*To cite this presentation: Pendyala, V.S. (2022) "Meta-algorithms in Machine Learning". PPT Presentation. IEEE Computer Society, Santa Clara Valley Chapter Webinar.*

*Video Recording: https://ieeetv.ieee.org/video/meta-algorithms-in-machine-learning*

# Meta-algorithms in Machine Learning

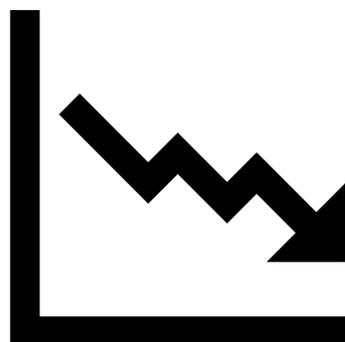Vishnu S. Pendyala, PhD

## How do we make the best of both?

PARAMETRIC MODELS LIKE LOGISTIC REGRESSION ARE IMPACTED BY BIAS

SOME OTHERS LIKE DECISION TREES ARE IMPACTED BY VARIANCE

How do you reduce the risk arising from fluctuations (variance) in a stock's price?
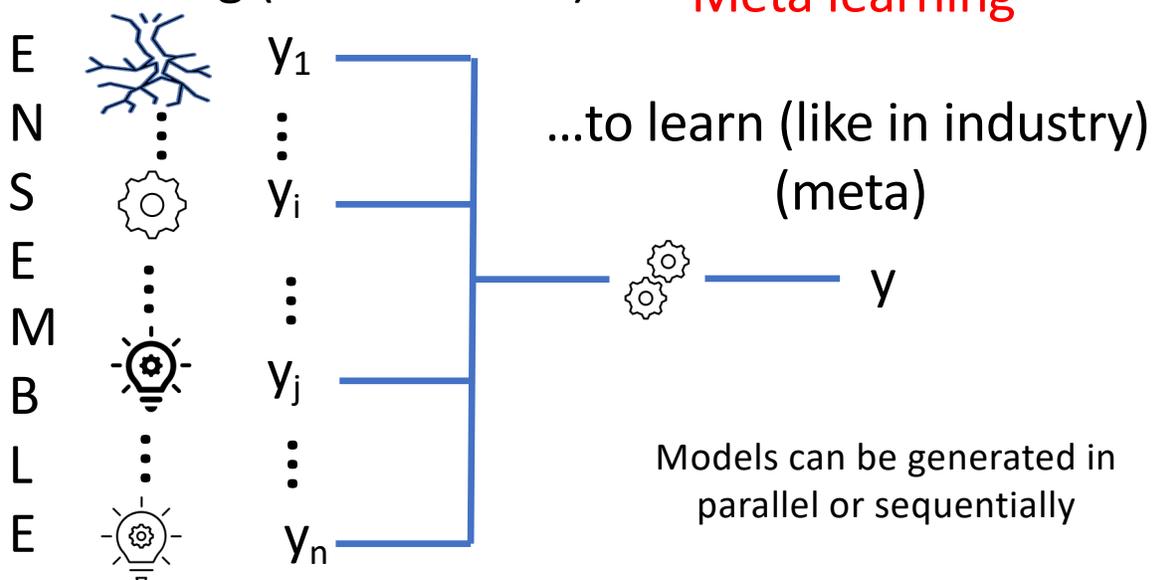
Diversify!

Learning (like in school)...

Meta learning

E
N
S
E
M
B
L
E

$y_1$

$y_i$

$y_j$

$y_n$

...to learn (like in industry) (meta)

$y$

Models can be generated in parallel or sequentially
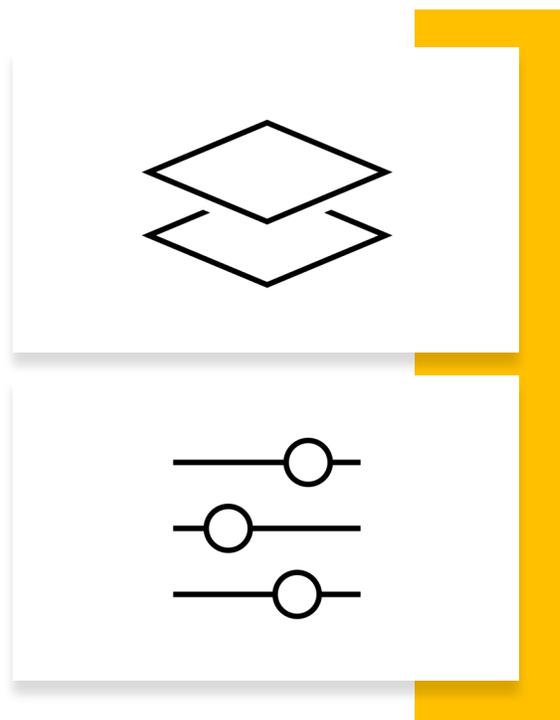
Base (often, weak) learners

## Meta Learning

- Learning is never perfect – incurs a loss
- Meta learning minimizes the losses from the previous round of learning
- Meta learning is a big deal for deep learning – few shot learning and more
- Deep learning itself is a kind of meta learning

## The rationale

- **Consider 256 weak learners for binary classification, $y_i \in \{-1, +1\}$**
- **Each uncorrelated classifier is a weak learner with error rate, say <u>0.45</u>**
- **For the ensemble to make a misclassification, majority (in this case 129 or more) base learners must misclassify**
- **This is like the coin-toss experiment, so we use binomial distribution for the probability of the ensemble making a misclassification**

  **P = P$_{129}$ + P$_{130}$ + … + P$_{256}$ where P$_i$ = $\binom{256}{i} p^i q^{256-i}$ p = 0.45 and q = 0.55**

  **<u>Cumulative probability</u>, P = $\sum_{i=129}^{256} \binom{256}{i} p^i q^{256-i}$ = <u>0.04</u>**

## Learning…

Bagging, Adaboost

$y_1$

$y_i$

$y_j$
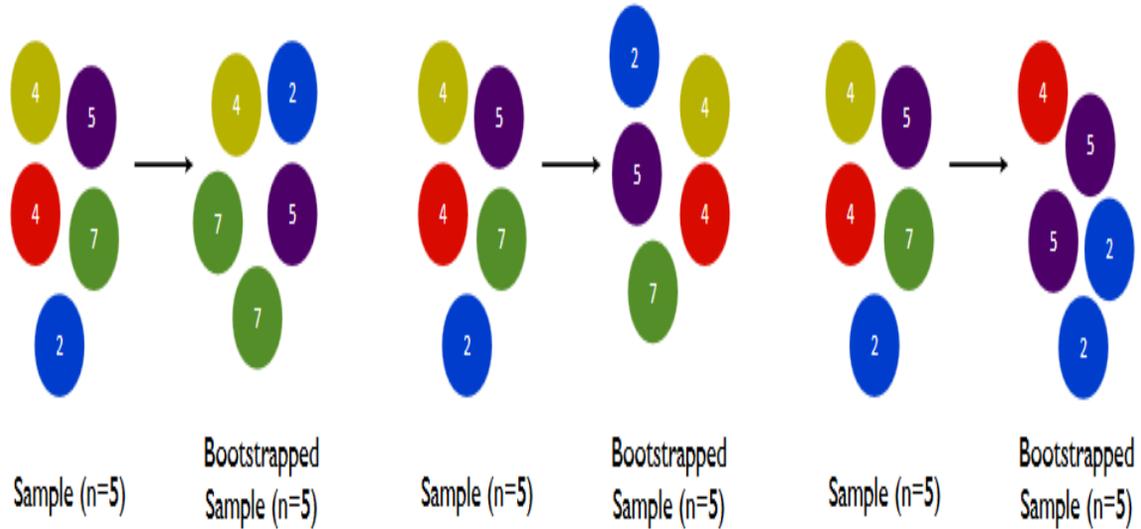
$y_n$

…to learn (meta)

$y$

But to diversify, we do not want to generate the same tree!

## How can we generate different decision trees from the same training dataset?

- Perturb X or Perturb Y
- Perturb X in two ways: row-wise or column-wise – randomly, no bias!
- Perturbation of X can be via
  - bootstrap sampling
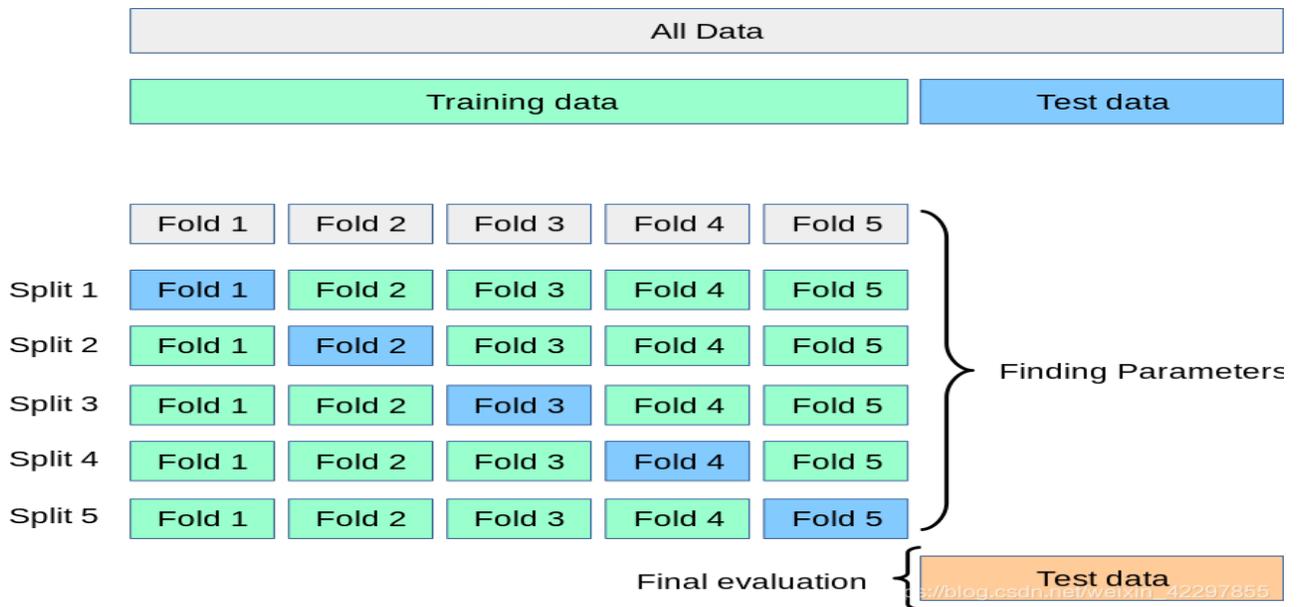  - k-fold sampling
  - weighted sampling
  - random subspaces

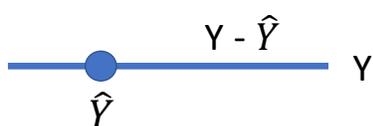| $x_1$ | $x_2$ | … | $x_n$ | Y |
|-------|-------|-----|-------|---|
|       |       |     |       |   |
|       |       |     |       |   |
|       |       |     |       |   |
|       |       |     |       |   |

# Bootstrap Sampling

# K-fold Sampling for Cross Validation

# How can we perturb Y?

- What aspect of Y typically changes with each iteration on a model in an ensemble?
  - The predicted value $\hat{Y}$
- Suppose we need to reach a target Y taking several steps (iterations).
- We can know where we are, as an absolute value $\hat{Y}$ or as a measure for how far (Y - $\hat{Y}$) we are from the target Y
- What if we use (Y - $\hat{Y}$) as the target variable instead of Y?
- Each iteration, we get new values for the target variable and Y is perturbed => we get a different model each time.
- Instead of trying to predict Y, we predict how far we are from Y; i.e. we fit the residuals instead of the target value.

Y - $\hat{Y}$

$\hat{Y}$

Y

E
N
S
E
M
B
L
E

$y_1$
$y_i$
$y_j$
$y_n$

## Next Question: How can we aggregate the models?

y

- Majority vote (mode)
- Average (mean)
- Weighted response
- Metamodel

## Dataset Perturbations and response aggregation for model ensembling

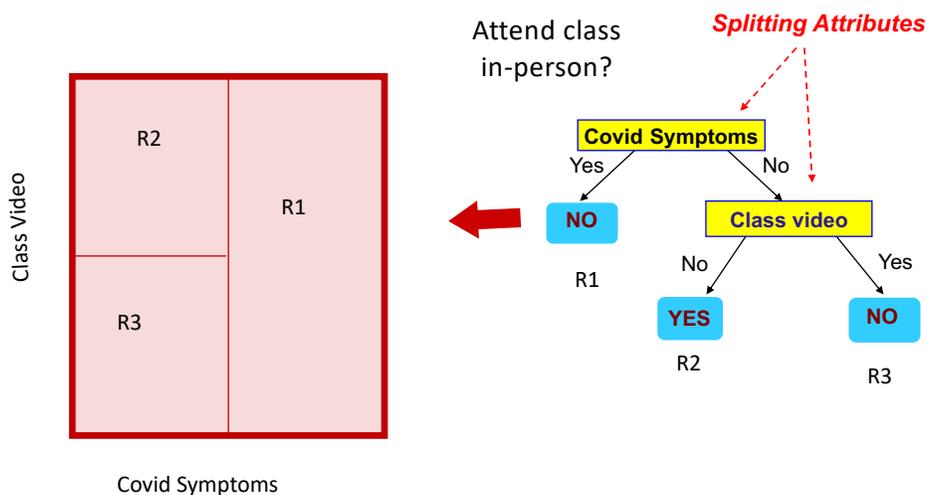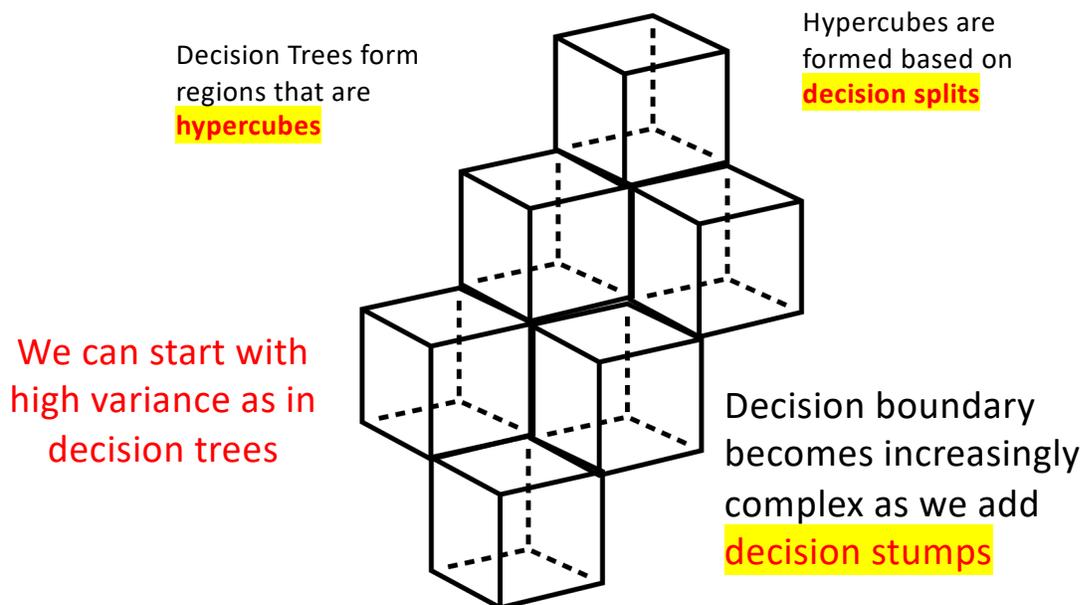| Ensemble Method | Perturb X row-wise using | Perturb X Column-wise using | Perturb Y using | Model generation | Aggregation strategy |
|---|---|---|---|---|---|
| Bagging | Bootstrap Sampling | None | None | In parallel | Mean or mode |
| Random Forest | Bootstrap Sampling | Random subspaces (at each tree/node) | None | In parallel | Mean or mode |
| Adaboost | Bootstrap Sampling with weighting | None | None | Sequential | Weighted response |
| Gradient Boosting | None | None | Pseudo-residuals | Sequential | Weighted response |
| Stacking | None (or) K-fold sampling | None | None | In parallel | Metamodel |

Next question: Where do we start?

High bias or high variance?

# Decision tree divides the instance space into regions
## High Variance, Complex Decision boundary

Attend class in-person?

**Splitting Attributes**

Class Video

R2

R1

R3

Covid Symptoms

**Covid Symptoms**

Yes        No

**NO**

R1

**Class video**

No                    Yes

**YES**

R2

**NO**

R3

Complex Decision Boundary

Decision Trees form regions that are **hypercubes**

Hypercubes are formed based on **decision splits**

We can start with high variance as in decision trees

Decision boundary becomes increasingly complex as we add **decision stumps**

# High Complexity of the Decision Boundary



income

# Low Complexity of the Decision Boundary

Y = Preferred Customer Status

## A Decision Stump



Income > 60k

Preferred        Not Preferred

**OR we can start with high bias as in a decision stump**



Preferred

Not Preferred

Considering just one feature out of many

income

## Two Approaches to Ensembles: Boosting (Sequential) and Bagging (Parallel)

We start here

We start here

High Bias
Underfitting

Boosting

Not impacted by high variance even when the complexity increases!

Test Error

Model Complexity (decision boundary)

High Variance
Overfitting

Not impacted by high bias even when the complexity decreases!

Bagging

Test Error

Model Complexity (decision boundary)

# Boosting

## Aggregating the model responses using weights

# Boosting: Baby steps to loss minimization

## Error, loss function, and the initial prediction

- Baby steps => errors at each step. How do we model the error or for the entire dataset, loss?
- The simplest error is the residual, r = (Y - $\hat{Y}$) called the residual
- This is same as the negative gradient of the popular loss function, the squared error, L = (Y - $\hat{Y}$)$^2$ ; r = - $\frac{1}{2}$ ($\frac{\partial L}{\partial \hat{Y}}$)
- To be agnostic to the loss function, ($\frac{\partial L}{\partial \hat{Y}}$), which may not always be in the form of a residual (Y - $\hat{Y}$) , is called pseudo-residual.
- For the loss to be minimum, ($\frac{\partial L}{\partial \hat{Y}}$) = 0
- For squared error, $\sum_{i=1}^{N}$ (Yi − $\hat{Y}_i$) = 0 => if we have to start with a good estimate for $\hat{Y}_i$ for all data items, N* $\hat{Y}_i$ = $\sum_{i=1}^{N}$ Yi
- First baby step is a single leaf with the mean of target values – not even a stump!
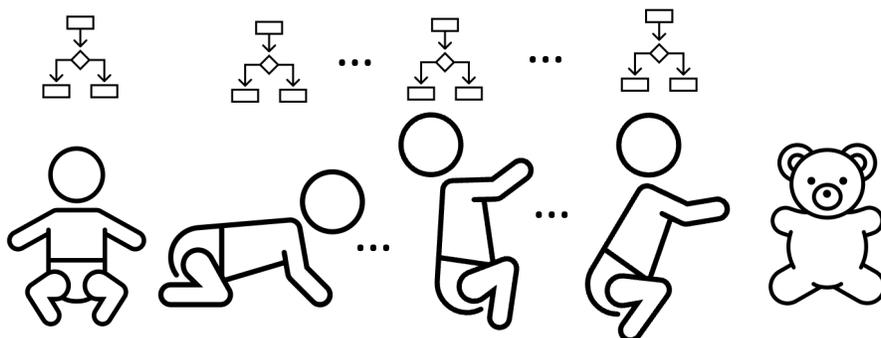
# Baby steps summing up to the target

- Residuals tell how bad the base learners are and how far the base learner is from the target variable.

- How can we bridge the gap (residual) left by the base learner?

- Why not generate a series of models $h_b(x)$ that try to bridge the gap, by predicting not the final target, but the subsequent residuals?

- The training dataset for $h_b(x)$ at each step is not $\{(x_i, y_i)\}$ but $\{(x_i, y_i - \hat{y}_i)\}$

- Then, $\widehat{Y} = H(\mathbf{x}) = \sum_{b=1}^{B} \lambda h_b(x)$ where $\lambda$ is the regularization parameter to slow the learning process and avoid overfitting.

$$h_0(x) \ + \ \lambda h_1(x) \ + \ \lambda h_2(x) \ + \cdots + \lambda h_i(x) \cdots + \lambda h_B(x) \qquad = H(x)$$

$$r_0 \qquad r_1 = r_0 - \lambda h_1(x) \qquad \cdots \qquad r_i = r_{i-1} - \lambda h_i(x) \qquad \cdots$$

## Boosting: Baby steps to loss minimization
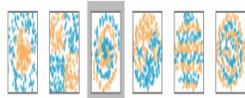
# Works for any differentiable loss function

replace r with the
pseudo-residual = $(\frac{\partial L}{\partial \widehat{Y}})$
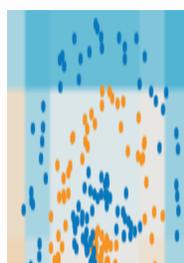
## Gradient Boosting Interactive Playground

Jul 5, 2016 • Alex Rogozhnikov •
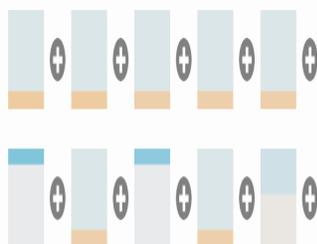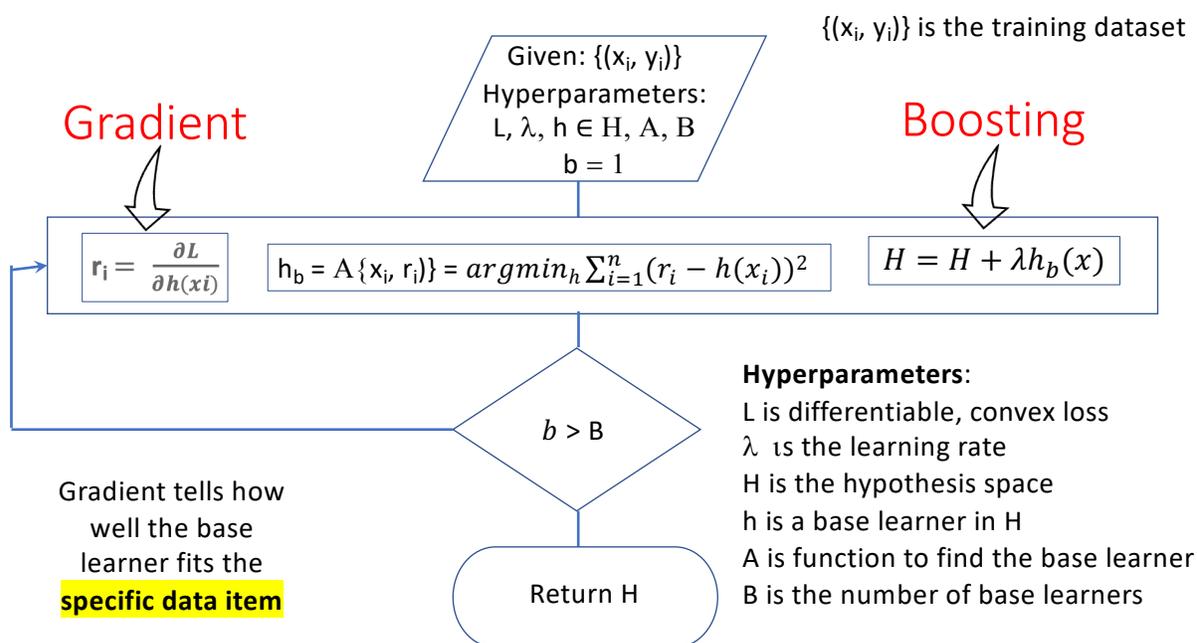
Dataset to classify:

Prediction:          Decision functions of first 30 trees
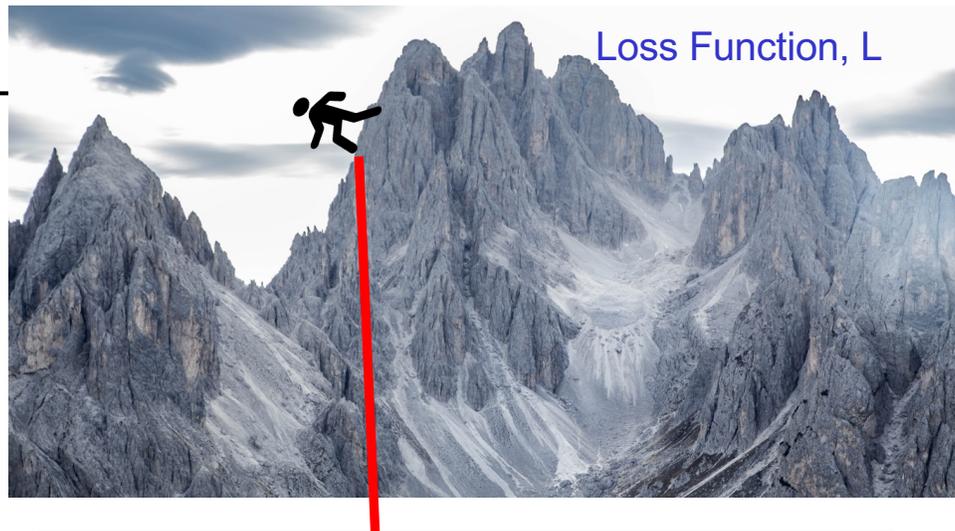
## Some notes

The base learner model, $h_i(x)$ is typically a decision tree as well and to regularize further, can just be a decision stump (depth=1).

In some sense the pseudo-residual indicates how difficult it is to fit the item using the base learner – can be thought of as a relative weight of the data item

{($x_i$, $y_i$)} is the training dataset

Gradient

Boosting

Given: {($x_i$, $y_i$)}
Hyperparameters:
L, $\lambda$, h $\in$ H, A, B
b = 1

$r_i = \dfrac{\partial L}{\partial h(xi)}$   $h_b$ = A{$x_i$, $r_i$)} = $argmin_h \sum_{i=1}^{n}(r_i - h(x_i))^2$   $H = H + \lambda h_b(x)$

$b > B$

Gradient tells how
well the base
learner fits the
**specific data item**

Return H

**Hyperparameters**:
L is differentiable, convex loss
$\lambda$ is the learning rate
H is the hypothesis space
h is a base learner in H
A is function to find the base learner
B is the number of base learners

# Boosting is yet another Gradient Descent!

Source: Mason, L., Baxter, J., Bartlett, P., & Frean, M. (1999). Boosting
algorithms as gradient descent. *Advances in neural information
processing systems*, *12*.

Loss Function, L

w ⟶ Slope is positive

Hiker must go left on the horizontal axis (decrease w)



Loss Function, L

w ⟶ Slope is negative

Hiker must go right on the horizontal axis (increase w)

# How do we express the last two slides in one line in math?
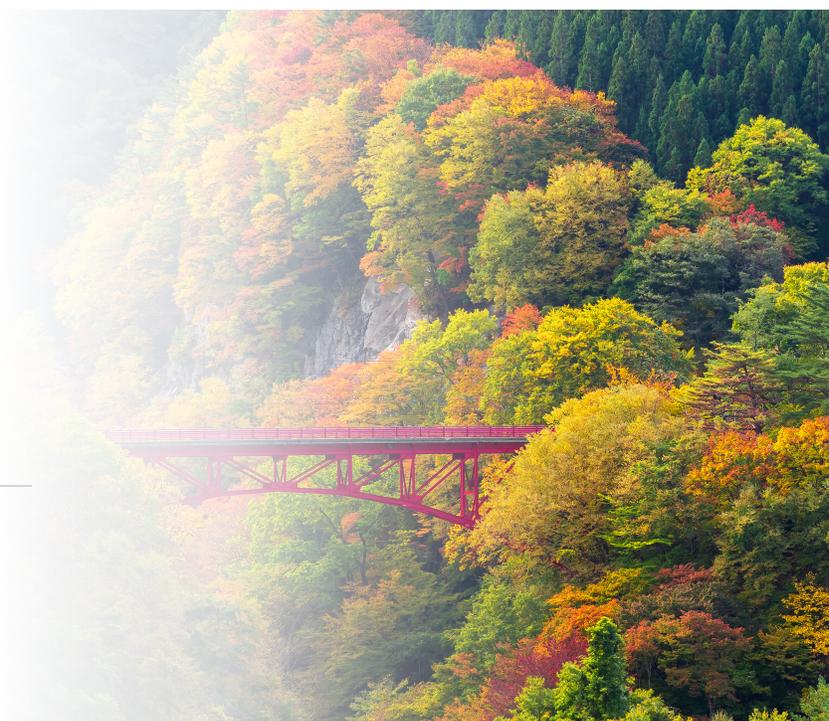
$$w_i = w_i - \lambda \frac{dL}{dw_i}$$

- Derivative is the ==slope== of the loss function
- $\lambda$ is the ==length of the strid==e in the direction of the slope
- If ==slope is positive, weights decrease== and vice versa
- If we extend it to ==multidimensions==, we use ==gradient== instead of derivative

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Adaboost

Adaptive boosting

# Adaboost vs Gradient boosting

- Both take baby steps; base learners are typically decision stumps
- Instead of shrinking the trees using a $\lambda$ that is constant for all the decision trees, in Adaboost, we use a varying $\alpha$ that is proportional to how well the base learner performs.

$$\widehat{Y} = H(\mathbf{x}) = \mathbf{sign}(\sum_{b=1}^{B} \alpha_b h_b(x))$$

- The gradient, which indicated the relative importance of a data item is now used to perturb X row-wise (Adaboost) instead of Y (GB)
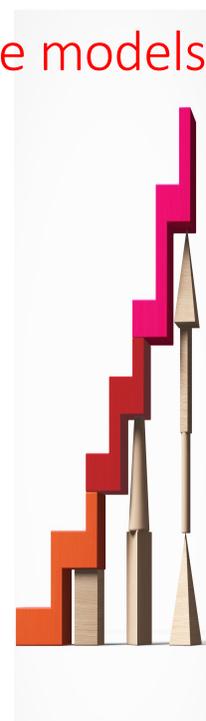
## Boosting for Classification: Weighting the base models

- What is the odds ratio of a model classifying correctly with probability p?

$$\frac{p}{1-p} = \frac{1-\epsilon}{\epsilon}$$ => better the model, higher the odds ratio

- Taking logarithm of ratios simplifies operations – multiplications convert to additions, divisions to subtractions

- The logit function, $\ln(\frac{1-\epsilon}{\epsilon})$ converts a probability p $\in (0,1)$ to a number r $\in (-\infty, +\infty)$ and = 2*arctanh(1-2$\epsilon$)

- $\ln(\frac{1-\epsilon}{\epsilon})$ is an indication of how well the base learner can classify, so can be used as a weight, $\alpha_b$ for the base learner

$$\alpha_b = \frac{1}{2} \ln\left(\frac{1-\epsilon}{\epsilon}\right) = \text{arctanh}(1-2\epsilon)$$

# Boosting for Classification: Weighting the data items

- Initially, all samples have the same weight: $\frac{1}{N}$

- We factor in $\alpha_b$ into the subsequent weight updates; $\alpha_b$ is a logarithm, so we use exponentiation:

$$w_i = w_i * e^{-\alpha_b y_i h(x_i)}$$
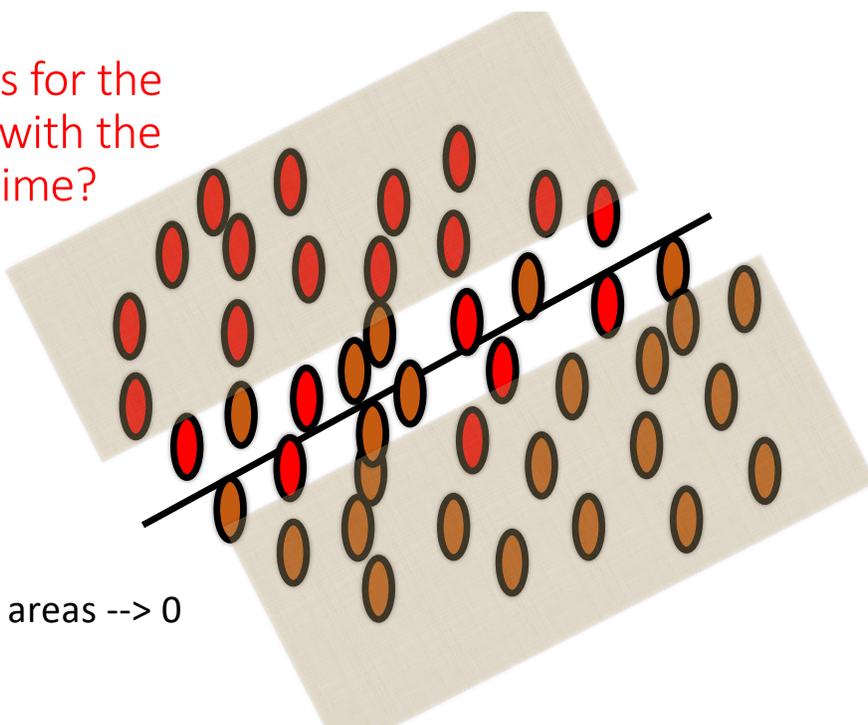
$\Rightarrow$ The weight for correctly classified data items, $e^{-\alpha_b}$ is exponentially low and for difficult ones, $e^{\alpha_b}$ is exponentially high

- Weights need to be normalized so that they sum up to 1 and become a probability distribution.

$$\text{Error } \epsilon = \frac{w_i * I(y_i \neq h_b(x_i))}{\sum_{i=1}^{N} w_i}$$

How do the weights for the data items change with the progression of time?

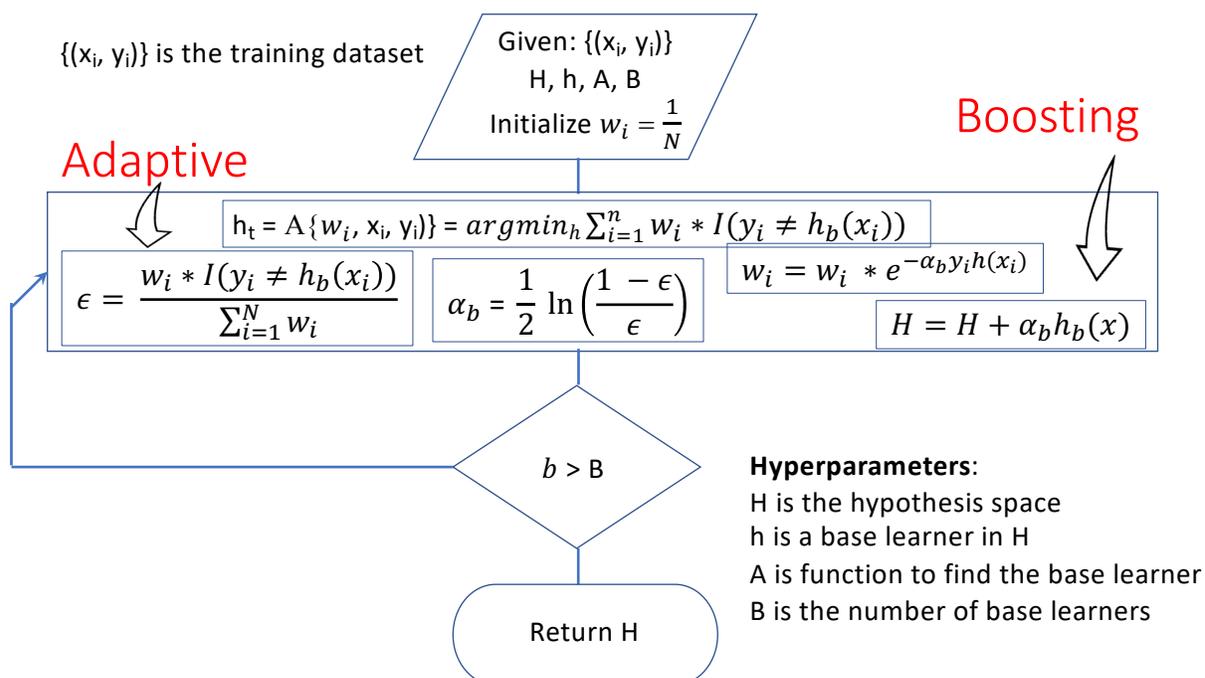For simplicity, here we are assuming a linear decision boundary
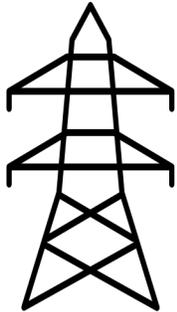
Weights in the shaded areas --> 0

## Loss function

- The loss function we are minimizing corresponding to $\alpha_b = \frac{1}{2} \ln\left(\frac{1-\epsilon}{\epsilon}\right)$ and $w_i = w_i * e^{-\alpha_b y_i h(x_i)}$ is exponential:

$\mathcal{L} = \sum_{i=1}^{N} e^{-y_i H(x_i)}$ => Gradient, $\left(\frac{\partial L}{\partial H(x_i)}\right) = -y_i\, e^{-y_i H(x_i)}$

- This time, we use the gradient to perturb X instead of Y and incorporate it into the weight update of the data items

{($x_i$, $y_i$)} is the training dataset

Given: {($x_i$, $y_i$)}
H, h, A, B
Initialize $w_i = \frac{1}{N}$

**Boosting**

**Adaptive**

$h_t = A\{w_i, x_i, y_i\}\} = argmin_h \sum_{i=1}^{n} w_i * I(y_i \neq h_b(x_i))$

$\epsilon = \dfrac{w_i * I(y_i \neq h_b(x_i))}{\sum_{i=1}^{N} w_i}$

$\alpha_b = \dfrac{1}{2} \ln\left(\dfrac{1-\epsilon}{\epsilon}\right)$

$w_i = w_i * e^{-\alpha_b y_i h(x_i)}$

$H = H + \alpha_b h_b(x)$

$b > B$

Return H

**Hyperparameters**:
H is the hypothesis space
h is a base learner in H
A is function to find the base learner
B is the number of base learners
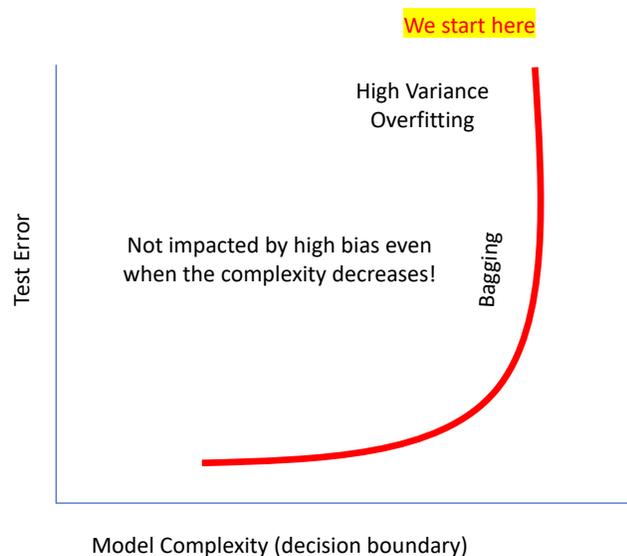
## Many more variants of boosting

- XGBoost, LightGBM, CatBoost are the popular ones
- Combine several techniques for scaling to big data, faster processing, handling categorical variables, missing values, textual data, etc
- Differ in tree generation, tree types, community support, hyperparameters, and naturally, performance

# Bagging – starting with high variance
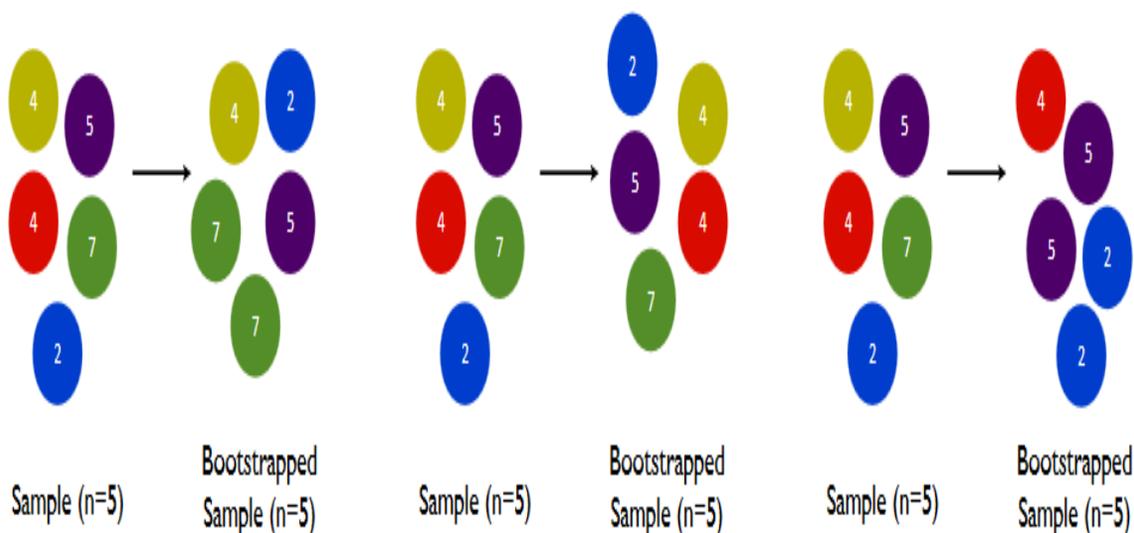
# Bagging - **B**ootstrap **Agg**regation

- Perturb X by bootstrap sampling
- Uses homogeneous base learners, typically decision trees
- No pruning – each tree is grown fully
- Tree growth can be easily parallelized
- Aggregation is by taking mean (regression) or mode (classification)
- Bootstrap samples generally leave out $1/e$ (=$\lim_{n \to \infty}(1-1/n)^n$) of the population
- Such samples are considered OOB
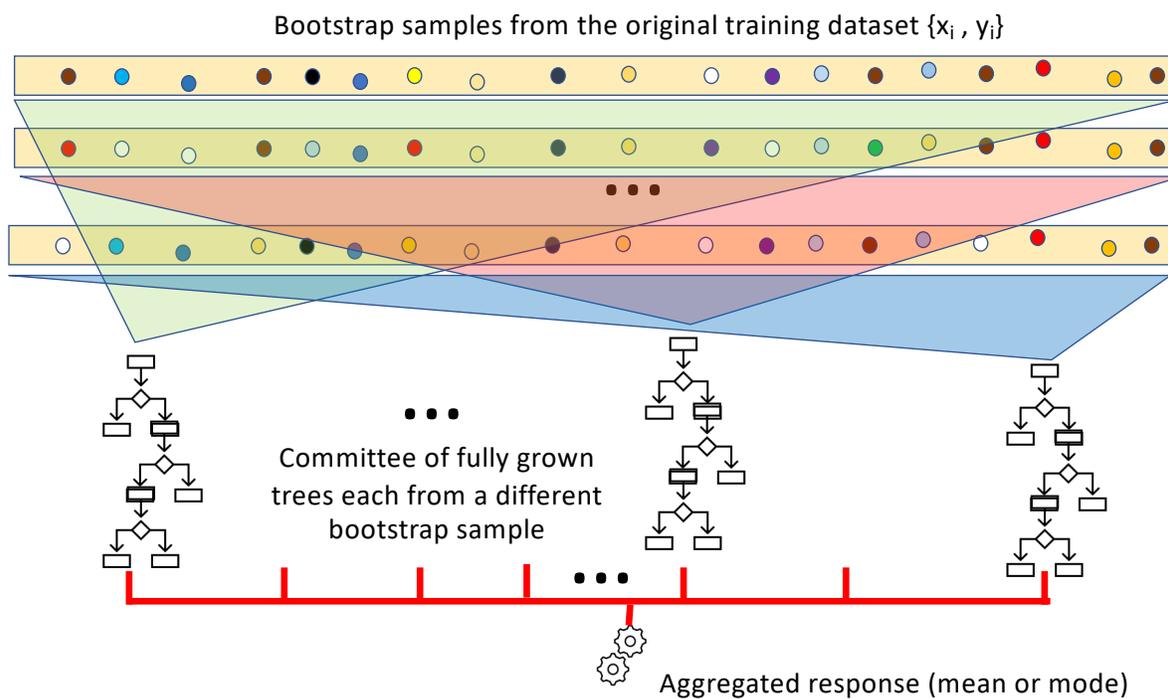- OOB samples automatically provide the validation dataset – no need for train-test split

We start here

High Variance
Overfitting

Not impacted by high bias even when the complexity decreases!

Test Error

Bagging

Model Complexity (decision boundary)

# Bagging uses Bootstrap Sampling



Sample (n=5) → Bootstrapped Sample (n=5)    Sample (n=5) → Bootstrapped Sample (n=5)    Sample (n=5) → Bootstrapped Sample (n=5)

Bootstrap samples from the original training dataset {$x_i$ , $y_i$}



Committee of fully grown trees each from a different bootstrap sample
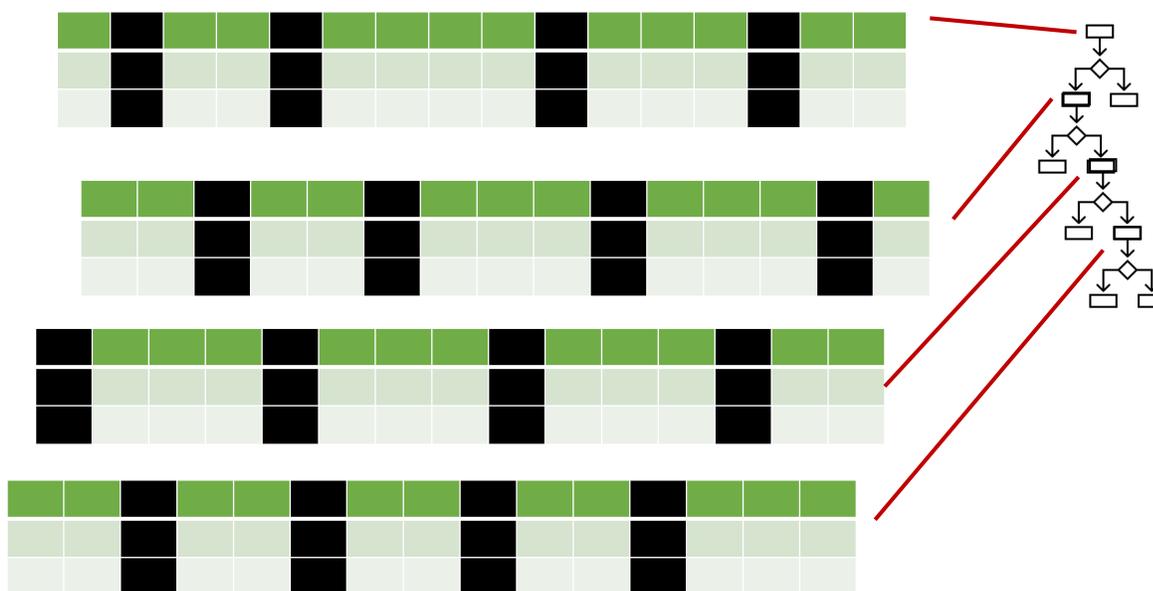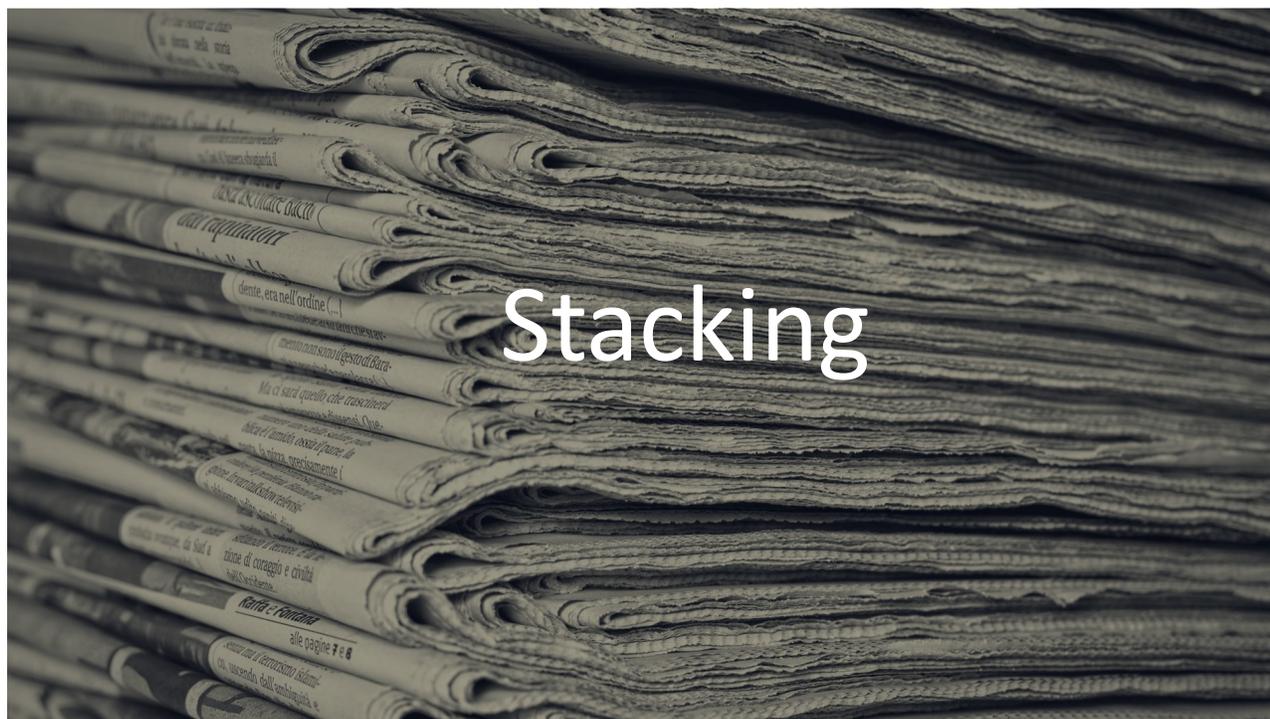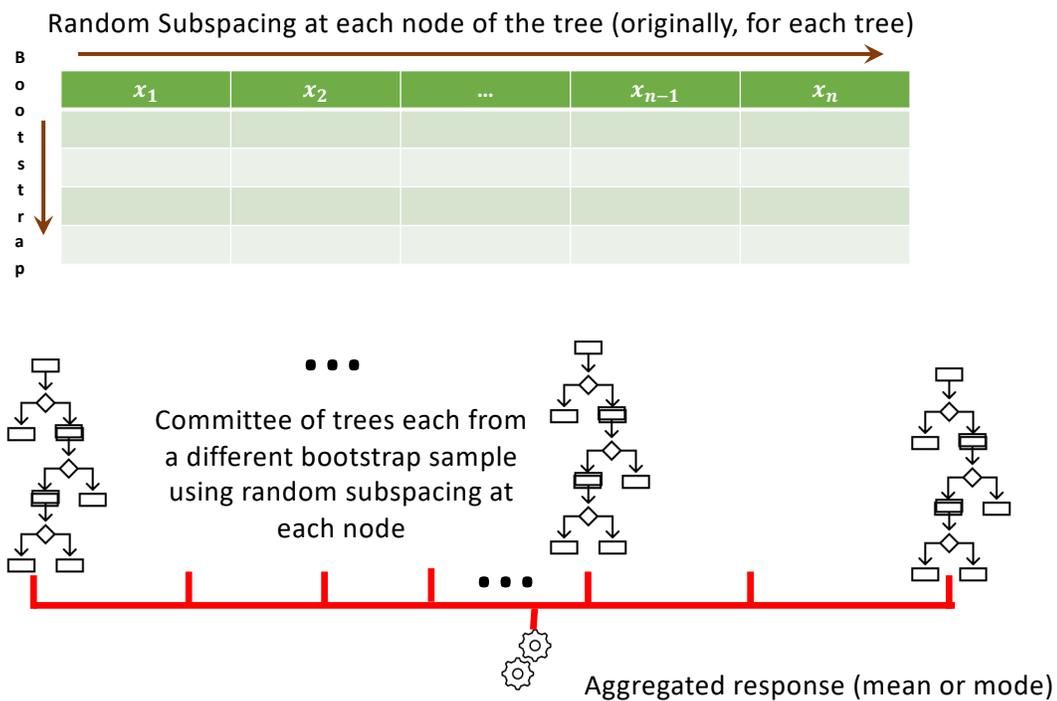
Aggregated response (mean or mode)

## Bagging: Problem

- Bootstrap sampling => Generated trees are identically distributed

- The relative importance of the features are preserved across the trees => trees are correlated => trees are i.d but not i.i.d

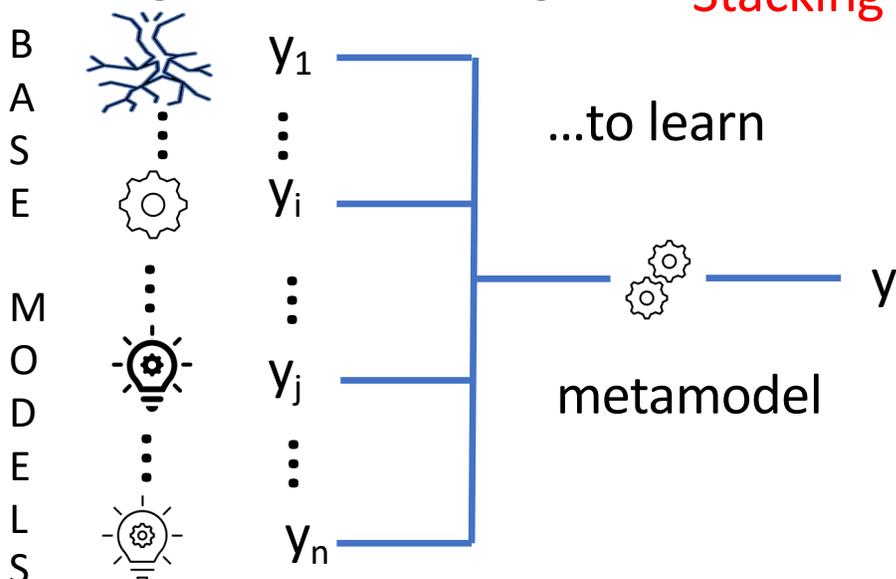- How do decorrelate trees and produce i.i.d trees?

# Random Forests

One more perturbation: Random subspacing

Random Subspacing at each node of the tree (originally, for each tree)

| | $x_1$ | $x_2$ | ... | $x_{n-1}$ | $x_n$ |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

B
o
o
t
s
t
r
a
p

Committee of trees each from a different bootstrap sample using random subspacing at each node

Aggregated response (mean or mode)

# Stacking

# Heterogeneous Learning…

## Stacking

B
A
S
E

$y_1$

$y_i$

M
O
D
E
L
S

$y_j$

$y_n$

…to learn

y

metamodel

## Stacking in a Nutshell

Original Dataset $\{x_i, y_i\}$ without any perturbations

| $x_1$ | $x_2$ | ... | $x_{n-1}$ | $x_n$ | $y$ |
|-------|-------|-----|-----------|-------|-----|
|       |       |     |           |       |     |
|       |       |     |           |       |     |

Base Learners $h_1 h_2 .. h_B$

Base Learner predictions are the new features; y remains the same

| $h_1$ | $h_2$ | ... | $h_B$ | $y$ |
|-------|-------|-----|-------|-----|
|       |       |     |       |     |
|       |       |     |       |     |

- Same data set => Base Learners are heterogeneous, so that they do not all learn the same
- No randomization of any sort of $\{x_i, y_i\}$ => Chances of overfitting are high
- Solution: Use different data sets for the base learners using k-fold sampling

**Final Prediction**

Meta Learner

# K-fold Sampling for Cross Validation

First fold of the training Dataset used to train base learners

**Step 1**

Base Learners $h_1 h_2 .. h_B$

Holdout Dataset (Second fold) used to train the meta learner

**Step 2**

Base Learners $h_1 h_2 .. h_B$

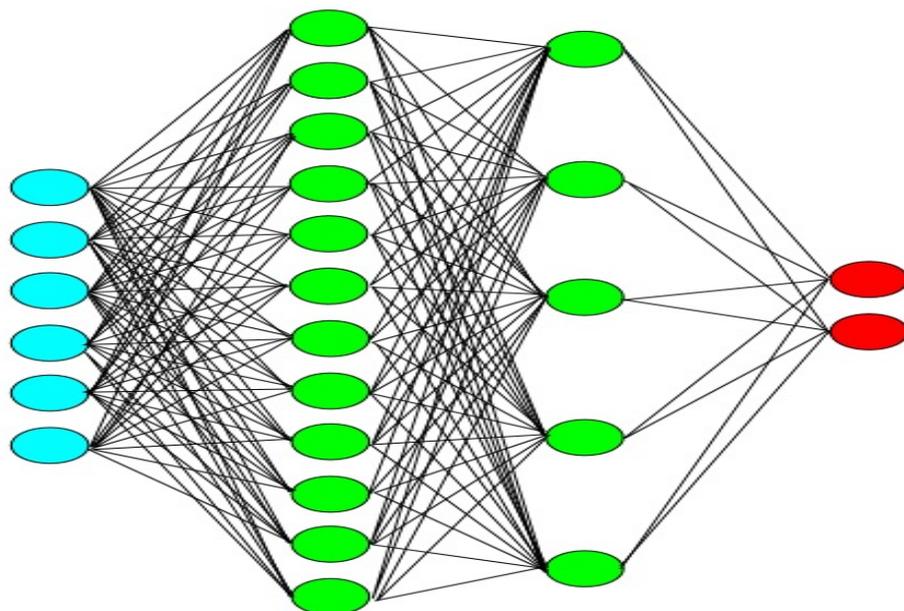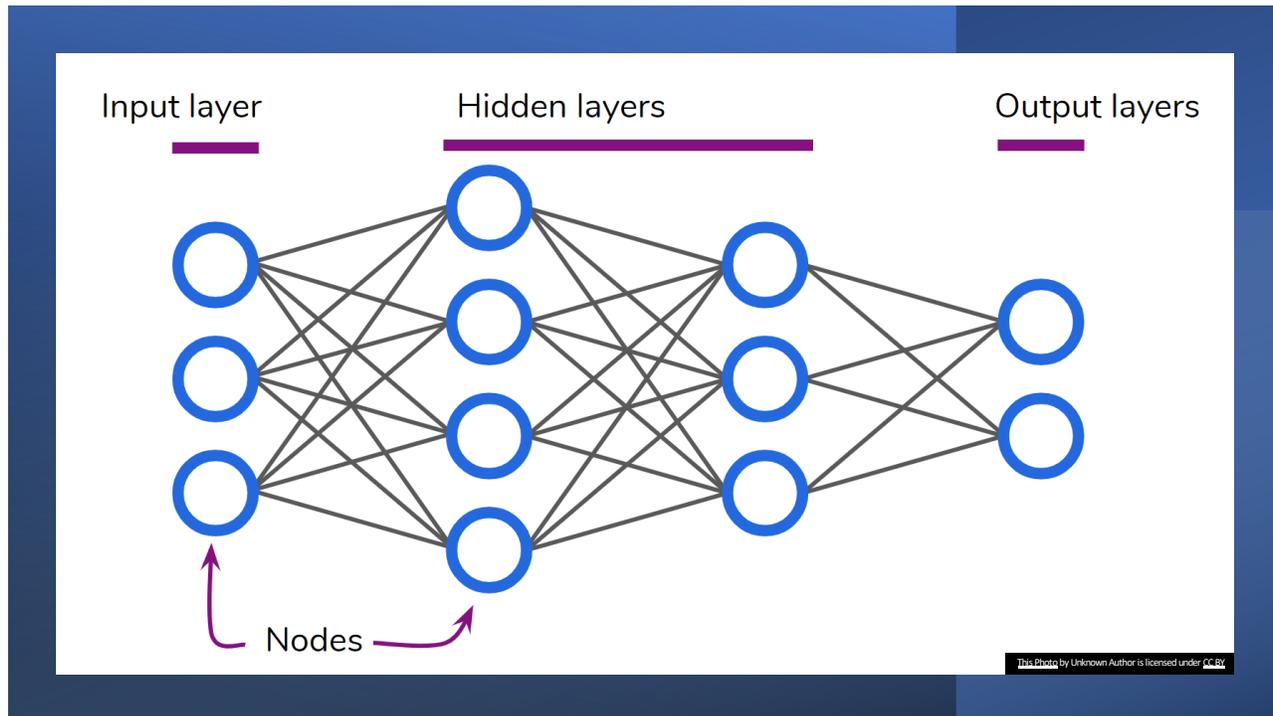If using k-folds, k-1 folds are used for training the base learners and the $k^{th}$ fold to train the meta learner

## Stacking with 2-fold dataset

Base Learner predictions are the new features; y remains the same

| $h_1$ | $h_2$ | ... | $h_B$ | $y$ |
|---|---|---|---|---|
| | | | | |
| | | | | |

Meta Learner

Multi-level Stacking

Artificial Neural Network

Multi-level Stacking

Input layer     Hidden layers     Output layers

Nodes

Extending Meta-learning to Deep Learning...

# Tell the language of the letter

 Malayalam

 Malayalam

 Hebrew

 Malayalam

 Gujarati

 Hebrew
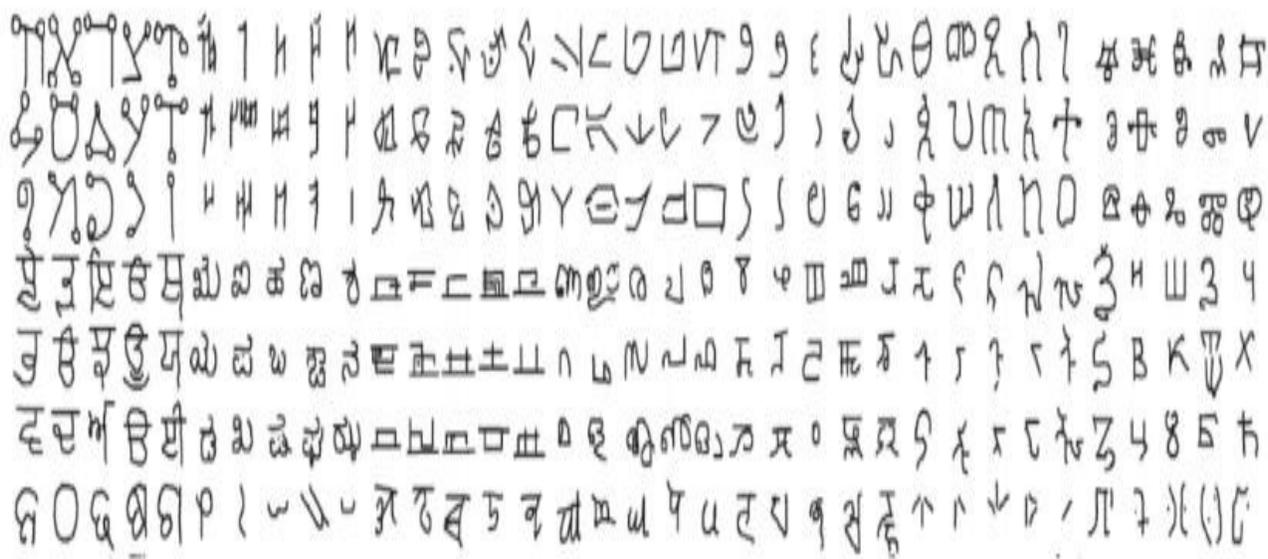
What's the ML algorithm that comes closest to doing this?

Can we make the machine learn in one-shot like you did?

# Suppose you have this omniglot labeled dataset…

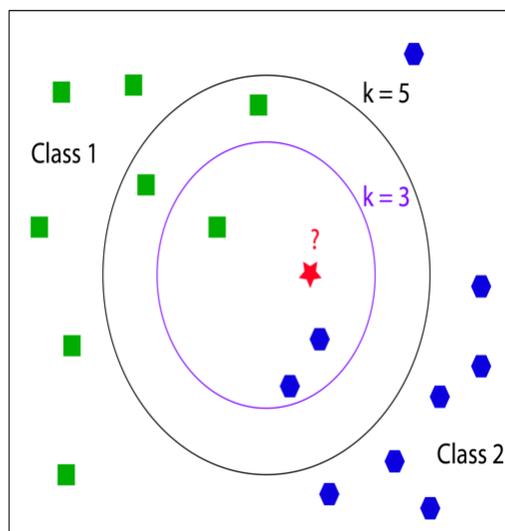How do you produce an ML model to recognize
a new letter not in the dataset?



One-shot Learning!

# K Nearest Neighbors: Classification

- In binary classification where y
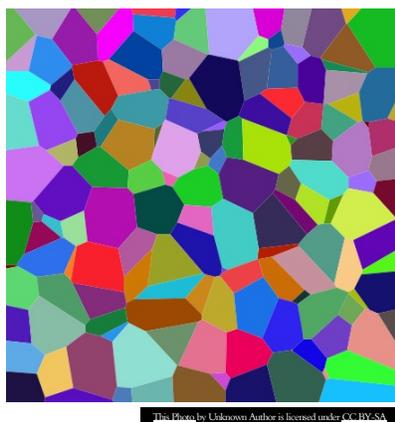  = +1 or –1, a test item is
  classified as

  $y_t = \text{sgn}(\sum_{i=1}^{k} y_i)$, where (xi, yi)
  are the k nearest neighbors of the
  test data item in the feature space

- The k nearest neighbors are
  determined based on a distance
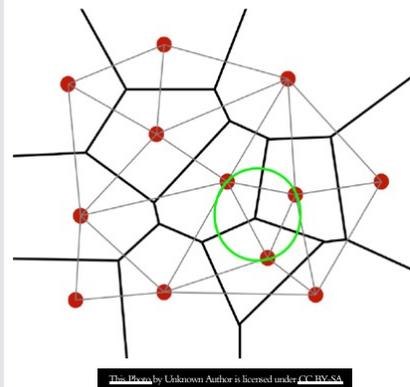  metric, **usually Euclidean**

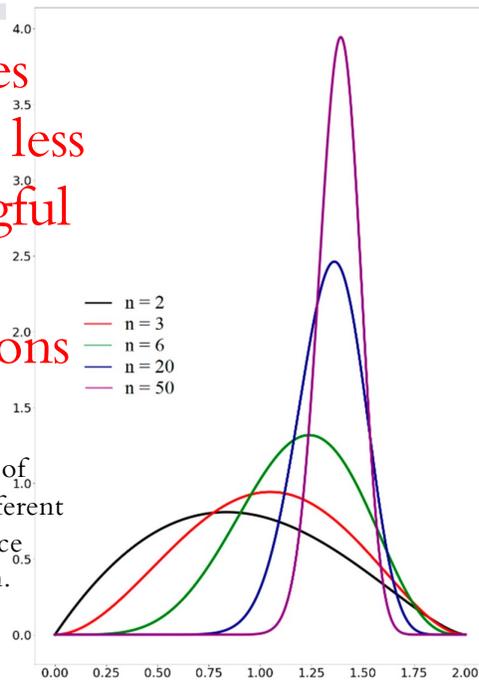# Visualization of the Induced Decision Boundary



1-NN:
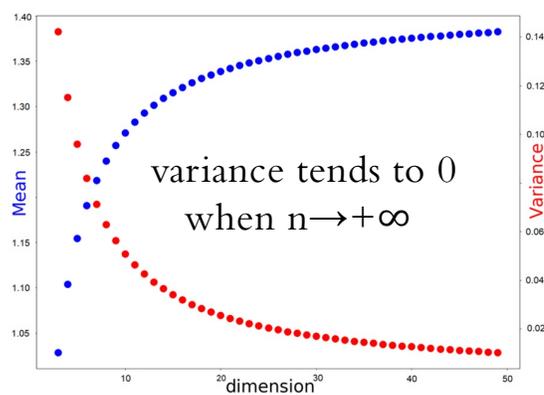VORONOI
DIAGRAM

◇

What's the problem with 1-NN?

## Distances become less meaningful in high dimensions

Distribution of distances for different values of space dimension n.



(a)

Lellouche, S., & Souris, M. (2019). Distribution of distances between elements in a compact set. *Stats*, *3*(1), 1–15.



variance tends to 0 when n→+∞

(b)

Can we build a neural network to learn to classify based on the distance in the feature space?

What if we compute the distances in the extracted feature space instead of in the original pixel space or learn a "metric space"?

*https://www.sjsu.edu/people/vishnu.pendyala/*
*@vishnupendyala*

Q&A