to y

- Thus the goal for the next round is to learn the “residual” $h(x) = y - f_m(x)$
- Residual is negative gradient of the loss w.r.t. $f(x)$ - thus called “gradient boosting”
- The final model $f_M(x)$, once the residual is sufficiently small, is what we will use
- The idea of gradient boosting is applicable to classification too
- **XGBoost** (eXtreme Gradient Boosting) is a very popular grad boosting algo

Based on sequentially constructing a DT

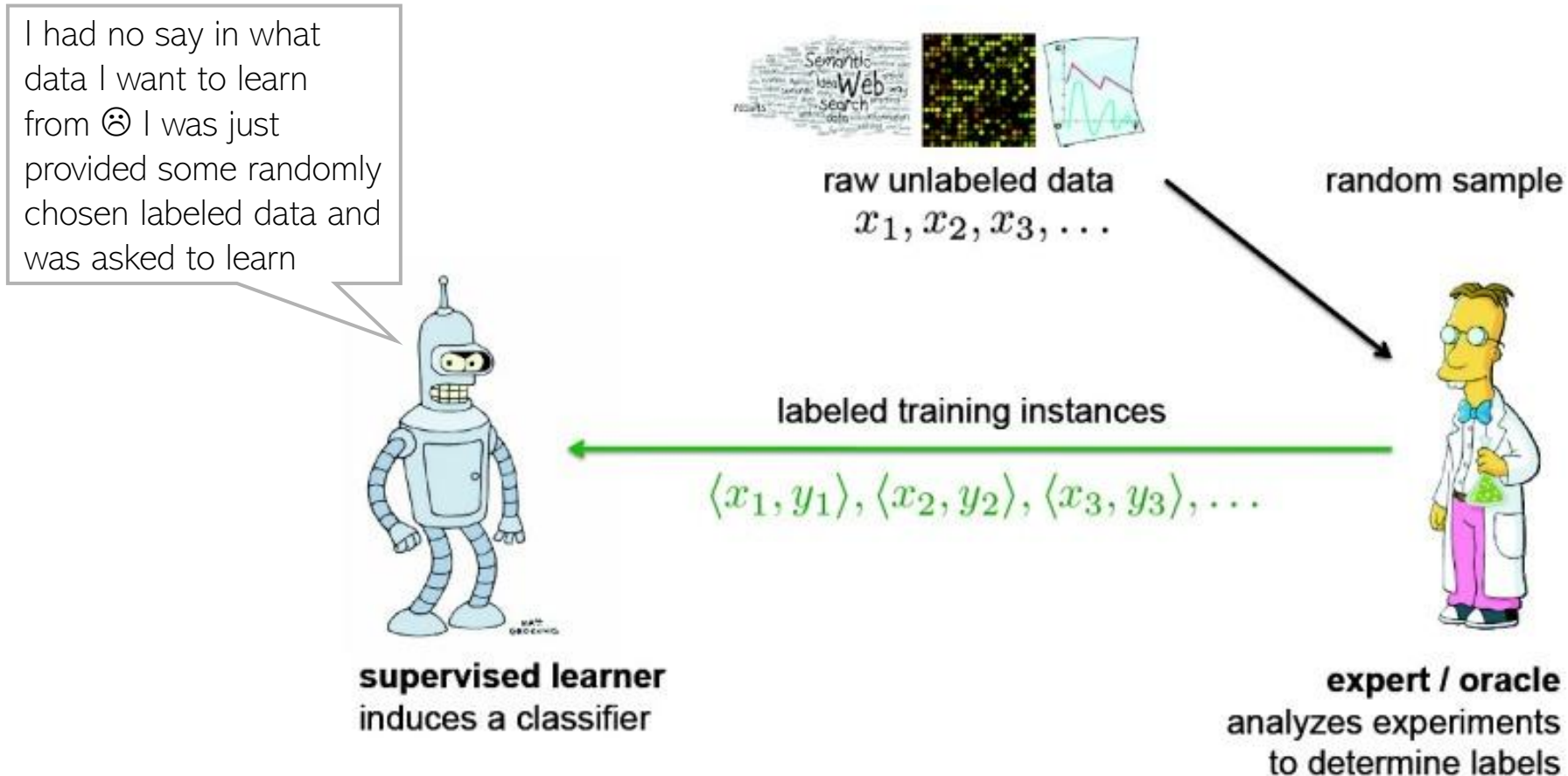
Active Learning

(an example of learning with human-in-the-loop)



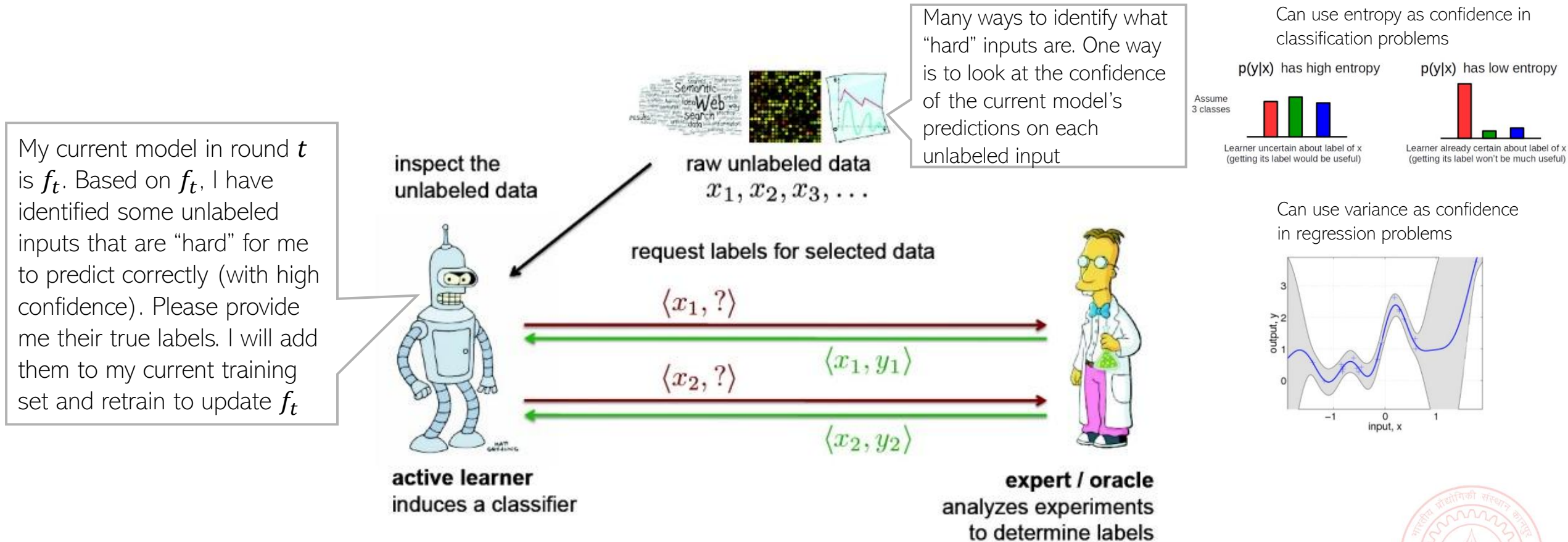
Active Learning

- Standard supervised learning is “passive” (learner has no control; we just give it data)
- We take a random sample of inputs, get them labelled by an expert, and train a model



Active Learning

- In active learning, learner can request what training examples it wants to learn from



Learning in the wild



Domain Adaptation

- We may have a “source” model trained on data from some domain
- We might want to deploy it in a new domain
- Performance of the source model will suffer
- To prevent this, we usually perform “domain adaptation” or “transfer learning”
- These are broad terms covering a variety of techniques that “finetune” the source model using labelled/unlabeled data from the new domain

We do expect some “commonality” (e.g., some common set of features) between the two domains otherwise we can’t hope to have any adaptation/transfer



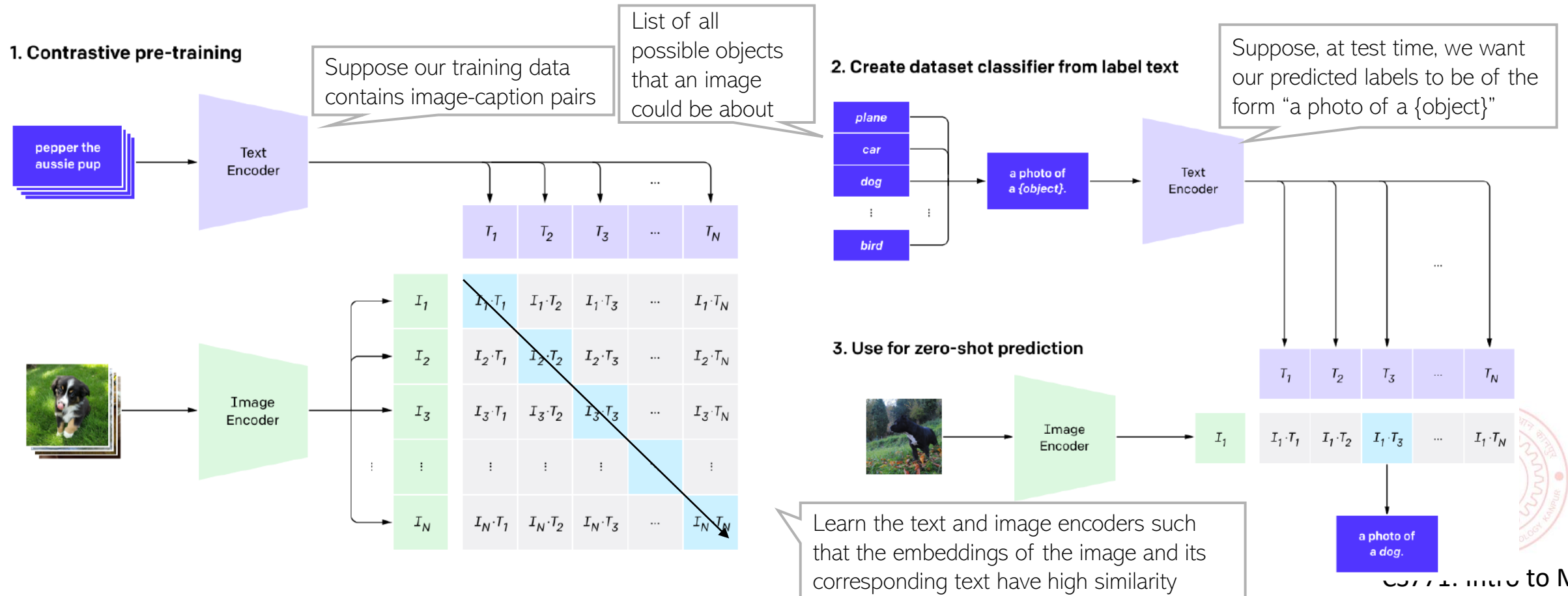
Zero-Shot Learning

25

Recall Homework 1 programming problem where we used some simple methods to solve it using attribute vector (an embedding) of each class

- What if our training data doesn't have the test data classes?
- Several methods to solve ZSL using deep learning. CLIP is a recent approach

CLIP: "Contrastive Language-Image Pre-training" (Radford et al, 2021)



The ending note..

- Good features are important for learning well
- The “classical” ML methods we studied in this course still continue to have high relevance
- Success of deep learning is largely attributed to (automatically learned) good features
- Deep learning is not a panacea – often simple classical models can do comparably/better
- First understand your data (plot/visualize/look at some statistics of the data, etc)
- Always start with a simple model that you understand well
 - Try to first understand if your data really needs a complex model
- Think carefully about your features, how you compute similarities, etc.
- Helps to learn to first diagnose a learning algorithm rather than trying new ones
 - Understanding of optimization algos, loss function, bias-variance trade-offs, etc is important
- No free lunch. No learning algorithm is “universally” good

