# Model Valuation and Performance Metrics

# Model Valuation and Performance Metrics

*K fold cross-validation*

**1.Data Preparation**: Splitting the dataset into training, validation, and test sets.

**2.Model Fitting**: Creating multiple models.

**3.Model Evaluation**:
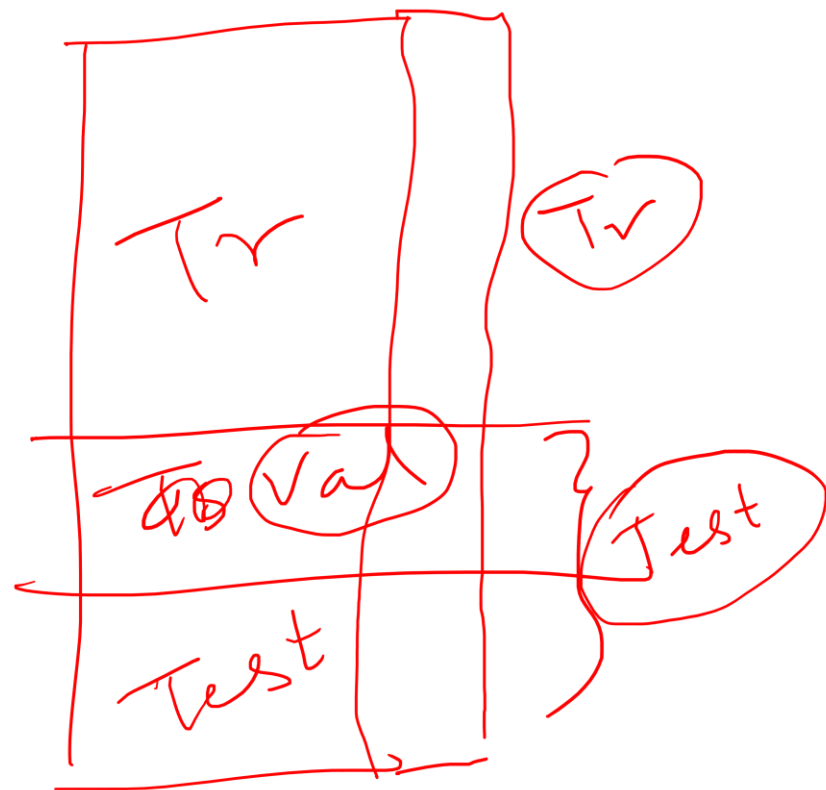
*LinReg / NonLin / KNN / SVR / DL / Reg tree*

- for Reg Problem: Mean Absolute error (MAE), Calculation of Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) on the validation set for both models.

*Log Reg / KNN / SVC / DT / DL*

- for Classification Problem: Using confusion matrix, accuracy, precision, recall (sensitivity), specificity, F1-score etc on the validation set.

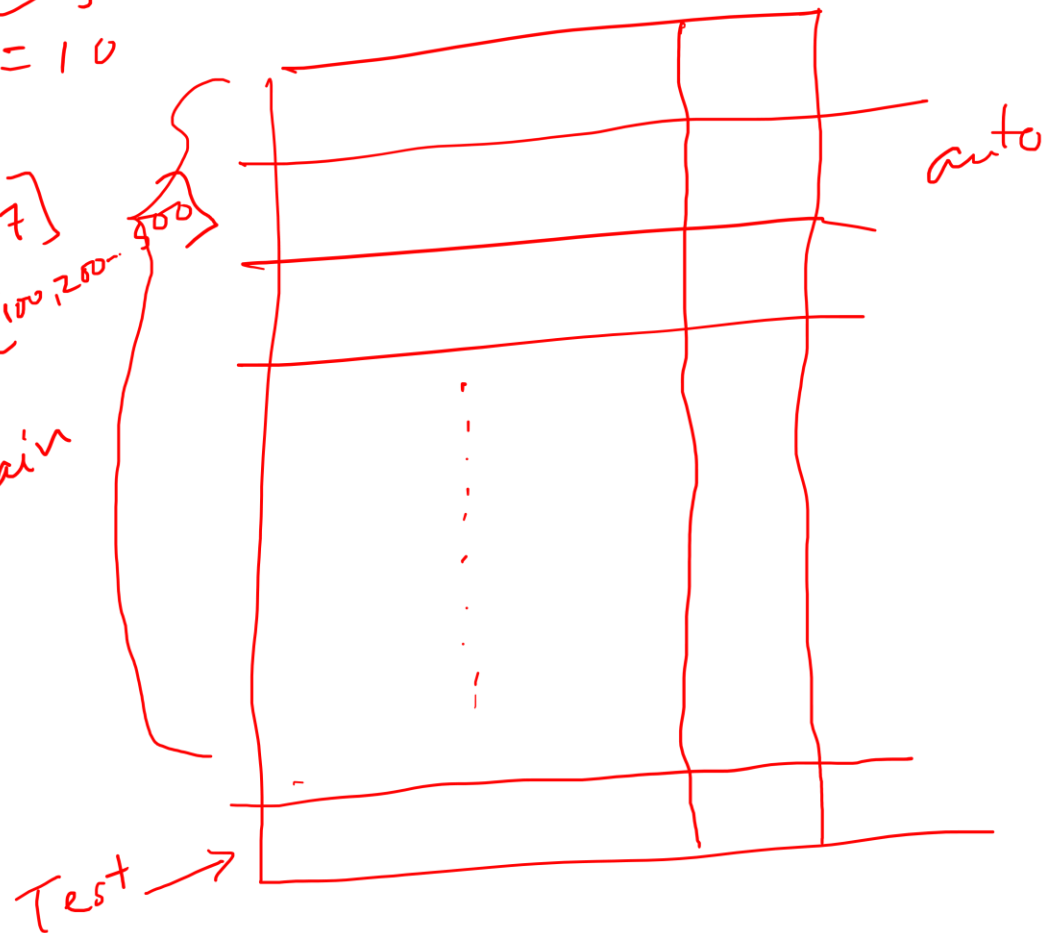**4.Cross-Validation**: Performing cross-validation on the training set.

**5.Final Model Selection and Evaluation**: Applying the best model on the test set.

Validation →



Tr

Tv

Val
db

Test

Test

K fold CV cross validation
= 10

DT
depth [3,7]
#tree[100,200—900]

Train

auto

Test →

KNN $\longrightarrow$ (K) $\longrightarrow$ how many neighbours

K Means $\longrightarrow$ (K) $\longrightarrow$ how many clusters

K fold $\longrightarrow$ (K) $\longrightarrow$ how many folds

# Model Valuation and Performance Metrics

- **Mean Squared Error (MSE)**

- MSE is a measure of the average of the squares of the errors—that is, the average squared difference between the estimated values and the actual value. It's a common measure of the estimation accuracy of a predictive model in regression tasks.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$

- **Root Mean Squared Error (RMSE)**

- RMSE is the square root of the MSE. It's a widely used measure of the differences between values predicted by a model or an estimator and the values observed. The RMSE represents the sample standard deviation of the differences between predicted and observed values.

$$RMSE = \sqrt{MSE}$$

# Model Valuation and Performance Metrics

- **Accuracy**

- Most commonly used metrics for evaluating classification models. It measures the proportion of total correct predictions (both true positives and true negatives) out of all predictions made.

- Accuracy=Number of Correct Predictions / Total Number of Predictions

Or, using the terms of the confusion matrix:

- Accuracy= (TP + TN) / (TP + FP + FN + TN)

- **Specificity**

- Specificity measures the proportion of actual negatives that are correctly identified as such (e.g., the percentage of healthy people who are correctly identified as not having the condition, in the medical context). It's a key metric when the cost of false positives is high.    Specificity=True Negatives (TN)/ (True Negatives (TN) + False Positives (FP))

- **Recall (Sensitivity)**

- Recall, also known as sensitivity, is the ratio of true positive predictions to the total actual positives. It answers the question: "Of all the actual positive instances, how many did we correctly classify as positive?"

      Recall=True Positives (TP) / (True Positives (TP) + False Negatives (FN))

# Model Valuation and Performance Metrics

- **Precision**

- Precision is the ratio of true positive predictions to the total positive predictions (including both true positives and false positives). It answers the question: "Of all instances classified as positive, how many are actually positive?"

    Precision=True Positives (TP)/(True Positives (TP) + False Positives (FP))

- **F1-Score**

- The F1-score is the harmonic mean of precision and recall. It provides a single score that balances both the precision and recall. It's particularly useful when you need to balance both precision and recall, such as in imbalanced datasets.

    F1-score=2×(Precision×Recall) / (Precision+Recall)

# Model Valuation and Performance Metrics



|  |  | True class | | Measures |
|---|---|---|---|---|
|  |  | Positive | Negative |  |
| Predicted class | Positive | True positive $TP$ | False positive $FP$ | Positive predictive value (PPV) $\dfrac{TP}{TP+FP}$ |
|  | Negative | False negative $FN$ | True negative $TN$ | Negative predictive value (NPV) $\dfrac{TN}{FN+TN}$ |
| Measures |  | Sensitivity $\dfrac{TP}{TP+FN}$ | Specificity $\dfrac{TN}{FP+TN}$ | Accuracy $\dfrac{TP+TN}{TP+FP+FN+TN}$ |

|  |  | Actual Value | |
|---|---|---|---|
|  |  | Positive | Negative |
| Result Obtained | Positive | True Positive (1- β) | False Positive Type-I Error (α) |
|  | Negative | False Negative Type-II Error (β) | True Negative |

Predicted class

|  | positive | negative |
|---|---|---|
| **Positive (P)** | **True Positive** (TP) | **False Negative** (FN) |
| **Negative (N)** | **False Positive** (FP) | **True Negative** (TN) |

True class

Row summary

| | |
|---|---|
| $TPR = \dfrac{TP}{P}$ | $FNR = \dfrac{FN}{P}$ |
| $TNR = \dfrac{TN}{N}$ | $FPR = \dfrac{FP}{N}$ |

Column summary

| | |
|---|---|
| $PPV = \dfrac{TP}{TP + FP}$ | $NPR = \dfrac{TN}{TN + FN}$ |
| $FDR = \dfrac{FP}{TP + FP}$ | $FOR = \dfrac{FN}{TN + FN}$ |

# Log Reg training and Thresholds



treshold prob = 0.5

Log Reg Curve

$$y = \frac{e^{\beta_0 + \beta_1 x_1}}{1 + e^{\beta_0 + \beta_1 x_1}}$$

Sigmoid $f^n$

Predicted as obese

The probability of obesity

# Perf of Log Reg on test data with diff Thresholds

Predictions on test data with tp = 0.5

$x \mid y$



Weight

|  |  | Actual | |
|---|---|---|---|
|  |  | Is Obese | Is Not Obese |
| Predicted | Is Obese | 3 | 1 |
| | Is Not Obese | 1 | 3 |

Sensitivity = 75%

Specificity = 75%

Acc → 75%

https://www.youtube.com/watch?v=4jRBRDbJemM&t=936s

# Perf of Log Reg on test data with diff Thresholds

Predictions on test data with tp = eg 0.1



|  |  | Actual | |
|---|---|---|---|
|  |  | Is Obese | Is Not Obese |
| Predicted | Is Obese | 4 | 2 |
|  | Is Not Obese | 0 | 2 |

$\dfrac{4}{4}$   $\dfrac{2}{4}$

Sensitivity =     100%     Acc

Specificity =     50%      75%

Think about an infectious disease. This is very important to correctly predict all the "yes" infected cases

# Perf of Log Reg on test data with diff Thresholds

Covid: Sensility

## Predictions on test data with tp = eg 0.9



This is better than 0.5 for sure

|  |  | Actual | |
|---|---|---|---|
|  |  | Is Obese | Is Not Obese |
| Predicted | Is Obese | 3 | 0 |
| | Is Not Obese | 1 | 4 |

$\frac{3}{4}$   $\frac{4}{4}$

Acc

Sensitivity =        75%

Specificity =        100%     Acc 75%.

But which threshold is the best?

https://www.youtube.com/watch?v=4jRBRDbJemM&t=936s

# Perf of Log Reg on test data with diff Thresholds

Predictions on test data with tp = eg 0.9



This is better than 0.5 for sure

| | | Actual | |
|---|---|---|---|
| | | Is Obese | Is Not Obese |
| Predicted | Is Obese | 3 | 0 |
| | Is Not Obese | 1 | 4 |

Sensitivity =      75%

Specificity =      100%

But which threshold is the best?

# ROC (Receiver Operator Curve) Curve and AUC (Area Under Curve)

# ROC (Receiver Operator Curve) Curve and AUC (Area Under Curve)

# KNN (K Nearest Neighbour)

# KNN (K Nearest Neighbour)



Cartesian/Manhattan/Cosine etc
2D or 3D or ND
Mean or voting
K = ? (HP tuning)
Which/How many IVs should we consider?
Any feature engineering? (stan/norm/unit transform etc)

$\cos \theta$

height

Reg Prob

$x_2$

$K = 5$

## 3-NN for regression

No. of Bedrooms

55
50
47
42
35
30

50+55+51 / 3 = 52
51
53
72

Class Prob

$(x_2, y_2)$

$(x_1, y_1)$

$\sqrt{(x_2-x_1)^2 + (y_2-y_1)^2}$

## 3-NN for classification

$x_2$

o: 3 x: 0 → predict o

$x_1$

Y

total sq. ft

House price prediction

$(402, 3)$
$478, 3$
$(602, 4)$
$(\quad , \quad)$

$x_1$

# SVM (Support Vector Machine)



$x_2$

Class 1
Class 2

1  2  3

margin

$x_1$

Lin

Vector

$x_1 \ldots x_n$

R → CC

Prompt

Vac

embedding

RAG

HR

GPT4
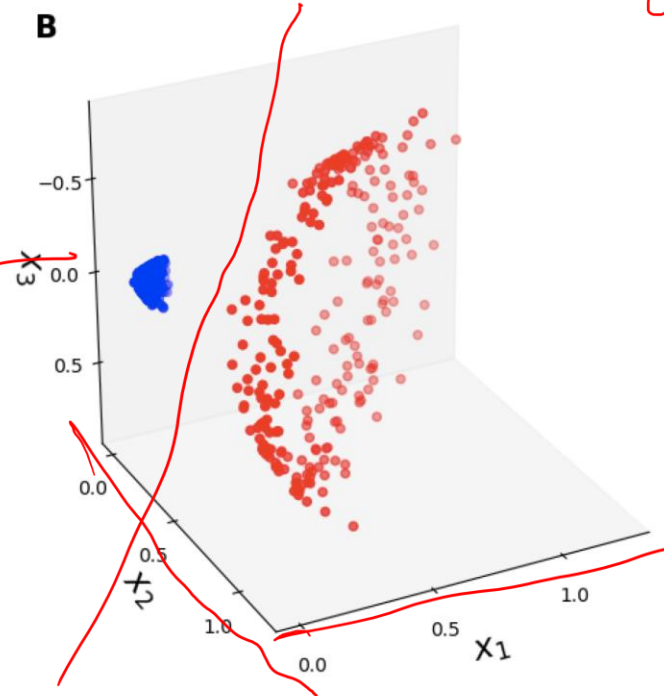
gemini

TS

Bard

→ Bert

FM

AI Chatbot

# SVM (Support Vector Machine)

# SVM (Support Vector Machine) – Kernel trick



https://gregorygundersen.com/blog/2019/12/10/kernel-trick/

# SVM (Support Vector Machine) – Kernel trick

$$x_3 = \sqrt{x_1 * x_2}$$
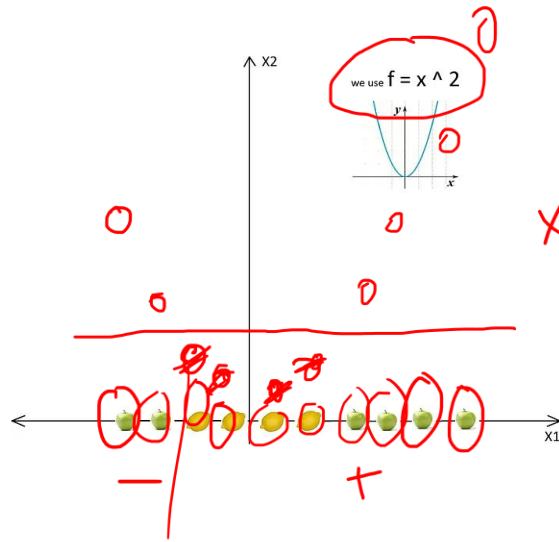
$$x^2 + y^2$$

$$x_3 = \sqrt{x_1^2 + x_2^2}$$
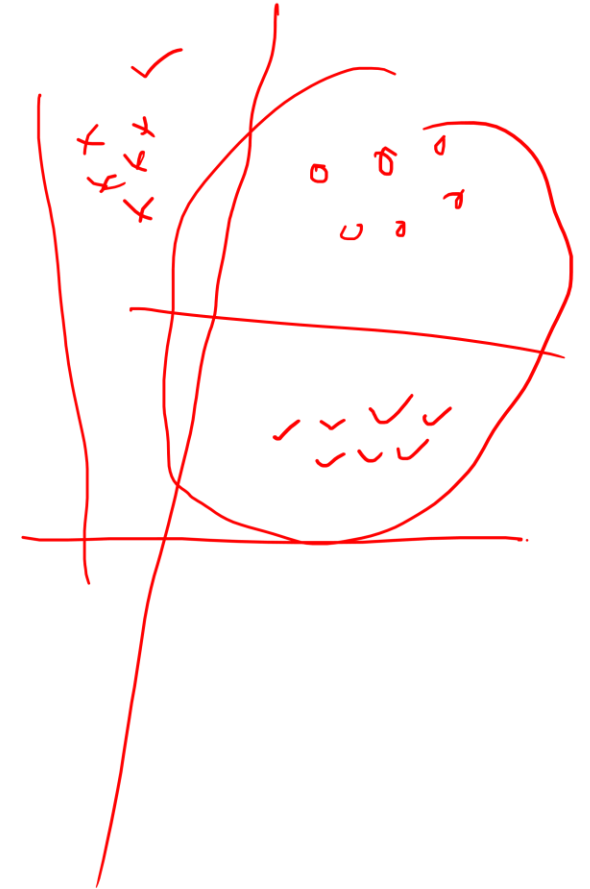
$$x_3 = 7$$

$$x_3 = x_1^2 + x_2^2$$

$$denm = 2^{2.25} + 2^{2.25}$$

$$= 4.5$$

Apple
$$x_3 = 4^2 + (-3)^2$$
$$= 16 + 9$$
$$= 25$$



https://towardsdatascience.com/svm-and-kernel-svm-fed02bef1200

# SVM (Support Vector Machine) – Kernel trick

lin / nonlinear / Radial

we use f = x ^ 2

$X_2 = X_1$

we use f = x ^ 2

https://towardsdatascience.com/svm-and-kernel-svm-fed02bef1200

# Random Forest



**X** d a t a s e t

N$_1$ feature     N$_2$ feature     N$_3$ feature     N$_4$ feature

Class N     Class O     Class M     Class N

MAJORITY VOTING
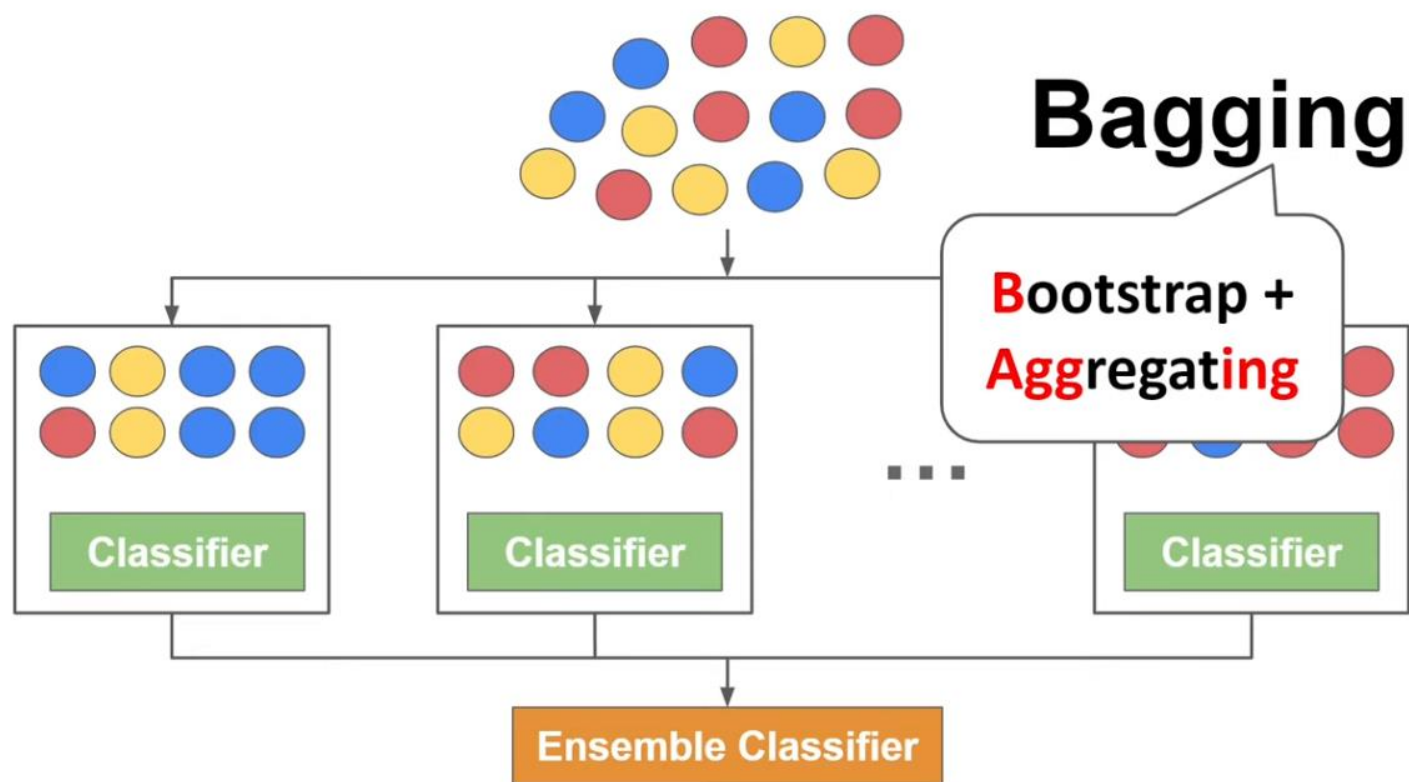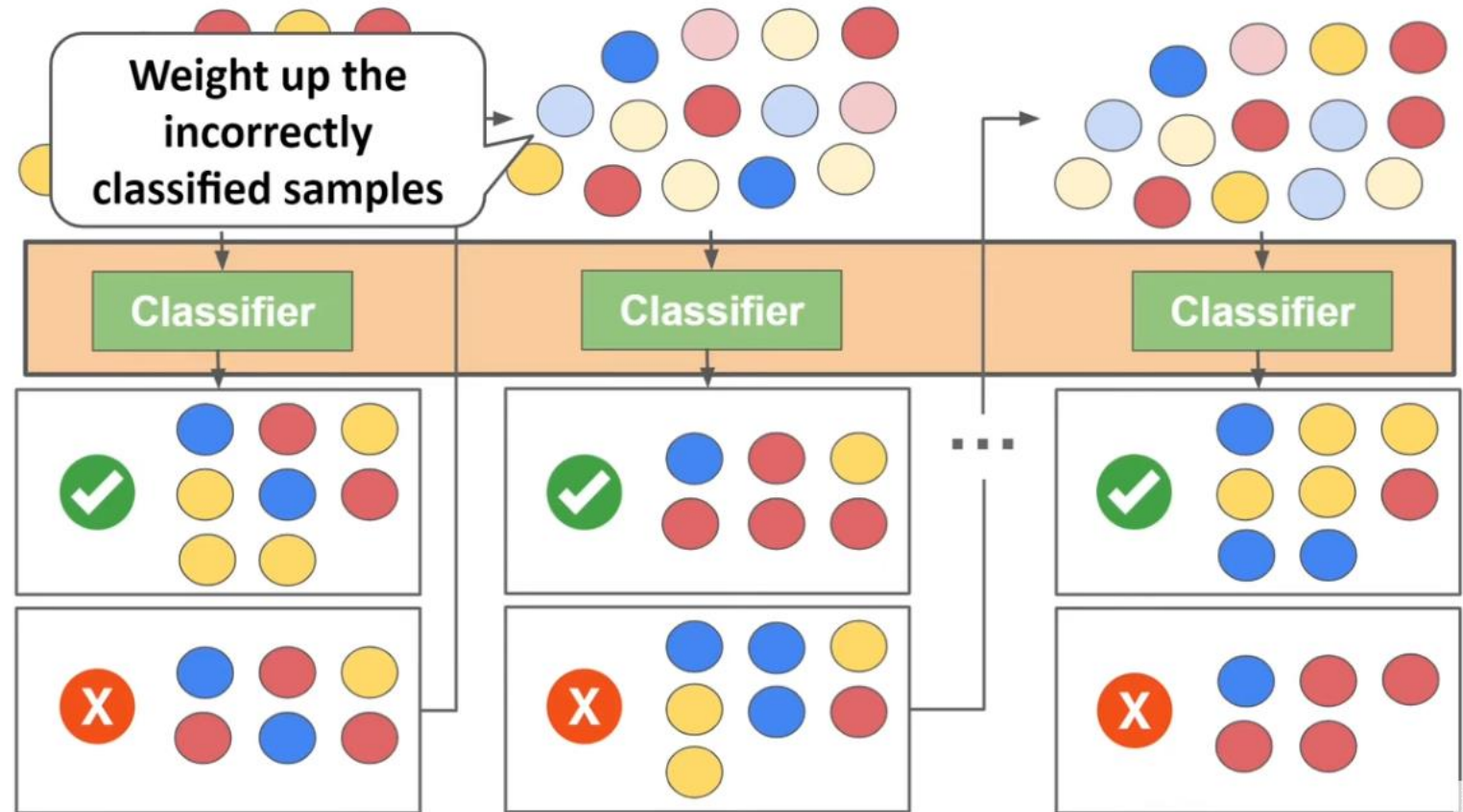
FINAL CLASS

# Bagging

# Boosting

# Regularization: Lasso vs Ridge vs Elastic

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 \right\} \quad \text{subject to} \quad \sum_{j=1}^{p} |\beta_j| \leq s \tag{6.8}$$

and

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 \right\} \quad \text{subject to} \quad \sum_{j=1}^{p} \beta_j^2 \leq s, \tag{6.9}$$