# BOOSTING

COMBINE WEAK LEARNERS THAT WERE TRAINED SEQUEN-
TIALLY AND PREDICTS BASED ON THE WEIGHTS OF EACH
MODEL (WEIGHT CALCULATED BASED ON MODEL PERFORMANCE).

## GRADIENT BOOSTING

$$dY(0) = Y - MEAN(Y)$$

FOR $k = 1 : (\#OF\ TREES / \#OF\ ITERATIONS)$:

$$LEARNER(k) = TRAIN\text{-}REGRESSOR(X, dY(k-1))$$
$$dY(k) = dY(k-1) - \alpha(k) * PREDICT(LEARNER(k), X)$$

"TRAIN OVER THE RESIDUAL"

## BAGGING

* GENERATE N DIFFERENT BOOTSTRAP TRAINING SAMPLE
  WITH REPLACEMENT
* TRAIN ALGORITHM ON EACH BOOTSTRAPPED SAMPLE
* COMBINE THEM ALL USING MAJORITY VOTE / MEAN

## RELATIVE VARIABLE IMPORTANCE

THE MEASURE IS MADE BASED ON #TIMES A VARIABLE
IS SELECTED FOR SPLITTING AND WEIGHTED BY THE IMPRO-
VEMENT TO THE MODEL AS A RESULT OF EACH SPLIT,
THEM AVERAGED OVER ALL TREES.

# CONFUSION MATRIX        POSITIVE CLASS = 0

|  |  | OBSERVED | |
|---|---|---|---|
|  |  | P | N |
| PREDICT | P | TP | FP |
|  | N | FN | TN |

$$\text{SENSITIVITY} = \frac{TP}{TP+FN} = \text{RECALL}$$
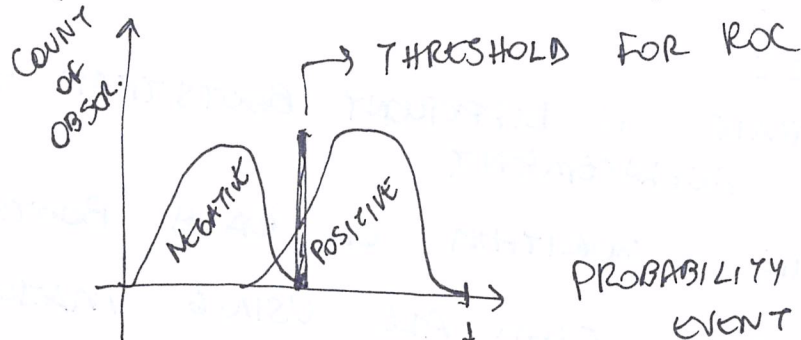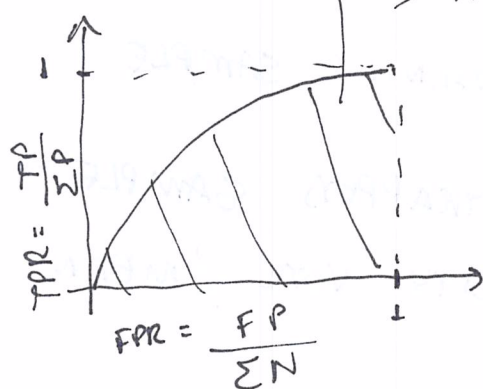
$$\text{SPECIFICITY} = \frac{TN}{TN+FP}$$

$$\text{PRECISION} = \frac{TP}{TP+FP}$$

$$\text{F1 SCORE} = 2 \cdot \frac{\text{PRECISION} * \text{RECALL}}{\text{PRECISION} + \text{RECALL}} = \frac{2}{\frac{1}{\text{PRECISION}} + \frac{1}{\text{RECALL}}}$$
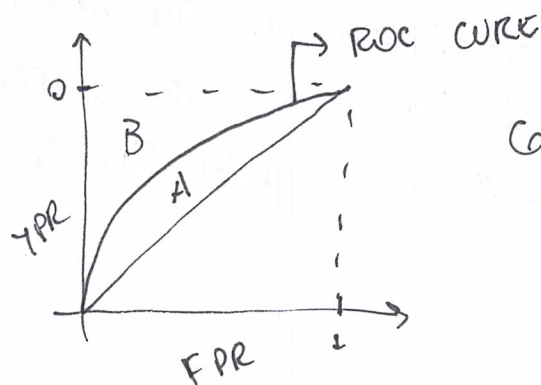
"HARMONIC MEAN"

$$\text{ACCURACY} = \frac{TP + TN}{P + N}$$
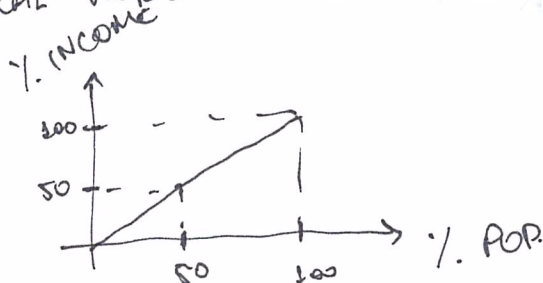
# ROC CURVE



→ AUC (AREA UNDER THE CURVE)

$$TPR = \frac{TP}{\Sigma P}$$

$$FPR = \frac{FP}{\Sigma N}$$

→ THRESHOLD FOR ROC

COUNT OF OBSR.

NEGATIVE   POSITIVE

PROBABILITY OF EVENT

# GINI COEFFICIENT   (INEQUALITY COEFFICIENT)



→ ROC CURVE

$$\text{GINI} = \frac{A}{A+B}$$

HIGHER A ⟹ BETTER MODEL

HIGHER INEQUALITY

REAL-WORLD EXAMPLE:

% INCOME

100

50

50   100

% POP.

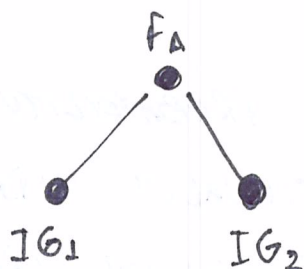"PERFECT DISTRIBUTION OF MONEY OVER POPULATION"

# CLASSIFICATION TREES

FOR ALL FEATURES TESTS INFORMATION GAIN USING ENTROPY:

$$IG = -\sum_{i=1}^{J} P_i \cdot \log_2 P_i$$

$P_i$ = PROPORTION OF CLASS $i$ ON NODE



$$IG_T = IG_1 + IG_2$$

(FOR FEATURE A THE INFORMATION GAIN WILL BE THE SUM OF $IG_1 + IG_2$)

# REGRESSION TREES

AT EACH ITERATION FOR EACH FEATURE $x_k$ FIND OPTIMAL $S$:

$$\min_{S} \left[ MSE(y | x_k < s) + MSE(y | x_k \geq s) \right]$$

(*$s$ IS THE ~~CLASSIFICATION~~ CUT OFF)

* FOR BOTH METHODS: VARIABLE IMPORTANCE GENERALLY BE COMPUTED BASED ON CORRESPONDING REDUCTION OF PREDICTIVE ACCURACY WHEN THE PREDICTOR OF INTEREST IS REMOVED OR SOME MEASURE OF DECREASE OF NODE IMPURITY.

# LINEAR REGRESSION

LINEAR APPROACH TO MODELLING RELATIONSHIP BETWEEN SCALAR RESPONSE TO EXPLANATORY VARIABLES "ORDINARY LEAST SQUARES"

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_n \cdot X_n$$

# LOGISTIC REGRESSION

LOG-ODDS OF PROBABILITY OF AN EVENT IS A LINEAR COMBINATION OF EXPLANATORY VARIABLES "MAXIMUM LIKELIHOOD ESTIMATION"
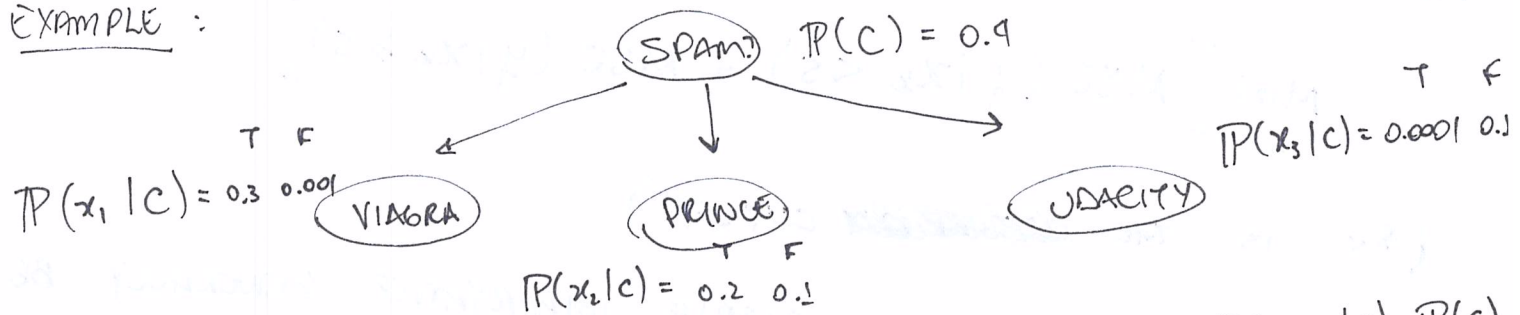
$$P(Y=1) = \frac{1}{1 + e^{-x}}$$

\* <u>VARIABLE</u> IMPORTANCE FOR REGRESSION CAN BE SET BASED ON THE COEFFICIENTS ONLY. IF THE FEATURES ARE NORMALIZED (CONTINUOUS) OR THEY ARE DISCRETE.

## NAIVE BAYES CLASSIFICATION

BASED ON BAYES THEOREM WITH CONDITIONAL PROBABILITY OF EVENT TO PREDICT. ASSUMES $x_i$ CONDITIONALLY INDEPENDENT OF EVERY OTHER FEATURE $x_j$ $(i \neq j)$ GIVEN CATEGORY $C_k$ =>

$$\boxed{P(C_k \mid x_1, ..., x_n) \propto P(C_k) \cdot \prod_{i=1}^{n} P(x_i \mid C_k)}$$

<u>EXAMPLE</u> :

SPAM?  $P(C) = 0.9$

$$\begin{array}{cc} & T \quad F \\ P(x_3 \mid c) = & 0.0001 \quad 0.1 \end{array}$$

$$\begin{array}{cc} & T \quad F \\ P(x_1 \mid c) = & 0.3 \quad 0.001 \end{array}$$  VIAGRA

PRINCE

UDACITY

$$\begin{array}{cc} & T \quad F \\ P(x_2 \mid c) = & 0.2 \quad 0.1 \end{array}$$

$$P(c \mid x_1 = T, x_2 = F, x_3 = F) \propto P(x_1 = T \mid c) \cdot P(x_2 = F \mid c) \cdot P(x_3 = F \mid c) \cdot P(c)$$
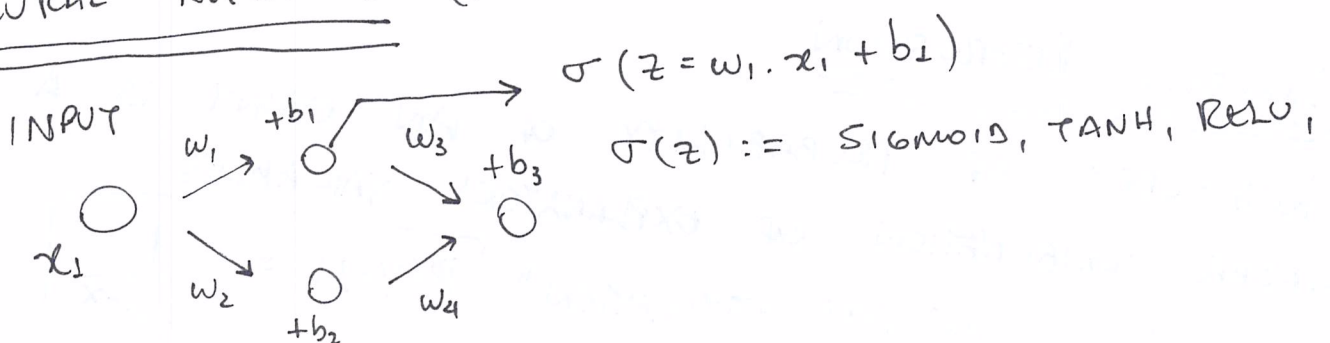
## RANDOM FOREST

SELECTING BAGGING SAMPLES FOR EACH TREE CHOOSE RANDOM FEATURES ($\sqrt{D}$).
   \* REDUCES VARIANCE BUT RANGE (AS ALL TREE MODELS) IS LIMITED.

## NEURAL NETWORKS (ARTIFICIAL NEURAL NETWORKS)

$$\sigma(z = w_1 \cdot x_1 + b_1)$$

INPUT

$w_1$  $+b_1$  $w_3$  $+b_3$

$\sigma(z) :=$ SIGMOID, TANH, RELU,

$x_1$

$w_2$  $+b_2$  $w_4$

$$\text{SIGMOID} := \sigma(z) = \frac{1}{1+e^{-z}} \quad ; \quad \sigma(z) \in (0,1)$$

$$\text{TANH} := \sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \text{TANH}(z) \quad ; \quad \sigma(z) \in (-1, 1)$$

$$\text{RELU} := \sigma(z) = \text{MAX}(0, z) \quad ; \quad \sigma(z) \in [0, \infty) .$$

$$\text{SOFT MAX} := \sigma_i(\vec{z}) = \frac{e^{z_i}}{\sum_{j=1}^{J} e^{z_j}} \quad , \quad i = 1, \dots, J \quad ; \quad \sigma(\vec{z}) \in (0,1)$$

## K-NN (K - NEAREST NEIGHBORS)

REGRESSION OR CLASSIFICATION OF K NEAREST NEIGHBORS
BASED ON DISTANCE FUNCTIONS:

EUCLIDEAN : $\sqrt{\sum_{i=1}^{K} (x_i - y_i)^2}$

MANHATTAN : $\sum_{i=1}^{K} |x_i - y_i|$

MINKOWSKI : $\left( \sum_{i=1}^{K} (|x_i - y_i|)^q \right)^{1/q}$

## K-MEANS

0. PLACE CONTROIDS AT RANDOM LOCATIONS (K CENTROIDS)
1. FIND NEAREST CONTROID TO EACH OBSERVATION
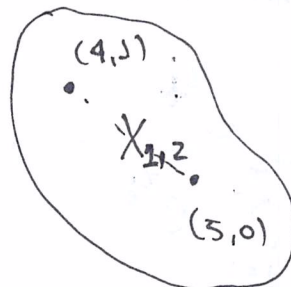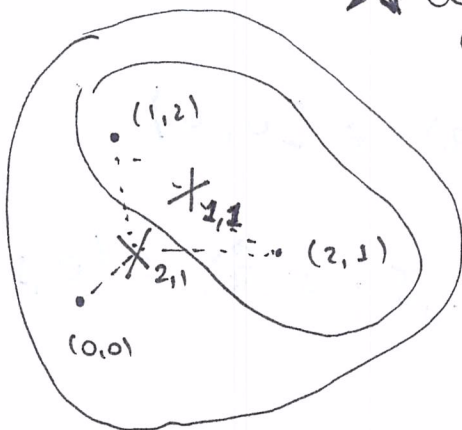2. ASSIGN OBSERVATION TO CLOSER CLUSTER
3. CALCULATE NEW CONTROID

\* REPEAT UNTILL DIFF. FROM PREVIOUS CONTROID IS
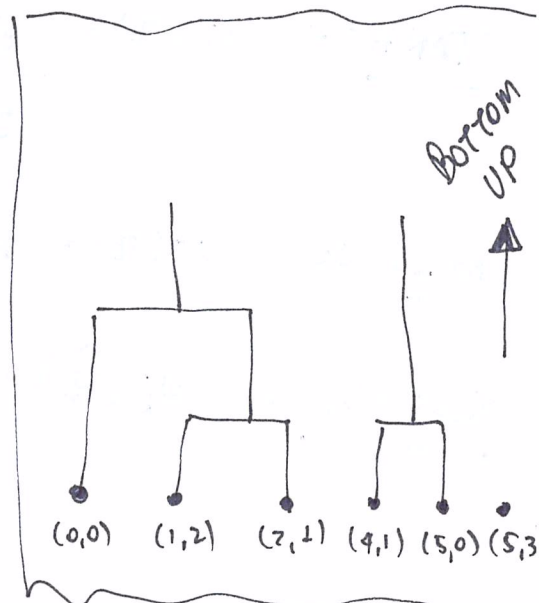MINIMUM (DIFFERENCE FROM DISTANCES)

→ K-MEANS ONLY USES EUCLIDEAN DISTANCE

R3

# HIERARCHICAL CLUSTERING



☆ CENTROIDS
CLUSTER
DISTANCE
MEASURES

(5,3)

(1,2)

$X_{1,1}$

$X_{2,1}$     (2,1)

(0,0)

(4,1)

$X_{1,2}$

(5,0)

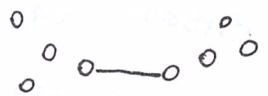BOTTOM UP

(0,0)  (1,2)  (2,1)  (4,1)  (5,0)  (5,3)

~~DISTANCES~~ DISTANCES TO BE USED:

EUCLIDEAN, MANHATTAN, MINKOWSKI, MAHALANOBIS := $\sqrt{(a-b)^T S^{-1} (a-b)}$

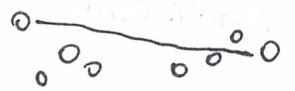\* $S$ IS THE COVARIANCE MATRIX

→ SINGLE LINKS: $D(c_1, c_2) = \min D(x_1, x_2)$

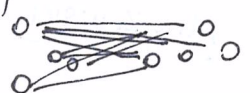DISTANCE BETWEEN CLOSEST ELEMENTS IN CLUSTERS

→ COMPLETE LINKS: $D(c_1, c_2) = \max D(x_1, x_2)$

DISTANCE BETWEEN FARTHEST ELEMENTS IN CLUSTERS

→ AVERAGE LINKS: $D(c_1, c_2) = \dfrac{1}{|c_1|} \dfrac{1}{|c_2|} \sum_{x_1} \sum_{x_2} D(x_1, x_2)$
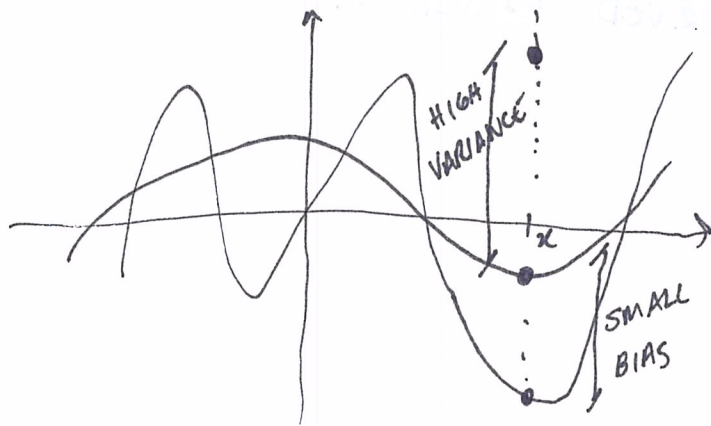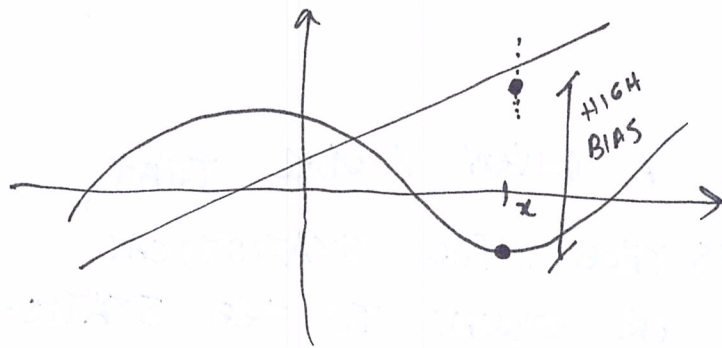
AVERAGE OF ALL PAIRWISE DISTANCES

→ CONTROIDS: $D(c_1, c_2) = D\left[ \left( \dfrac{1}{|c_1|} \sum_{x_1} \vec{x} \right), \left( \dfrac{1}{|c_2|} \sum_{x_2} \vec{x} \right) \right]$
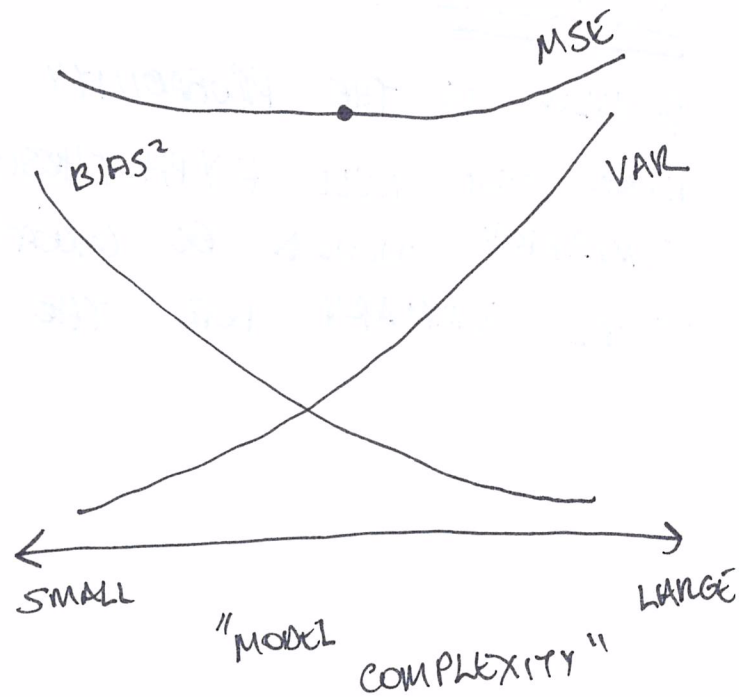
☆ DISTANCE BETWEEN CONTROIDS
(MEANS) OF TWO CLUSTERS

## BIAS - VARIANCE DECOMPOSITION

- BIAS IS THE ERROR FROM ERRONEOUS ASSUMPTIONS IN THE MODEL
- VARIANCE IS THE ERROR FROM SENSITIVITY TO SMALL FLUCTUATIONS IN THE TRAINING SET.

$$MSE(\hat{\theta}) = BIAS^2(\hat{\theta}) + VAR(\hat{\theta})$$

$$\begin{cases} BIAS\uparrow = UNDERFITTING \\ VAR\uparrow = OVERFITTING \end{cases}$$
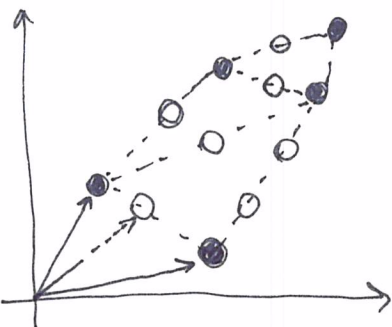
## CROSS- VALIDATION

PREVENTS OVERFITTING BECAUSE MODEL THAT FITS RANDOM NOISE ON TRAINING DATA WON'T PERFORM GOOD ON VALIDA-TION DATASET.

K-FOLD : SEPARATES DATA ON K FOLDS, TRAINING WILL BE K-1 FOLDS AND VALIDATION WILL BE 1

LOOCV : LEAVE-ONE-OUT CROSS-VALIDATION IS THE SAME AS K-FOLD BUT K = N.

## SMOTE (SYNTHETIC MINORITY OVER-SAMPLING TECHNIQUE)

→ CREATES NEW "SYNTHETIC" OBSERVATIONS
- IDENTIFY FEATURE VECTOR AND NEAREST NEIGH.
- TAKE THE DIFF. BETWEEN TWO
- MULTIPLY DIFF. BY RANDOM BETWEEN 0 AND 1
- IDENTIFY NEW POINT ON LINE SEGMENT BY ADDING RANDOM NUMBER TO FEATURE VECTOR

- REPEAT @



R4

# P- VALUE

P-VALUE IS THE PROBABILITY FOR A GIVEN MODEL THAT, WHEN THE NULL HYPHOTHESIS IS TRUE, THE STATISTICAL SUMMARY WOULD BE GREATER OR EQUAL TO THE STATISTICAL SUMMARY FOR THE OBSERVED RESULTS.

## P-VALUE

P-VALUE IS THE PROBABILITY FOR A GIVEN MODEL THAT, WHEN THE NULL HYPHOTESIS IS TRUE, THE STATISTICAL SUMMARY WOULD BE GREATER OR EQUAL TO THE STATISTICAL SUMMARY FOR THE OBSERVED RESULTS.

## REGULARIZATION

: REDUCE VARIANCE AT THE COST OF INTRODUCING SOME BIAS

$$\left( \begin{array}{l} \uparrow BIAS \Rightarrow UNDERFITTING \\ \uparrow VAR \Rightarrow OVERFITTING \end{array} \right)$$

LINEAR REGRESSION MODEL

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

ORDINARY LEAST SQUARES (OLS) $\Rightarrow$ ESTIMATE $\hat{\beta}$ SUCH A WAY THAT SUM OF SQUARES OF RESIDUALS IS AS SMALL AS POSSIBLE

$$MIN \, L_{OLS}(\hat{\beta}) = MIN \sum_{i=1}^{n} (y_i - x_i^T \hat{\beta})^2 = MIN \| y - X\hat{\beta} \|^2$$

IN ORDER TO OBTAIN $\hat{\beta}_{OLS} = (X^T X)^{-1}(X^T Y)$

- ## RIDGE REGRESSION (L2 PENALTY)

OLS LOSS FUNCTION IS AUGMENTED IN A WAY WE PENALIZE THE SIZE OF PARAMETER ESTIMATES :

$$L_{RIDGE}(\hat{\beta}) = \sum_{i=1}^{n} (y_i - x_i^T \hat{\beta})^2 + \lambda \sum_{j=1}^{m} \hat{\beta}_j^2 = \| y - X\hat{\beta} \|^2 + \lambda \| \hat{\beta} \|^2$$

  \* ALSO CALLED AS L2 PENALTY

- ## LASSO REGRESSION (L1 PENALTY)

SIMILAR TO RIDGE REGRESSION BUT LOSS FUNCTION IS:

$$L_{LASSO}(\hat{\beta}) = \sum_{i=1}^{n} (y_i - x_i^T \hat{\beta})^2 + \lambda \sum_{j=1}^{m} |\hat{\beta}_j|$$

   * ALSO CALLED L1 PENALTY

- ## ELASTIC NET

A COMBINATION OF BOTH RIDGE REGRESSION AND LASSO REGRESSION, LOSS FUNCTION IS:

$$L_{eNET}(\hat{\beta}) = \sum_{i=1}^{n} \frac{(y_i - x_i^T \hat{\beta})^2}{2n} + \lambda \left( \frac{1-\alpha}{2} \sum_{j=1}^{m} \hat{\beta}_j^2 + \alpha \sum_{j=1}^{m} |\hat{\beta}_j| \right)$$

WHERE $\alpha$ IS THE MIXING PARAMETER BETWEEN RIDGE $(\alpha=0)$ AND LASSO $(\alpha=1)$.

# GRID-SEARCH

GRID-SEARCH IS USED TO FIND THE OPTIMAL HYPERPARAMETERS OF A MODEL WHICH RESULTS IN THE MOST 'ACCURATE' PREDICTIONS.

   IT CAN BE CHOOSE THE FOLLOWING PARAMETERS:

- PENALTY : L1, L2, ELASTIC NET
- LEARNING RATE
- METRIC : ACCURACY, RECALL, PRECISION, F1-SCORE

   * DISCUSSIONS ABOUT RANDOM-SEARCH IS BETTER (AND FASTER) THAN GRID-SEARCH. SAME CONCEPT BUT INSTEAD OF SETTING SOME VALUES/INPUTS THEY ARE RANDOMLY CHOSEN.

# (PROBABILITY / ODDS / LOG ODDS)

PROBABILITY of $\boxed{80\%}$ OF RAIN TODAY

ODDS RATIO IS 80% OF CHANCE OF RAIN DIVIDED BY 20% OF CHANCE OF NOT RAIN $\Rightarrow 80\% / 20\% = \boxed{4}$

LOG ODDS IS THE LOGARITHM OF ODDS $\Rightarrow \boxed{LN(4)}$

* ODDS RATIO $= \dfrac{P(A)}{P(-A)} = \dfrac{P(A)}{1 - P(A)}$   $\therefore$ LOG ODDS $= LN\left(\dfrac{P}{1-P}\right)$

AND   LOG ODDS $= LN\left(\dfrac{P}{1-P}\right) = \beta_0 + \sum_{i=1}^{m} \beta_i x_i$