# BIAS-VARIANCE TRADEOFF (Supervised Learning)



Bias



Variance

✓ Error = Bias$^2$ + Variance + Irreducible Error

Increasing Bias will decrease variance
Increasing variance will decrease bias.

Bias : Assumptions about the form of target function

Variance : Prediction error if different training set is used

Parametric algorithms have high bias & low variance
E.g. Linear Regression, logistic regression, LDA

Non-parametric algorithms have low bias & high variance
E.g. Decision trees, KNN, SVM

Examples: If value of k is high in KNN → increased bias
∨ Depth of Decision tree inc → variance increased

Bias → Under-fitting on training data (under-catching signal)

Ways to address it:

1) Train longer
2) Use more complexity model
3) Add features
4) Decrease regularization

Variance → over-fitting in validation/dev set (over-catching signal ie noise too)

Ways to address it

1) Add more data
2) Decrease #features
3) Increase regularization
4) Use simpler model (less complex)
5) Early-stopping

memorizing training data and not well "generalizing on unseen/test data

Example: Linear Regression models may under-fit the training data but can generalize well (can perform well in validation or test set)

Polynomial Regression model may fit the training data well but can fail to generalize in valid/test data.

Some situations

Questions: What is bias?
What is variance?
What is bias-var tradeoff?
What is underfitting? Ways to address it?
What is overfitting? Ways to address it

# HANDLING IMBALANCED DATASETS

1) Extract test set before applying any techniques to handle imbalance (test set can be imbalanced as data coming for scoring may be imbalanced)

(Remember: It is not "necessary" to balance classes e.g. using tree based approaches)

2) Get the evaluation metrics right (if data is imbalanced)

- Not accuracy ✗
- Kappa : class accuracy normalized by imbalance of classes
- Precision ↑Recall (of both classes individually)
- AUCPR (if +ve class is more important)
- AUCROC (if both classes are important but can be optimistic for +ve minority class)

3) Fixing in Data
✓ (Making classes balanced)
- Oversampling of Minority Class
  - Random oversampling (duplicating)
  - SMOTE (synthetic minority oversampling technique)
  - ADASYN (may add noise)
- Undersampling of Majority Class
  - Random undersampling
  - Edited Nearest Neighbor (cleans, not undersample(s))
  - Tomek links (cleans, not undersample(s))

If data is fixed for balance → regular classification metrics

4) Letting ML Algo take care of class imbalance

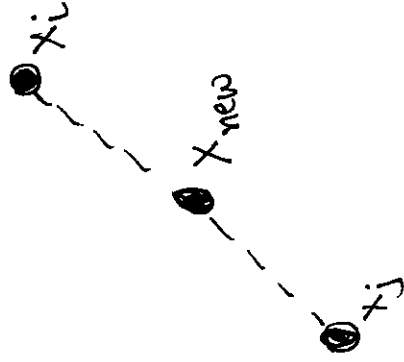One Class SVM (trains on only minority class may class can be used for testing/evaluation)

Cost Sensitive Algos: (penalizes minority
— Custom Obj/Loss fn (example class mistakes heavily)
    class dependent cost)
— Specify class wt. (Class-dependent cost)
    classifying diff wt to samples from diff. classes)

Tree-based Algos: Tree based Algos
generally fare well on
imbalanced datasets as
yike tree addresses
both classes (better than
parametorized algos)

# OVERSAMPLING TECHNIQUES

**1) SMOTE:** (Synthetic Minority Over-Sampling Technique)

→ For a given obs $x_i$, a new (synthetic) observ. is generated by interpolating between one of $k$ nearest neighbors of $x_i$,
e.g. $x_j$



$$x_{new} = x_i + \lambda (x_j - x_i)$$
where $\lambda = [0,1]$

→ To select $x_i$, three options:

[random] 1) regular : randomly select $x_i$

[KNN] 2) borderline : Separate all $x_i$ into 3 classes using $k$ nearest neighbor

    a) noise : all nearest-neighbors are of diff. class than $x_i$

    b) in danger: atleast half of nearest neighbors are of same class as $x_i$

    c) safe : all nearest neighbors are of same class as $x_i$

[SVM] 3) SVM: Uses SVM to identify samples close to decision boundary and selects $x_i$ from these points

→ **Disadv:** Can generate noisy samples by interpolating new pts b/w marginal outliers and inliers.

**Soln:** After oversampling, we under-sampling techniques such as Edited Nearest Neighbor or Tomek links to clean up

---

2) **ADASYN:** (Adaptive Synthetic Sampling)

→ For a given $x_i$, the number of new samples generated $\alpha$ number of nearby samples which do not belong to the same class as $x_i$.

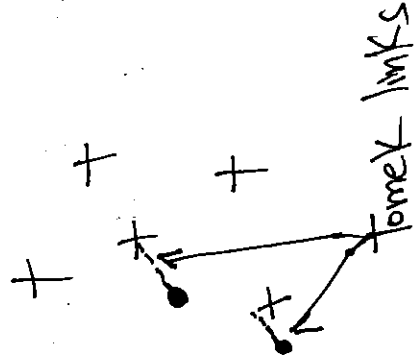$$\Rightarrow \text{\#new samples} \, \alpha \, 3$$

→ **Disadv:** Focuses on generating new samples which are outliers

**soln:** After oversampling we undersampling tech such as Edited Nearest Neighbor or Tomek links to clean up

UNDERSAMPLING TECHNIQUES

1) TOMEK LINKS:

- A tomek link exists b/w 2 observations if they are of different classes and nearest neighbors of each other



tomek links

- For undersampling we remove any obs. from majority class for which tomek links are identified

- It doesn't create balance between classes, it simply "cleans" the dataset by removing noisy obs. of majority class, which may result in easier classification problem

## 2) EDITED NEAREST NEIGHBORS:

— For the majority class samples, $k$-nearest neighbors
  ✓ are computed and if the samples do not agree
    ("enough") to their $k$-nearest neighbors, they are
    removed

— This simply "cleans" the noisy samples from
  ✓ majority class and does not necessarily "balance"
    the classes

Question: ways to handle imbalanced dataset?

# SPARSE DATA / FEATURES

## SPARSE vs. MISSING:

✓ — Sparse: many of the values are zero and you "know" they are zero

✓ — Missing: Value is "unknown" → may be zero or non-zero

## How TO HANDLE SPARSE DATA:

— Features with high sparsity % → remove

— Encode sparse data with some const value (can be zero) so that the ML algo will pick signals (corresponding i.e.m. labels) to this const value, if any

✓ — Convert to lower dimensions e.g. PCA ↓dense

✓ — Embeddings

✓ — Use Algos that can handle sparse data e.g. xgboost, FM/DeepFM, wide/Deep Model

✓ — Discretization (for continuous)

— Feature-hashing (for categorical)

— Feature-crossing (for categorical)

# ENCODING NUMERICAL FEATURES

- Normalization
- Standardization
- Bucketing / Discretization

# ENCODING CATEGORICAL FEATURES

- One-hot    (if #categories is less)
- Embeddings
- Feature-hashing .    (multiple categories hashed to a single bucket)

Question: Ways to encode?
1. Sparse data
2. Categorical data
3. Numerical data

## MISSING DATA

- Some Algorithms handle missing data e.g. tree based methods (robust to missing data) (Decision tree)
  - ✓.

- Decide b/w Row-wise Imputation or Column wise.

### Column wise
1) Below threshold % data are missing → Discard
2) Mean, median, mode imputation
   - Nearest neighbor fit (resource intensive)
   - Impute
3) Predictive models - similar to predicting data
   - constant for time-series data
4) Use a global constant
5) Add a column to signify missing values
6) Impute using adjacent values (time-series data)
   - In presence of collinearity, correlated features that have high % of missing values can be discarded (entirely)
   - Discard
   - below threshold %

- Imputf. Impute based on other predictors

- Practically, impute by any method then use cross-validation to check if evaluation metric is not giving desired results.

- It is important to note that as soon as we impute
  - ✓ → introducing bias

Question: How to handle missing data?

# THEORY BEHIND MISSING VALUES

1) Missing Completely at Random (MCAR): The propensity of a data point to be missing is completely random i.e. it doesn't depend on any values in the dataset, missing or observed.
   → 'Little' test to find out MCAR

2) Missing At Random (MAR): The propensity for a data point to be missing is not related to missing data per se, but to other observed data in the dataset.
   E.g. in a dataset for a survey of wt → women are less likely to fill the values of wt → MAR
   (depends on other field 'sex')

   → Create a dummy var. and plot other predictors against this dummy var → if relationship exists = MAR

3) Missing Not At Random (MNAR): The propensity of a data point to be missing depends on the value of depressed individuals → bez of depression level value missing data itself in that column
   E.g. in a survey of depression level value containing depression levels → people who have high depression levels are more likely i.e. to leave it blank to detect it → No way to detect it

Backdecials → Regression Mean

| | Mean | Regression | EM (like MLE) | MI (Multiple Imputation — Most probable → R package |
|---|---|---|---|---|
| MCAR | ✓ | ✓ | ✓ | ✓ → cannot detect |
| MAR | ✗ | ✓ | ✓ | ✓ |
| MNAR | | | | ✓ → cannot detect |

GENERAL STRATEGY: 0: Study relationship of missing values to other predictors & causally of predictors. If model do not have predictors.
1: Impute OR really high or low value to the missing entry
2. Create a new column, put 1 if the original column had /so/144 no value else put 0

GENERATIVE MODEL : → mean imputation is seen
(or) median if some outlier in available data)

Reason: Generative models estimate parameters from the distribution of data, imputing mean does not change the existing dist. of data

DISCRIMINATORY MODEL : → 1) Impute with a global constant value (high or low)

2) Create a new derived column put 1 if original column is empty

Reason: Adding a discriminatory feature so that the model will learn the discriminatory decision boundary

# OUTLIERS

When not to worry about Outliers?
- Tree based Model
- Non-parametric test
  ✓

When to worry about Outliers?
- Especially regression
  ✓

# DETECTING OUTLIERS:
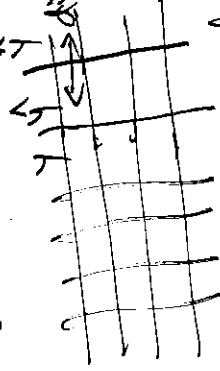
1) Univariate approach:
   (column wise)
   ✓

$$IQR = Q_3 - Q_1$$
1.5 times below $Q_1$
1.5 times IQR above $Q_3$

2) Multivariate: Compute the distance of observations from euclidean origin & plot
   (Row - wise)
   ✓

Cook's distance: Difference in predicted values
for all observations with & without $i$
where $i$ = each obsev.

3) Multivariate
   (Row - wise)
   especially useful in regression

$Y_i - Y'_i$ = difference b/w these predicted values with and without $i$

Generally mean difference in predicted values $Y_i$
to observation $i$
= Outlier

Outliers package in R

# HANDLING OUTLIERS:

1. Use a robust error metric : Instead of MSE, use MAE
2. Unless sure about that the outlier is entry error, do not drop it.
3. Cap the outlier
4. Transforming data — log or sqrt
5. If outliers form a major chunk in data — ~~they~~ it may point of an underlying phenomenon
   — Model them separately with and without outliers
6. Compare results from there

7. Use a robust Algo for outliers e.g.
   RANSAC (Random Sample Consensus)
   — It takes a subset of data that are inliers
   — Fit the model
   — Test all other data points against the fitted model and add those points that fall within user given tolerance
   e.g. residual user given threshold
   — Iterates until a user given iterations are reached or the tolerance is exceeded.

Question: Ways to handle outliers?

# MULTI-COLINEARITY

## WHEN NOT TO CARE ABOUT MULTI-COLINEARITY?

i) If it's a prediction problem (and not interpretation)

ii) Model is non-linear (e.g. tree based models OR NN)
   ⇓
   By def collinear means "linear dependence" b/w predictors

## WHAT IS MULTI-COLINEARITY?

- Multi-collinearity occurs when predictors are co-related with other predictors

## DIAGNOSIS / DETECTION OF MULTI-COLINEARITY

1. Variance Inflation Factor (VIF) (VIF $\geq 5$ or $> 10$)

$$VIF = \frac{1}{1-R^2}$$

- VIF is defined for individual predictor variable (conceptually), one uses the predictor as the dependent var in a regression model on all other predictors and calculates $1-R^2$ from the regression as the "usable" fraction of that predictor in the full regression model.

- If $1-R^2$ i.e. usable fraction is low e.g 0.)
   then $\frac{1}{1-R^2}$ will become big
   ⇓
   ⇓ coeff is large

2. Standard error of coeff is large

- VIF is not natively defined for categorical data but after transforming categorical data to numerical representation e.g. one hot encoding, VIF can be applied

- One-hot encoding introduces multi-collinearity between diff levels of the categorical data
  e.g. if column = sex/gender then it's value is either M or F
  One-hot encoding Gender-M or Gender-F

| Gender | One-hot encoding Gender-M | Gender-F |
|---|---|---|
| M | 1 | 0 |
| F | 0 | 1 |
| M | 1 | 0 |
| F | 0 | 1 |

These 2 columns are -ve co-related since presence of one makes another absent

WAYS TO HANDLE MULTI-COLINEARITY

1. Correlation Matrix : i) Remove one of the highly co-related vars (correlated vars) (between 2 vars)

   ii) Cluster based on co-relation (e.g. dist b/w var A & var B = correlation coeff) and take 1 from every cluster

2. L1 : Will select one b/w 2 correlated var

3. PCA : Will combine all predictors as orthogonal components

NOTE: Testing collinearity b/w categorical var: Chi-square test
Testing collinearity b/w continuous & categorical var: t-test (if category=2) ANOVA (if category >2)

# FEATURE SELECTION WAYS

1. FILTER: Features are selected based on statistical tests e.g. information gain, correlation coeff, mutual information etc.

2. L1 Regularization

3. WRAPPER : Features are selected based on their performance in a model. E.g. feature importance

HYPERPARAMETER TUNING METHODS

1. RANDOM SEARCH:

- We manually specify a set of possible values for each hyperparam and num_of_iterations

- Random search "randomly" selects one value from each of the possible values of each hyperparam and builds a model and calculates score (offline metric)

  This is done num_of_iterations times and the combination of hyperparams which yields the best score is selected.

- Not all combinations of hyperparam values are searched, limited to num_of_iterations

2. GRID SEARCH

- We manually specify a set of possible values for each hyperparam

- Grid search "exhaustively" selects all the combination of hyperparam values, builds model, calculates score and picks the combination of hyperparam values that yields the best score

- All combinations of hyperparams are searched exhaustively i.e. it is time-consuming and becomes impractical for large number of hyperparams

3. BAYESIAN (`Hyperopt` package uses Bayesian approach)

— Issue with Random & Grid Search:
  ⓪ Each time we try diff hyperparams, we have to train the model, make predictions on validation and then calc. metric
    i.e. evaluating the objective fn to find the score is extremely expensive
  ② Grid & Random Search are completely "uninformed" by past evaluations

— So^n in Bayesian:
  ① Use a "surrogate" model that maps hyperparam values to score   i.e.  P(score | hyperparams)
     The "surrogate" model is much easier to optimize
  ② keeps track of past evaluation results to "inform" the next set of hyperparams to try (explore-exploit)

— Steps:
  1 → Build a "surrogate" prob. model of the objective fn
  2 → Find hyperparams that perform best on this "surrogate" model
  3 → Apply these hyperparams to the true obj. fn
  4 → Update the surrogate model based on results from step #3
  5 → Repeat steps #2-#4 until max iteration is reached

— Bayesian approach can find better hyperparams in less time bc↓ they take into account past evaluations

# NORMS

## Norm of a Vector:
— measures the size/length of vector

— measures the distance from origin to the point x

— maps vectors to non-neg. values (scalar) (x being the vector)

— Mathematically: $L^p$ norm

$$\|x\|_p = \left(\sum_i |x_i|^p\right)^{1/p} \quad \text{for } p \in R, \ p \geq 1$$

## NORM:
— A norm is a function from vector space to non-neg. real numbers that behave like the distance from the origin that satisfies:

— A norm is any function $f$

1) $f(x) = 0 \Rightarrow x = 0$ (definiteness)

2) $f(x+y) \leq f(x) + f(y)$ (triangle ineq.)

3) $\forall \alpha \in R, \ f(\alpha x) = |\alpha| f(x)$ (Absolute homogeneity)
   ↑ scalar

4) $f(x) \geq 0$ for all $x$ (Non-negativity)

## L1 norm:
Let $x = [1, 2, 3]$

$\|x\|_1 = 1 + 2 + 3$

$= 6$

— It is the Manhattan dist. from origin to point identified by x

— L1 norm is commonly used in machine learning when the difference between zero and non-zero elements is very imp. Every time an element of x moves away from zero by $\epsilon$, L1 norm increases by $\epsilon$

## L2 Norm :-

Let $x = [1, 2, 3]$

$$\|x\|_2 = \sqrt{1^2 + 2^2 + 3^2}$$

$$= \sqrt{14}$$

- Also called Euclidean norm

- It is euclidean distance from origin to point identified by $x$

- L2 norm is not preferred in settings where it is important to discriminate b/w elements that are exactly zero and elements that are small but non-zero.

    L2 norm increases very slowly near origin.

E.g. in above if $x = [0.92, 3]$

$$\|x\|_2 = \sqrt{2 + 3^2 + (0.9)^2}$$

$$= \sqrt{13 + (0.9)^2}$$

very small diff to $\sqrt{14}$

For L1 : $[1, 2, 3] \rightarrow 6$

$\quad [0.9, 2, 3] \Rightarrow 5.9 \quad \uparrow$ bigger diff.

- The squared L2 norm is more convenient to work with mathematically and computationally than L2 norm itself. Eg. the derivative of squared L2 norm w.r.t. each element of $x$ each depend only on the corresponding element of $x$ while all of the derivatives of L2 norm depend on corresponding element of $x$ can be calc. with vector entire vector.

## L0 NoRM:

- Count of non-zero elements e.g [0,3,0] → 1
- Not technically a norm since scaling the vector by $\alpha$ does not change the number of non zero entries

## MAX NoRM: ($L\infty$ norm)

- max.abs.magnitude of element in vector e.g. [1,9,3] → 9

## FROBENIUS NoRM

- Used in context of matrix, not vector

- $$\|A\|_F = \sqrt{\sum_{ij} A_{ij}^2}$$

- Analogous to L2 norm but for matrix

- Intuitively, if the matrix is rolled into a vector (1-d) then euclidean norm of this vector = Frobenius Norm

# REGULARIZATION USING L1/L2

$\nearrow$ Ridge

Lasso $\nearrow$

→ To solve over-fitting

→ Types: → Lasso (L1)
→ Ridge (L2)
→ Elastic.Net    (linear combination of L1 & L2)

→ Penalizes complexity of model by adding a penalty term

→ The penalty term in a way acts as a "bias"
(For given error, increasing bias reduces variance which in turn reduces overfitting)

→ Works with any parametric algo  (e.g. neural net)

## L1 Regularization:

$$Loss_{L_1} = \sum (y_i - \hat{y_i})^2 + \lambda \sum |\beta_i|$$

→ will depend on $\beta_i$

SSE

$\lambda$ times sum of abs. value of coeff

$\lambda$ ← hyperparam

→ L1 regularization shrinks some parameters to zero : feature selection

→ Increasing $\lambda$ will cut features one by one until no variable remains ($\beta(coeff.)$)

## L2 Regularization:

$$Loss_{L2} = \underbrace{\sum (y_i - \hat{y}_i)^2}_{\substack{\text{SSE} \\ \text{will depend} \\ \text{on } \hat{\beta}_i}} + \underbrace{\lambda}_{\text{hyperparam}} \underbrace{\sum \hat{\beta}_i^2}_{\substack{\lambda \text{ times } \\ \text{sum of} \\ \text{squares of coeff}}}$$

→ L2 regularization will force the parameters (coeff.) to be relatively small/shrink (by increasing $\lambda$)

→ Practically performs better at prediction than L1

→ L2 is preferred over L1 when
   i) prediction performance is concerned and
   ii) Features have multi-collinearity

When two predictors are highly co-related,
L1 simply picks one of them leaving other.
However in L2, it keeps both of them and
jointly shrinks the corresponding coeff

## ELASTIC NET:

$$Loss_{elastic net} = \sum (y_i - \hat{y}_i)^2 + \alpha \underbrace{\lambda \sum |\beta_i|^2}_{\text{Ridge Penalty}} + \underbrace{(1-\alpha) \lambda |\beta_i|}_{\substack{\text{Lasso} \\ \text{Penalty}}}$$

## SELECTING $\lambda$ & $\alpha$: Cross-Validation

# WHY L1 ENCOURAGES ZERO COEFF. BUT L2 SHRINKS COEFF

Multiple ways of explain:

1) Using loss function optimization

2) Using contour plot

## i) Using Loss Function Optimization:

### L2:

Consider a model with single coeff $\beta$, then

$$L_2 = (y - x\beta)^2 + \lambda\beta^2$$

$$= y^2 - 2xy\beta + x^2\beta^2 + \lambda\beta^2$$

To minimize this eqn, take derivative w.r.t $\beta$ and equate to zero (to get coeff's optimal value)

$$\frac{\partial L_2}{\partial \beta} = 0$$

$$\frac{\partial(y^2 - 2xy\beta + x^2\beta^2 + \lambda\beta^2)}{\partial \beta} = 0$$

$$\Rightarrow 0 + (-2xy) + 2x^2\beta + 2\lambda\beta = 0$$

$$\Rightarrow -2xy + 2x^2\beta + 2\lambda\beta = 0$$

$$\beta(x^2 + \lambda) = xy$$

$$\Rightarrow \beta = \frac{xy}{x^2 + \lambda}$$

To make $\beta = 0$, $\lambda \to \infty$ (since it is addition). And $\therefore \beta$ will be as low as feasible but will not become 0

## $L_1$:

$$L_1 = (y - x\beta)^2 + \lambda|\beta|$$

$$= y^2 - 2xy\beta + x^2\beta^2 + \lambda|\beta|$$

To minimize, take derivative w.r.t. $\beta$ and equate to 0
(to get optimal value of $\beta$)

For demonstration purpose, let $\beta > 0$

$\Rightarrow \quad \dfrac{\partial L_1}{\partial \beta} = 0$

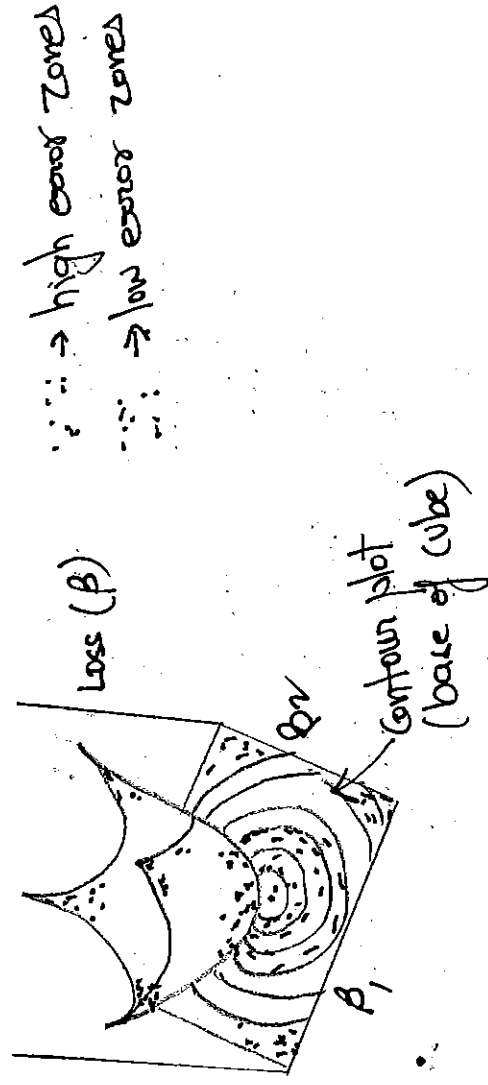$\Rightarrow \quad -2xy + 2x^2\beta + \lambda = 0$

$\Rightarrow \quad 2x^2\beta = 2xy - \lambda$

$\Rightarrow \quad \beta = \dfrac{2xy - \lambda}{2x^2}$

To make $\beta = 0$, anytime $\lambda = 2xy$, the condition will be satisfied.
& this can happen a lot

ii) USING CONTOUR PLOT:

Gradient Descent represented as Contour Plot



Loss (β)

βv

Contour plot
(base of cube)

β₁

$\because \rightarrow$ high Loss Zone

$\because \rightarrow$ low Loss Zone

L2 Regularization Term: $\quad \lambda \beta^2$

— For 2 variables/features, this will translate to $\beta_1^2 + \beta_2^2$

— For illustration/simplicity, let's take $\lambda = 1$

∴ Reg. term for L2 $= \beta_1^2 + \beta_2^2 \quad$ (if $\lambda = 1$)

L1 Regularization Term: $\quad \lambda |\beta|$

— For 2 var/features, this will translate to $\lambda(|\beta_1| + |\beta_2|)$

— For simplicity/illustration, let's take $\lambda = 1$

∴ Reg. term for L1 $= |\beta_1| + |\beta_2| \quad$ (if $\lambda = 1$)

- Error/Loss = bias² + variance + irreducible error

To reduce variance/over-fitting for a given constant error ⟹ bias needs to be increased

In fact, that's what the regularization terms are doing
i.e. increasing bias

∴ Regularization Term = Bias term

- Regularization parameter $\lambda$:
  - Variance depends on the w⃗s $(\beta\text{'s})$
  - Bias depends on the w⃗s $(\beta_0\text{'s}$ bias wⱼ)
  - Now since we are trying to reduce variance (both of which depend on $\beta\text{'s}$) we have to add an additional parameter that can regulate the size of bias term.

  - This regularization param is a hyper param
    else gradient descent will set it to 0
    and travel to the global minimum.
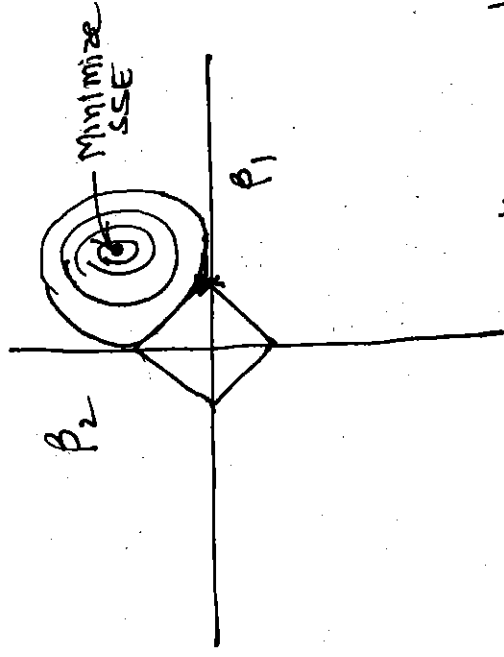    Hence, control of $\lambda$ cannot be given to gradient descent and needs to be kept out of gradient descent.

- If bias is constant/same (i.e. comb. of $\beta_1$ & $\beta_2$ that generate the same bias)
  L2:  $\beta_1^2 + \beta_2^2 = $ constant → Circle of radius constant
  L1:  $|\beta_1| + |\beta_2| = $ constant → Diamond
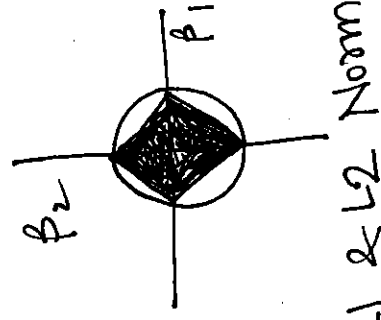
L1 Norm with GD contour



L2 Norm with GD contour



L1 & L2 Norm

2 Forces at play here:

i) Bias term pulling $\beta_1$ & $\beta_2$ to lie on circle (L2) or diamond (L1) of cost fn (without bias term) indicated by ^ minimum

ii) GD trying to travel to global minimum indicated by green dot

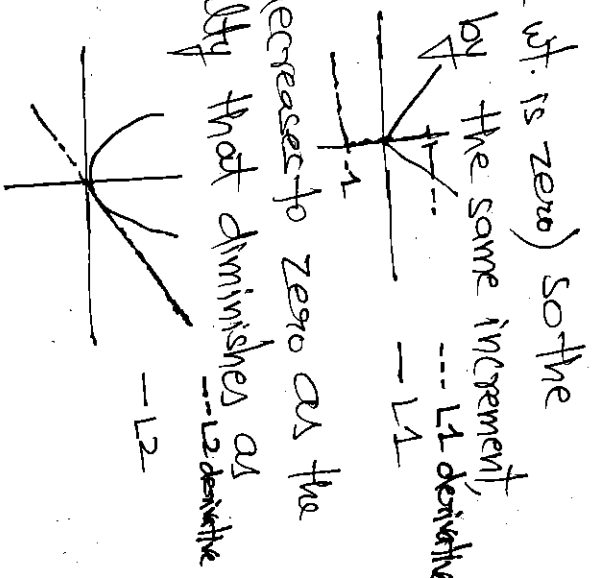Both the forces pull and finally settle near the point of intersection indicated by 'black cross' (x).

iii) For the same amount of bias term generated, the area occupied by L1 norm is smaller and closer to axes (compared to L2 norm. This is what causes the point of intersection b/w L1 norm & GD contour to converge/intersect on the axes leading to feature selection

# DIFFERENCE BETWEEN L1 & L2 REGULARIZATION:

- L1 reg. penalizes the sum of abs. values of wts whereas L2 reg. penalizes the sum of squares of wts

- L1 reg. solution is sparse whereas L2 reg. solution is not sparse

- L1 reg. performs feature selection whereas L2 reg. doesn't perform feature selection (shrinks wts close to zero, not zero)

- L1 reg. is robust to outliers whereas L2 is not robust to outliers

## WHY L1 make wts. sparse but L2 does not?

- Gradient of L1 is -1 or 1 (except where the wt is zero) so the penalty moves the value closer to zero by the same increment, regardless of wt's current value
  whereas

- Gradient of L2 is a linear fn that decreases to zero as the wt approaches zero, resulting in a penalty that diminishes as the wt value becomes smaller



- -- L1 derivative
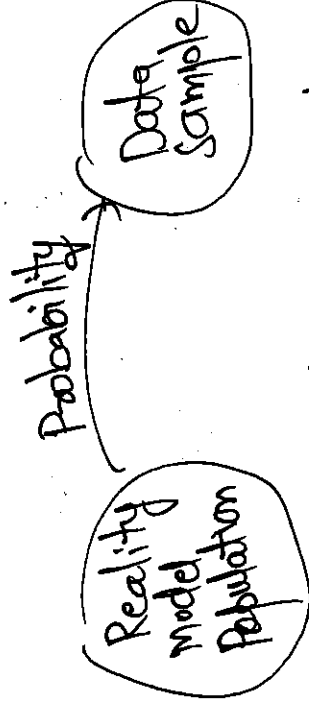- -- L1

- -- L2 derivative
- -- L2

## Which is better L1 or L2?

- L1 is more robust : L1 takes abs. value of wts so the cost of outliers only increases linearly whereas L2 takes square of the wts so cost of outliers increases exponentially

- L2 is computationally less expensive : L2 takes square of wts has closed soln and can be solved in terms of matrix math whereas L1 takes abs. value of wts has no closed form soln (non-differentiable) and cannot be solved in terms of matrix math
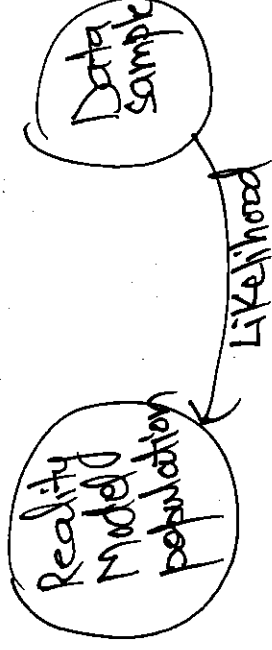
# LIKELIHOOD & PROBABILITY:

## PROBABILITY:



Reality / Model / Population →(Probability)→ Data Sample

e.g. What is the chance of observing particular data or sample given a specific model or population (whose parameters we know)

Knowing parameters → Prediction of outcome

Prob. (data | distribution)

## LIKELIHOOD:



Reality / Model / Population ←(Likelihood)— Data Sample

e.g. Given observed data, what is the chance that a given reality or model is true?

If you observe x, what is the best/most likely distribution with its parameters? possible

Observation of data → Estimation of parameters

L (distribution | data)

The likelihood of true parameters being a certain value given data
= Prob. of observing the data given some true parameters θ

$$L(\theta|x_i) = \prod_{i=1}^{n} P(x_i|\theta)$$

probability of observing $x_1, x_2, \ldots x_n$ given θ

$$P(x_1, x_2, x_3, \ldots x_n|\theta) = \text{# of samples are independent}$$
$$= P(x_1|\theta) \cdot P(x_2|\theta) \ldots$$
$$= \prod_{i=1} P(x_i|\theta)$$

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

# MLE (MAXIMUM LIKELIHOOD ESTIMATION)

- Maximizing the likelihood fn i.e. Max. $L(\theta | x_i) = \prod_i f(x_i | \theta)$
- Find the parameter values that make the observed data most likely.
- E.g. if we assume the p.d.f. is normal then given data what is the best estimate of $\mu$ and $\sigma^2$ (parameters of normal distribution)

- To make a new prediction, we simply evaluate p.d.f using the best parameters found
- The p.d.f. can be binomial, gaussian, exponential etc.
- MLE does not tell "how" to find the optimal value of $\theta$
  → it just tells how one value of $\theta$ is "more likely" than other
  - optimization algos like $\hat{\theta}$ MLE find this optimal value ("how")

# EXPECTATION MAXIMIZATION:

- EM is a technique to compute MLE
- EM is MLE both hidden states
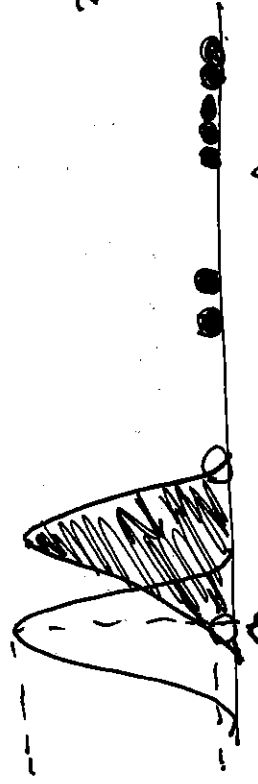- EM can work with any distribution

- MLE defines the objective function whereas EM solves it in iterative way.

# GAUSSIAN MIXTURE MODELS:

- Probabilistic model for representing normally distributed
  - Hint for using Mixture model: data looks multimodal i.e. there is more than one peak in the chart of data.
- Mixture models in general don't require knowing which subpopulation a data point belongs to allowing the model to learn the subpopulations automatically
- Unsupervised

# EM & GMM

EM is used to estimate the parameters of GMM

1) Start with arbitrary **initial** distribution parameters $(\mu, \delta^2)$
   for a given no. of distribution

2) Compute the likelihood that each distribution produced that data point → Expectation step

This point has higher likelihood of belonging to white (its pdf)

Likelihood of white = pdf on white gaussian distribution

Likelihood of black = pdf on black gaussian distribution

Total likelihood = Likelihood of white + Likelihood of black

black wt = likelihood of black / Total likelihood

white wt = likelihood of white / Total likelihood

3) Now take all data points having black position
   & compute new $\mu$ & $\delta^2$ → Maximization step

$$\mu = \frac{\sum wt \times data\ point}{\sum wt}$$

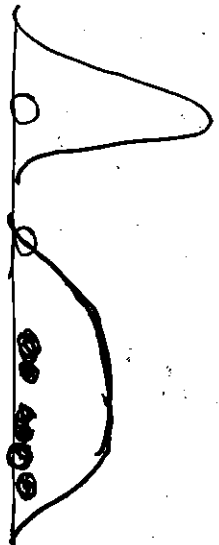$$\delta^2 = \frac{\sum (wt \times (data - \mu))}{\sum wt}$$

Similarly do it for white data points

4) The distribution will shift

5) Now remove/forget what position of which data point belonged to which distribution i.e. start fresh with updated distribution parameters

6) Again calculate likelihood of each data point with respect to updated distributions

ie) Repeat steps 2 & 3 until a threshold of ... for likelihood is reached



## RELATION TO K-MEANS:

K-Means performs "hard" assignment of data points to clusters, in which each data pt is associated uniquely with one cluster

EM algo makes a "soft" assignment of data pts to all clusters probabilistically

## BAYESIAN or MAXIMUM A POSTERIORI (MAP):

- We have some knowledge of data (prior)

- There is no point estimate which best explains the data instead there are multiple values of a parameter e.g. $\mu$ & $\sigma^2$

- This gives us multiple models with the same prior

- To predict a new example/data point we have to compute "weighted sum" of these predictions

# MLE & MAP

→ Both MLE & MAP are used to estimate some variable in the setting of probability distributions.

→ The likelihood fn, in case of MLE

$$\theta_{MLE} = \arg\max_{\theta} \Pi_i P(x_i|\theta)$$

$\theta \to$ parameter we want to infer

$x_i \to$ data

Taking log, as each $P(x_i|\theta)$ is a no. between 0 and 1 and multiplying them together when no. of $x_i$'s are large (approach to infinity) → will result in underflow.

$\Pi \to$ product changed to $\Sigma \to$ summation bcz of log

$$\log xy = \log x + \log y \quad ①$$

$$\theta_{MLE} = \arg\max_{\theta} \log \Pi_i P(x_i|\theta)$$

$$= \arg\max_{\theta} \Sigma \log P(x_i|\theta)$$

→ The Posterior fn, in case of MAP

Using Bayes Rule:

$$\underbrace{P(\theta|x)}_{\text{Posterior}} = \frac{P(x|\theta) P(\theta)}{P(x)}$$

ignoring the normalization constant in denominator

$$\text{Maximizing} \quad \propto P(x|\theta) P(\theta)$$

likelihood · Prior

$$\theta_{MAP} = \arg\max_\theta \prod_i P(x_i|\theta) P(\theta)$$

Taking log,

$$= \arg\max_\theta \log \prod_i P(x_i|\theta) P(\theta)$$

$$= \arg\max_\theta \sum_i \log P(x_i|\theta) P(\theta) \quad ②$$

① and ②, the only difference is
① prior $P(\theta)$ in MAP

⇒ This means that the likelihood is now weighted
⇒ some wt coming from prior
with

∴ MLE is a special case of MAP

→ if $P(\theta)$ = constant or uniform [but not some distribution like gaussian where depending on the region of distribution probability is high or low ie. never always the same]
Then we can ignore the constant term

↳ $= \theta_{MAP}$

↳ MLE equivalent to Ridge regression when i) weight prior belong to Normal dist. ii) error belong to Normal dist.

NOTE: It is important to remember that MLE or MAP assumes samples/distributed are independent then only
$P(x_1, x_2, x_3...x_n|\theta) = P(x_1|\theta) \cdot P(x_2|\theta)...P(x_n|\theta) = \prod_i P(x_i|\theta)$

NOTE: If you have some idea about the parameter → use MAP with that you are estimating → idea that you have abt the parameter the prior where prior ≠ idea abt the parameter

NOTE: Maximizing log likelihood = Minimizing deviance = -2 log-likelihood fn (Minimizing squared error (if the errors belong to normal dist))

→ Many of the penalized max. likelihood techniques.
(regularization)

$\parallel$

MAP with certain parameter priors

i) Quadratic wt. decay (shrinkage, L2) → Gaussian prior

ii) Absolute wt. decay (lasso, L1) → Laplace prior

# GENERALIZED LINEAR MODELS

- extension of linear regression

- linear regression assumes the error terms/residuals are normally distributed

  vs

  Generalized linear models assume the outcome var. $y_i$ are not normally distributed

  ✓ $y_i$ generated from a particular dist. $y_i$

  ✓ exponential family ( e.g. Normal, Poisson, Binomial/etc.)

- GLM generalizes linear regression by allowing the error terms/residuals the outcome var. $y_i$ linear model to be related to the response via a "link function"

─⌐ "Linear" bcz : predictors affect the dist. of outcome only thru linear

  comb. of $x_i \cdot \beta$

─⌐ Link Fn :   Link fn transforms this linear (usually non-linear)   comb. of predictors into outcome's

  space     e.g. logistic fn

- Example: logistic Regression, Linear Regression

— Linear Regression may be viewed as a special case of generalized linear model with identity link and regressed normally distributed and error terms are normally distributed

# GENERATIVE vs DISCRIMINATIVE MODELS

→ learns the joint probability distribution $P(x,y)$

→ learns the conditional probability distribution $P(y|x)$

Example:

$(1,0) \cdot (1,0) \quad (2,0) \quad (2,1)$

$P(x,y)$:

|       | $y=0$ | $y=1$ |
|-------|-------|-------|
| $x=1$ | $2/4$ | $0$   |
| $x=2$ | $1/4$ | $1/4$ |

$= \dfrac{\text{No. of times } (x,y) \text{ appears}}{\text{Total no. of } (x,y) \text{ pairs}}$

$P(y|x)$:

|       | $y=0$ | $y=1$ |
|-------|-------|-------|
| $x=1$ | $2/2$ | $0/2$ |
| $x=2$ | $1/2$ | $1/2$ |

$= \dfrac{\text{No. of times } (x,y) \text{ appears}}{\text{Total no. of } (x) = \text{given y}}$

→ Models how the data was generated in order to categorize a signal

Does not care how the data was generated, tries to find the "discrimination" between classes.

→ E.g. Naive Bayes
       GMM

E.g. Linear Regression
     Logistic Regression
     Decision Trees
     SVM
     Neural Networks

we are only interested in $P(y|x)$

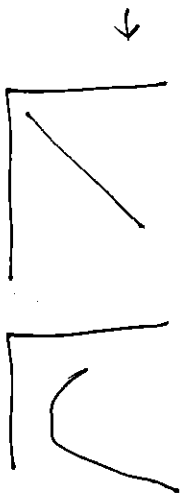NOTE: $P(x,y)$ can be transformed into $P(y|x)$ using Bayes conditional probability rule :

$$P(x,y) = P(y|x) \cdot P(x)$$

but $P(x)$ calculation can be an extra step and irrelevant for the task in hand, so directly calculate $P(y|x)$
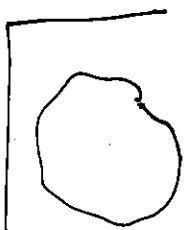
# LINEAR VS NON-LINEAR MODELS

→ Linear in "parameters" (not predictors)

→ E.g. $y = \beta_0 + \beta_1 x_1 + \beta_2 x_1 x_2$

→ E.g. Linear Regression, Logistic Regression

→ Decision boundary is linear in case of classification
→ E.g. Linear Regression, Logistic Regression

→ everything else

E.g. $y = e^{ax}$

Tree based Models, Neural networks (bcz of non linear activation fns)

---

## PARAMETRIC VS NON-PARAMETRIC MODELS

→ No. of parameters are finite

→ E.g. Linear Regression, Logistic Regression, Linear SVM

→ "There is a distribution that the data follows"

→ No. of parameters are infinite
i.e. Complexity of model grows with training data

E.g. KNN, Decision Trees, RBF kernel SVM, Neural networks

→ "Distribution-free model"

# CLASSIFICATION METRICS

Actual

|  | 1 | 0 |
|---|---|---|
| Predicted 1 | TP | FP |
| 0 | FN | TN |

**1) Accuracy :** $\dfrac{(TP+TN)}{TP+FP+TN+FN}$

**2) Precision:** $\dfrac{TP}{TP+FP}$ : If your classifier predicts 1, how likely it is to be true

**3) Recall :** $\dfrac{TP}{TP+FN}$ : Out of all the +ve classes, how much did your classifier catch

Sensitivity/ TPR

**4) F1 score :** $\dfrac{2}{\frac{1}{P}+\frac{1}{R}} = \dfrac{2PR}{P+R}$

Why harmonic mean, not simple average?

1) If $P=1$ & $R=0$

| | 1 | 0 |
|---|---|---|
| 1 | 100 | 0 |
| 0 | 10,000 | some value |

$Avg = 0.5$
which is not good, as there are lots of FN

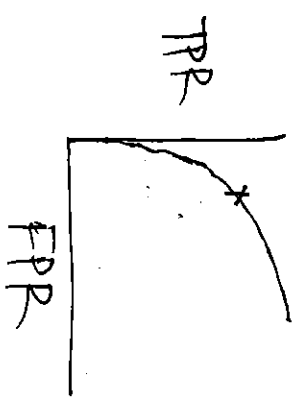2) Since numerator is same for both precision and recall, it makes sense to take avg of denominator

3) Harmonic mean (encourages similar values for precision and recall ie more the precision and recall deviate from each other, the worse the harmonic mean)
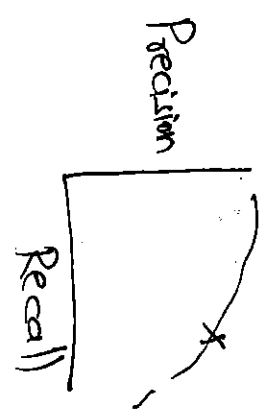
5) Specificity: $\dfrac{TN}{TN+FP}$

6) FPR : $1 -$ specificity $= \dfrac{FP}{TN+FP}$

7) TPR : $\dfrac{TP}{P} = \dfrac{TP}{TP+FN}$    ROC vs PR

ROC Curve



Point that maximizes area under ROC curve → top left
$[\,tpr - (1-fpr) \geq 0\,]$

Both are derived from confusion matrix by varying the thresholds, so every point on ROC curve has a corresponding point on PR curve

PR Curve



Point that maximizes area under PR curve → top right
$[\,precision - recall \geq 0\,]$

Recall = TPR to the ... & Precision

PR curve must be preferred over ROC curve when:

1) only concerned with +ve class (no TN in PR curve))
→ ROC curve is immune to imbalanced data (doesn't change)

2) class imbalance : if +ve class is less in number and -ve class dominates in dataset
→ FPR is negative class only metric → change in FP+TN instances
ie. change $\overline{FPR = \frac{FP}{FP+TN}}$ instances $(\alpha 1 - ve \text{ instance})$
→ Thus, if underlying FP+TN change ...
data doesn't change, RoC doesn't change ... in imbalanced datasets TPR

# ROC Curve

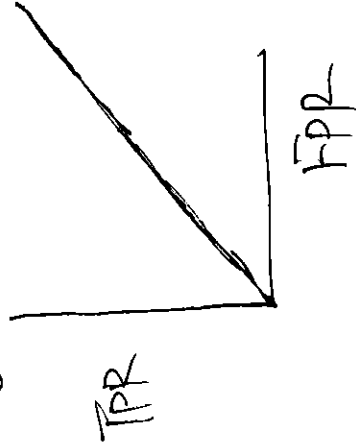TPR & FPR are defined at every threshold

ROC curve is plotted by calculating TPR & FPR at every threshold

Interpretation (AUC): probability that a random +ve is assigned a higher score than a random -ve

AUCROC is used to compare classifiers independent of threshold

Classifier with no discrimination ability b/w +ve & -ve class will have ROC curve = diagonal (even in imbalanced datasets)
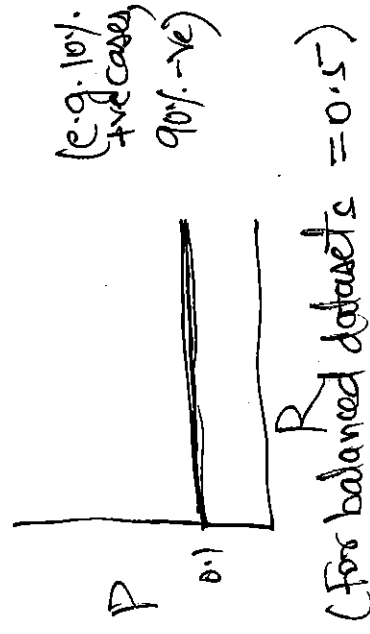


TPR vs FPR

# PR Curve

P & R are defined at every threshold

PR curve is plotted by calculating P & R at every threshold

Interpretation (AUC): Average precision, where the avg. is taken across all thresholds

AUCPR is used to compare classifiers independent of threshold

Classifier with no discrimination ability b/w +ve & -ve class will have PR curve: = horizontal line at = ratio. proportional to +ve cases in dataset



P vs R, line at 0.1 (e.g. 10% +ve cases, 90% -ve)
(for balanced datasets = 0.5)

# UNDERSTANDING THE DIFFERENCE B/W ROC & PR Curves

ROC Curve used : TP, TN, FP & FN ie. all elements of
✓      confusion matrix. ⇒ ROC curve focused on
both classes

PR Curve used : TP, TN & FP ie. one element FN is
✓      conveniently left out of confusion matrix.
+ve & -ve class is left out of confusion matrix
imbalanced datasets is in abundance for
focused on "minority" class ⇒ PR curve
focused on "minority" class

Even if distribution of +ve class & -ve class changes
(ie. proportion of +ve to -ve instances), ROC AUC does not
change as it employs all elements of confusion
matrix and we use TP rate & FP rate which is
static columnar ratio in confusion matrix so
do not depend on class distributions.

However, when class distribution changes, PR AUC
changes as it does not take into account FN

Even if distribution in imbalanced datasets, ROC AUC
is high but PR AUC is low since PR AUC
only focuses on "minority" class.

It is possible that in imbalanced datasets, ROC AUC
is high but PR AUC is low since PR AUC
only focuses on "minority" class.

Note that: PR curve & ROC curve have one axis in
common ie. TPR = Recall

# RANKING METRICS

1. PRECISION @ K = $\dfrac{\text{Number of recommended item @ K that are relevant}}{\text{Number of recommended items @ K}}$

E.g.

| Rank-1 | 2 | 3 | 4 | 5 |
|--------|---|---|---|---|

— Relevant
— Not relevant

$P@1 = \dfrac{1}{1}$

$P@2 = \dfrac{1}{2}$

$P@3 = \dfrac{1}{3}$

$P@4 = \dfrac{2}{4}$

$P@5 = \dfrac{3}{5}$

—Disadv : Fails to take into account the relative ordering of relevant docs in top K (order unaware metric)

e.g. in above example, first 3 items could be relevant and last 2 could be not relevant but still P@K would be 3/5

— Range [0,1]

— Mostly used for binary relevance. Can be adapted to graded relevance (i.e. non-binary relevance, e.g. relevance score) by thresholding and converting to binary relevance

2. RECALL @ K = $\dfrac{\text{No. of recommended Item@k that are relevant}}{\text{Total no. of relevant items}}$ (in the recommended list since we are mostly interested in recommended top-K items)

— relevant
— not relevant

E.g.
Rank = 1    2    3    4    5

Recall @1 = $\dfrac{1}{3}$

Recall @2 = $\dfrac{1}{3}$

Recall @3 = 2/3

Recall @4 = 3/3

Recall @5 = 3/3

→ Disadv: 1) Fails to take into account relative ordering of relevant docs in top K (order unaware metric)
   E.g. in above example, if first 2 docs were relevant and next 1 doc were not relevant then also Recall @3 would be 2/3

2) Another disadv: By increasing K to N i.e. K = Total No. of recommended items(N) we can return a perfect score everytime harder to score well with Recall @K metric

3) Smaller K value makes it

— Range [0,1]
— Mostly used for binary relevance. Can be adapted to graded relevance by thresholding and converting to binary relevance.

# 3. AVERAGE PRECISION:

— Metric that tells how much of the relevant docs are
  √ concentrated in the highest ranked predictions

—
$$\frac{\sum\limits_{k=1}^{n} P@K \times relevance(K)}{No. \ of \ relevant \ docs}$$

where, $relevance(K) = \begin{cases} 0 \\ 1 \end{cases}$
depending on
whether the doc at
rank k is relevant
or not

  √

— This metric is able to give more wt. to errors that
  happen high up in the recommended list. Conversely,
  √ it gives less wt to errors that happen deeper in
  the recommended list

— One of the methods to calculate Area Under P-R Curve
  (other) method is using trapezoidal rule) for one class [obj: detern
                                           use case]

— Defined for 1 user / list.

E.g.

— relevant
— not relevant

User 1.
list 1

$\frac{1}{1} \times 1$   $\frac{1}{2} \times 0$   $\frac{2}{3} \times 1$   $\frac{3}{4} \times 1$   $\frac{3}{5} \times 0$

$$\frac{\frac{1}{1} + \frac{2}{3} + \frac{3}{4}}{3} = 0.8$$

Now, if it changes to

$\frac{1}{1} \times 1$   $\frac{2}{2} \times 1$   $\frac{2}{3} \times 0$   $\frac{3}{4} \times 1$   $\frac{3}{5} \times 0$

$$\frac{\frac{1}{1} + \frac{2}{2} + \frac{3}{4}}{3} = \frac{2.75}{3}$$

$$= 0.916 \quad (increased)$$

— Range: [0, 1]
— Can be adapted to graded relevance by thresholding

# 4. MAP (Mean Avg. Precision)

- Avg. precision Over a set of Users/lists

$$\frac{\sum_{u=1}^{U} \text{Avg. Precision } (u)}{\text{Number of Users/lists}} \quad \text{where } U = \text{\# of Users}$$

- ▨ relevant
- ☐ not relevant

User 1:   ▨ ▨ ▨ ▨ ▨
$\frac{1}{1}\times 1$   $\frac{1}{2}\times 0$   $\frac{2}{3}\times 1$   $\frac{3}{4}\times 1$   $\frac{3}{5}\times 0$

$$= \frac{1+\frac{2}{3}+\frac{3}{4}}{3} = 0.8$$

User 2:   ▨ ▨ ☐ ▨ ▨
$\frac{0}{1}\times 0$   $\frac{0}{2}\times 0$   $\frac{0}{3}\times 0$   $\frac{1}{4}\times 1$   $\frac{2}{5}\times 1$

$$\frac{\frac{1}{4}+\frac{2}{5}}{2} = 0.325$$

$$\text{MAP} = \frac{0.8 + 0.325}{2} = 0.56$$

- Defined for a group of users/lists

- **Mean Area Under PR curve** for all classes of each obj. [obj. detection be case]
- **Con:)** Since all classes are not of interest (e.g. class 0)
  on equal in imbalanced dataset) it may not give complete/accurate picture
  2) Primarily for binary relevance
- Range: [0,1]
- Can be adapted to graded relevance by thresholding

# 5. RECIPROCAL RANK: Inverse rank where the first relevant doc is found

E.g.

▨ ▢ ▢    RR = $\frac{1}{3}$

▢ ▨ ▢    RR = $\frac{1}{2}$

▢ ▢ ▢    RR = 0

✓ Con: Only cares about the first relevant doc's rank

— [0,1] Range
— Can be adapted to graded relevance by thresholding

# 6. MEAN RECIPROCAL RANK

✓ Avg. of RR for multiple users / lists

▨ ▢ ▢    RR = 1

▢ ▢ ▨    RR = $\frac{1}{3}$

▢ ▨ ▢    RR = $\frac{1}{2}$

$$MRR = \frac{1 + \frac{1}{3} + \frac{1}{2}}{3} = 0.61$$

✓ — Con: Only cares about first relevant doc's rank

— [0,1] Range
— Can be adapted to graded relevance by thresholding

∴ All the previous 6 metrics are usually used for binary relevance

# 7. NORMALIZED DISCOUNTED CUMULATIVE GAIN

- Used when relevancy is not binary, instead it is a real number (can also be used for binary relevance)

- Gives more importance to correctly predicted ranks
  → at top and "discounts" mistakes as you go down the ranks

- Compared to MAP, NDCG further tunes the recommended lists evaluation. Since relevance is a real number, instead of binary, it is able to use the fact that some documents are "more" relevant than others.

$$CG_p = \sum_{i=1}^{p} rel_i$$

$$DCG_p = \sum_{i=1}^{p} \frac{rel_i}{\log_2(i+1)}$$

DCG at position $p$
(not normalized hence difficult to compare)
e.g. longer list will have higher DCG compared to shorter list

$$DCG_p = \sum_{i=1}^{p} \frac{2^{rel_i} - 1}{\log_2(i+1)}$$ (used in industry)

where $p$ = # of elements in the recommended list

$rel_i$ = graded relevance of result at position $i$

$\log_2(i+1)$ = logarithmic reduction factor to penalize proportionally to the position of result

$$IDCG_p = \sum_{i=1}^{|REL_p|} \frac{rel_i}{\log_2(i+1)}$$

ideal DCG at position $p$

or $$\sum_{i=1}^{|REL_p|} \frac{2^{rel_i} - 1}{\log_2(i+1)}$$

$|REL_p|$ = list of relevant docs sorted by relevance upto position $p$ (high rel to low rel.)

- IDCG is "ideal DCG" so the formula should remain consistent

$$\rightarrow NDCG_p = \frac{DCG_p}{IDCG_p}$$

- Per User/list metric, we need to avg it out for all users in test set

- Range $[0,1]$

E.g.
Predicted Ranking with True Ratings

True Rating $\boxed{3}$

Predicted Rank=1

| | 2 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | $\boxed{3}$ | $\boxed{2}$ | $\boxed{3}$ | $\boxed{0}$ | $\boxed{1}$ | $\boxed{2}$ |

$$CG_6 = 3+2+3+0+1+2 = 11$$

$$DCG_6 = \frac{rel_i}{\log_2(i+1)}$$
(using)

$$DCG_6 = \frac{3}{\log_2(1+1)} + \frac{2}{\log_2(2+1)} + \frac{3}{\log_2(3+1)} + 0 + \frac{1}{\log_2(5+1)} + \frac{2}{\log_2(6+1)}$$
$$= 6.8$$

To get $IDCG_6$, sort the list according to true rating

Rank=1
| 3 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|
| $\boxed{3}$ | $\boxed{3}$ | $\boxed{2}$ | $\boxed{2}$ | $\boxed{0}$ | $\boxed{0}$ |

$$IDCG_6 = \frac{rel_i}{\log_2(i+1)}$$
(using)

$$IDCG_6 = \frac{3}{\log_2(1+1)} + \frac{3}{\log_2(2+1)} + \frac{2}{\log_2(3+1)} + \frac{2}{\log_2(4+1)} + \frac{1}{\log_2(5+1)} + 0$$
$$= 7.14$$

$$NDCG_6 = \frac{DCG_6}{IDCG_6} = \frac{6.8}{7.14} = 0.961$$

CONS OF NDCG: (occurs when unequal size of returned list)

1) Does not penalize for <u>missing docs</u> in the recommended list

Query 1 returns : 1,1,1

Query2 returns : 1,1,1,1,1

$ND(G_3 \ (query 1) = NDCG_5 \ (query 2)$

To fix this: i) enforce fixed size of result set
ii) use minimum scores for missing docs

Query 1 : 1,1,1,0,0

Query 2 : 1,1,1,1,1

$NDCG_5 \ (query 1) < NDCG_5 \ (query 2)$

2) Does not penalize <u>bad docs</u> in the recommended list

Query 1 returns : 1,1,1

Query 2 returns : 1,1,1,0

$NDCG_3 \ (query 1) = NDCG_4 \ (query 2)$

# NLP METRICS (~~~~~~~~~~~~~~) → should be defined based on task (task-specific)

## 1. BLEU (Bilingual Evaluation Understudy)

✓ → n-gram overlap b/w O/P sentence & Reference sentence

- n-gram can be unigram, bigram, trigram, ...
- precision kind of metric, range [0,1] (from O/P to ref)

✓ → CLIP = $\dfrac{\text{No. of words matched}}{\text{No. of words in O/P}}$ (from O/P to ref)

Clips no. of words matched by no. of times the word appeared in reference

$\dfrac{O/P}{\text{school school}}$

### Example

1. एक लड़का स्कूल जाता है

Ref: I am going to school

If no CLIP in definition = $\dfrac{2}{2} = 1$ → misleading

If CLIP, BLEU-1 = $\dfrac{1}{2} = 0.5$

2. एक लड़का स्कूल जाता है

Ref: I am going to school

BLEU-1 = $\dfrac{2}{3}$

for unigram matches

$\dfrac{O/P}{\text{school here}}$

→ penalizes extra word in O/P

→ penalizes correct words from O/P is not penalized

for unigram matches

— Issues: 1) missing correct words from O/P ✓
  2) Semantics ✗
  3) Order of words is not taken into account for unigram matches (length penalty)

soln: n-gram

— Usually BLEU score = $\dfrac{\text{geometric mean of all 4 n-gram precisions i.e unigram, bigram, trigram & quad-gram}}{}$

$= \sqrt[4]{P_1 P_2 P_3 P_4}$

2. METEOR:  Overcomes drawbacks of BLEU, which are
   1) Does not take recall into account
   2) Allows only exact n-gram matching
   METEOR does synonym/stemmed match + computes recall as
   finally computed F-score well

3. ROUGE
   - Recall related metric [0,1]
   - ROUGE-1 = $\dfrac{\text{No. of word matches}}{\text{No. of words in reference}}$
     recall

   ROUGE1 = $\dfrac{\text{No. of word matches}}{\text{No. of words in O/P}}$
   precision

   - In practice, ROUGE1 is the F1 score calculated on top of
     ROUGE-1 precision & recall

   - Similarly for bigram, trigram : ROUGE-2, ROUGE-3

   - ROUGE-L : - Least Common subsequence matching
               - no exact m-gram matching
               - same order but not necessarily continuous

4. PERPLEXITY : - Used in lang. generation tasks
   - How confused/uncertain the model is in generating
     next token/text
   - Lower is better, range : [0,1]
   - $e^{\text{cross-entropy}}$

5. (cosine similarity b/w
   o/p & reference) can be
   used as a metric as well

   Let's say vocab has 6 words and prob. distribution of
   next word is available given a word.
   e.g. P(w$_2$|a)
   ∴ Prob. of sentence a red fox
   = P(a) * P(red|a) * P(fox|a red)
   finally, normalize it by length of sentence/text

|       | a | fox | red |
|-------|------|------|------|
|       | 0.01 | 0.4  | 0.2  |
| a fox | 0.08 |      |      |

# VECTORS:

$$\begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$= 3\hat{i} + 2\hat{j}$$



# BASIS VECTORS:

Vectors that define a co-ordinate system $\hat{i}$ & $\hat{j}$ are the basis vectors of XY co-ord. system

technical defn: → set of linearly independent vectors that span the full space



Hence whenever we specify a vector eg $\begin{bmatrix} 3 \\ 2 \end{bmatrix}$, implicitly there is a notion of basis vectors and the vector we specified is a "scaled" version of basis vectors

$$\begin{bmatrix} 3 \\ 2 \end{bmatrix} = 3 \text{ times } \hat{i} + 2 \text{ times } \hat{j}$$

↓ basis vector in x-dir

↓ basis vector in y-dir

# LINEAR COMBINATION OF VECTORS:

Let $\vec{v}$ & $\vec{w}$ be two vectors then their linear combination is defined as

$$a\vec{v} + b\vec{w}$$

scalar

(note: scalar multiplication + vector addition)

If $a$ & $b$ can take any values (real numbers) 3 possibilities

i) $a = 0$, $b = 0$

$\Rightarrow a\vec{v} + b\vec{w} = $ origin

ii) $a \neq 0$, $b \neq 0$ and $\vec{v}$ & $\vec{w}$ do not line up

$\Rightarrow$ you can reach all pts on a plane

(iii) $a \neq 0$, $b \neq 0$ and $\vec{v}$ & $\vec{w}$ line up

$\Rightarrow$ you can reach all pts on the same line as $\vec{v}$ & $\vec{w}$

# SPAN OF VECTORS:

The span of $\vec{v}$ & $\vec{w}$ is the set of all its reachable via their linear combinations

$$a\vec{v} + b\vec{w}$$

Scalars that can be any real number

If $\vec{v}$ & $\vec{w}$ do not line up ( point in diff. directions)

⇒ Span → 2d plane

Similarly,
For 3 vectors $\vec{v}$, $\vec{w}$ & $\vec{u}$ , if they do not line up

⇒ Span → 3d space

If two of these vectors line up

⇒ Span → 2d plane.

If all three of these vectors line up

⇒ Span → line

LINEAR DEPENDENCE: In a group of vectors, if a vector does not add anything to the span or is redundant ( since it is a scaled version by any other vector or linear comb. of other vectors) → Linearly dependent vector

$$\vec{u} = a\vec{v} + b\vec{w}$$

Linearly dependent → linear comb.

# LINEARLY INDEPENDENT:

If a vector adds another dimension to span imagrb of vectors or cannot be expressed as a linear combo. of other vectors in the grph

→ Linearly Independent vector

$$\vec{u} \neq a\vec{v} + b\vec{w}$$

For all values of a & b

# MATRICES AS LINEAR TRANSFORMATION OF SPACE:

Recall that in a co-ord. system, $\hat{i}$ & $\hat{j}$ are the basis vectors



in matrix format

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$\hat{i}$ co-ord.   $\hat{j}$ co-ord

Now a matrix $\begin{bmatrix} 1 & 3 \\ -2 & 0 \end{bmatrix}$ represents a linear

transformation of $\hat{i}$ & $\hat{j}$   where new $\hat{i} = 1(\text{old } \hat{i}) - 2(\text{old } \hat{j})$

new $\hat{j} = 3(\text{old } \hat{i}) + 0(\text{old } \hat{j})$



$\text{new} (\hat{i}) \begin{bmatrix} 1 \\ -2 \end{bmatrix}$   $\text{new } \hat{j} \begin{bmatrix} 3 \\ 0 \end{bmatrix}$

## Note:

"Linear" in linear transformation: 1) origin is fixed before
& after transformation of
space
2) new $\hat{i}$ & new $\hat{j}$ is still a

3) grid lines are parallel & evenly spaced    line and not curved
even in the transformed space

# MULTIPLICATION OF VECTOR & MATRIX

(in the original position)

— Can be thought of as where the vector lands $\hat{}$ after the transformation of space represented by matrix

Let $v = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$ & $M = \begin{bmatrix} 1 & 3 \\ -2 & 0 \end{bmatrix}$

$\vec{v} \downarrow$ means

$\vec{v} = -1\hat{i} + 2\hat{j}$

new $\hat{i}$ lands $\xleftarrow{}$ new $\hat{j}$ lands

new $\hat{i} = \begin{pmatrix} 1 \\ -2 \end{pmatrix} \quad \downarrow \text{means}$

new $\hat{i} = 1(\text{old } \hat{i}) + -2 (\text{old } \hat{j})$
new $\hat{j} = 3(\text{old } \hat{i}) + 0 (\text{old } \hat{j})$

→ Transformed $\vec{v} = -1 \, (\text{transformed } \hat{i}) + 2 \, (\text{transformed } \hat{j})$

→ multiplication of matrix = $-1 \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 2 \begin{bmatrix} 3 \\ 0 \end{bmatrix}$

$\Downarrow$ multiplication of matrix by a vector

$= \begin{bmatrix} -1 \\ 2 \end{bmatrix} + \begin{bmatrix} 6 \\ 0 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$

# MULTIPLICATION OF MATRIX BY MATRIX:

Since matrix represents a transformation of space

matrix multiplication can be thought of as

consecutive transformation of space i.e. one transformation

followed by another

## Condition:

$m \times n \quad n \times p \implies m \times p$

$$\underline{\#cols} = \underline{\#rows}$$

Dot product ↗↘

Eg.

$$\begin{bmatrix} 2 & 3 & 1 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 1 \\ 8 & 9 \\ 0 & 10 \end{bmatrix} = \begin{bmatrix} & \\ 58 & \end{bmatrix}$$

$2 \times 3 \qquad 3 \times 2 \qquad 2 \times 2$

FORMAL DEF$^n$ of LINEARITY (Linear Transformation)

Let L be the transformation:

i) Additivity: $L(\vec{v} + \vec{w}) = L(\vec{v}) + L(\vec{w})$

ii) Scaling: $L(c\vec{v}) = cL(\vec{v})$

# DETERMINANT OF A MATRIX

We know matrix represents a transformation of space

Now this transformation can increase the space or decrease the space (w.r.t. original space)

✓ The factor by which the original space increases or decreases $\Rightarrow$ determinant of matrix.

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \rightarrow$$



new space defined by $\hat{\imath}$ & new $\hat{\jmath}$

original space defined by $\hat{\imath}$ 2 $\hat{\jmath}$

In 2d terms, space = area

original area = 1

new area = $3 \times 2 = 6$

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \rightarrow 3 \times 2 - 0 \times 0 = 6$$

✓ So in 2d, |determinant of a matrix| = factor by which unit area in original space changes

Now determinant can be -ve : if the new transformed space is obtained by flipping original space. Eg. $\begin{bmatrix} 3 & 0 \\ 0 & -2 \end{bmatrix}$

# ZERO DETERMINANT:

When will the area in the transformed space $= 0$



$$\begin{bmatrix} 0 & 0 \\ 4 & 4 \end{bmatrix}$$ or $$\begin{bmatrix} 1 & 1 \end{bmatrix}$$

or

$$\begin{bmatrix} 1 & 3 \\ 0 & 0 \end{bmatrix}$$

$\Rightarrow$ If the columns in the matrix are linearly dependent ∴ determinant becomes zero

$\Rightarrow$ "Singular Matrix": When the determinant is zero (i.e. inverse does not exist)

# IDENTITY MATRIX:

Recall, matrix is a transformation of space

Identity matrix is a special type of transformation that does nothing.

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{new } \hat{i} = 1(\text{old } \hat{i}) + 0(\text{old } \hat{j}) = \text{old } \hat{i}$$
$$\text{new } \hat{j} = 0(\text{old } \hat{i}) + 1(\text{old } \hat{j}) = \text{old } \hat{j}$$

a) Identity Matrix is a square matrix (# of rows = # of columns)

b) Identity Matrix has 1s on the main diagonal & 0s elsewhere

Multiplying a vector by identity matrix = leaves the vector unchanged

Multiplying a matrix by inverse of the matrix = identity matrix = equivalent to doing nothing at all

$$A^{-1} A = I$$
$$A A^{-1} = I$$
$$(A \times I = A \text{ or } I \times A = A)$$

# INVERSE OF A MATRIX:

Inverse of a matrix exists only when

1) matrix is a square matrix

2) determinant of the matrix $\neq 0$

To understand why det (matrix) ≠ 0:

det (matrix) = 0 when transformation of space represented by
matrix squishes the original space

A → represents a transformation of space

$A^{-1}$ → represents a reverse transformation of space

So essentially $A^{-1} \cdot A$ ⇒ does not change the space

$A^{-1} \cdot A$ ⇒ a reverse transformation to A

However when det (A) = 0, the transformation
squishes the space e.g. 2d space becomes
1d (line) and you cannot recover 2d space
back from line.

2×2 Matrix:

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

1) Swap the positions of a and d (a & d) put -ve in front of
   b & c.
2)
3) Divide everything by ad-bc (determinant)

Why do we need Inverse? Because with matrices, we don't
divide! (no concept of dividing by a matrix)
But we can multiply by an Inverse, which achieves the same
thing

To find matrix X:

XA = B
$XA A^{-1}$ = $BA^{-1}$
X = $BA^{-1}$

# SYSTEM OF EQUATIONS

$2x + 3y + 9z = 10$

$10x + 11y + 12z = 20$

$39x + 14y + 27z = 62$

This can be thought of as:

$x, y, z \rightarrow$ w.t.s

$2, 3, 9, etc. \rightarrow$ feature values

In matrix form:

$$\begin{bmatrix} 2 & 3 & 9 \\ 10 & 11 & 12 \\ 39 & 14 & 27 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 10 \\ 20 \\ 62 \end{bmatrix}$$

$A \qquad x \qquad = b$

To have a "unique" sol$^n$

For a system of $m$ linear eq$^n$ with $n$ unknowns:

i) $m = n$

ii) inverse of matrix $A$ should exist

$Ax = b$

$A^{-1} A x = A^{-1} b$

$I x = A^{-1} b$

$x = A^{-1} b$

This is called "analytically" solving system of linear eq$^n$

Eg. A group took a trip on a bus at $3 per child and $3.20 per adult for a total of $118.40.
They took the train back at $3.50 per child and $3.60 per adult for a total of $135.20.
How many children & how many adults?

Let $x_1$ = # of children $x_2$ = # of adults

$$3x_1 + 3.20 x_2 = 118.40$$
$$3.5 x_1 + 3.60 x_2 = 135.20$$

$$\begin{bmatrix} 3 & 3.20 \\ 3.5 & 3.60 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 118.4 \\ 135.2 \end{bmatrix}$$

$$A \; X \quad = \quad B$$
$$A^{-1} A X \quad = A^{-1} B$$
$$I X \quad = A^{-1} B$$
$$X \quad = A^{-1} B$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 & 3.20 \\ 3.5 & 3.60 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 118.4 \\ 135.2 \end{bmatrix}$$

$$= \begin{bmatrix} 16 \\ 22 \end{bmatrix}$$

Solving $x_1 = 16$
$x_2 = 22$

# ELEMENT-WISE OPERATION / HADAMARD PRODUCT

1. Operands (vector/matrix) should be of same dimension

   = Operator (corresponding) applied to elements of each operand

2. Exception : M = Matrix, V = vector
   ⟹ rows should be of same dimension

$$M = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad V = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$2\times3 \qquad\qquad 2\times1$$

$$= \begin{bmatrix} 1 & 2 & 3 \\ 8 & 10 & 12 \end{bmatrix} \quad 2\times3$$

↳ broadcast

3. \* or np.multiply()

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \end{bmatrix}$$

Now perform element-wise prod

---

## DOT-PRODUCT : ✓

$$\vec{v} \cdot \vec{u} = |v||u|\cos\theta$$

1. no. of columns of first operand = no. of rows of second operand

   $$n\times m \quad m\times x \implies n\times x$$

2. If $\vec{a} = a_x + b_y + c_z$, $\vec{b} = b_x + b_y + b_z$

   $$np.dot(a,b) = a_x b_x + a_y b_y + c_z z$$

3. If both are matrix ⟹ matrix multiplication

4. One V & One M

$$\begin{bmatrix} 3 & 0 & 2 \\ 2 & 0 & -2 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} (3\times1+0\times2+2\times3) \\ (2\times1+0\times2-2\times3) \\ (0\times1+1\times2+1\times3) \end{bmatrix} = \begin{bmatrix} 9 \\ -4 \\ 5 \end{bmatrix}$$

$$3\times3 \qquad 3\times1 \qquad\qquad\qquad 3\times1$$

# DERIVATIVE OF SIGMOID

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$\frac{d}{dx}\left(\sigma(x)\right) = \frac{d}{dx}\left((1+e^{-x})^{-1}\right)$$

$$= -(1+e^{-x})^{-2}(-e^{-x})$$

$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \frac{e^{-x}}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}}$$

$$= \frac{1+e^{-x}-1}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}}$$

$$= \left(1 - \frac{1}{1+e^{-x}}\right) \cdot \frac{1}{1+e^{-x}}$$

$$= (1 - \sigma(x))(\sigma(x))$$

tanh is rescaled sigmoid fn

$$\tanh(x) = 2\,\sigma(2x) - 1$$

---

$$\frac{d}{dx}(x^n) = nx^{n-1}$$

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(e^{-x}) = -e^{-x}$$

$$\frac{d}{dx}(ax) = a$$

$$\frac{d}{dx}(e^{ux}) = e^{ux} \cdot \frac{d}{dx}(ux)$$
$$= e^{ux} \cdot u$$

$$\frac{d}{dx}(1+x) = 1$$

# CLOSED FORM SOLN : (ANALYTICAL SOLN)

Closed form soln of equation $\rightarrow$ finite number of elementary operations $(+,-,\times,\div),\sqrt{},$ ... $\therefore$ soln to get soln at

Eg. $ax^2 + bx + c = 0$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

✓ Closed form soln are not preferred due to very expensive operations. Hence gradient descent methods are used

e.g. for linear reg. closed form soln

$$\beta = (X^T \cdot X)^{-1} \cdot X^T \cdot y$$

(inverse )expensive

bcz

$X = K$ columns ( K predictors / Attributes )

$N$ rows ( no. of observations )

Can be very big matrix itself

# KOLMOGOROV COMPLEXITY:

— length of optimal specification of an algo or object (string)

— ||||||||||| → Less kolmogorov complexity as there is pattern

— a59c3de22a → High kolmogorov complexity as there is no obvious pattern

— kol. complexity of $x$ = shortest program outputting $x$

If you fix a lang, e.g. python & write program to output above two strings then since first string has an obvious pattern, it will require shorter program to o/p it

∴ Kt describes how "compressible" a string is

# EIGEN VALUES & EIGEN VECTORS

Recall matrix is a transformation of space

$$\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$$

new î , new ĵ

Now think what happens to a vector when multiplied by this matrix meaning where will this vector land in the transformed space (w.r.t. original space)



span of the vector in original space

new vector

span of the vector in original space

If the vector in original space was on î, it will still remain on î-axis as î is now $3\hat{\imath}$ (transformation) but will be stretched by a factor of 2

Similarly the vector $\begin{bmatrix} -1 \\ 1 \end{bmatrix}$ when multiplied by $\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$

$$\Rightarrow \begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}\begin{bmatrix} -1 \\ 1 \end{bmatrix} = -1\begin{bmatrix} 3 \\ 0 \end{bmatrix} + 1\begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} -3 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \end{bmatrix} = 2\begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

will remain on its "span" but will be stretched by a factor of 2

# Eigen Vectors:

The vectors which remain on their original span even after the matrix transformation are called eigen vectors

Note: eigen vectors are generally normalized to unit length

## Eigen Values:

The factor by which eigen vectors are stretched or contracted in the new transformed space are called eigen values.



---

## Eigen Vectors = Axis of Rotation

If a vector stays in the same orientation (span) when the space is transformed, it is in effect the axis of rotation

Now if this cube rotates but the vector still points in the same dir (or opp dir but on the same span) → axis of rotation

# FINDING EIGEN VECTOR & EIGEN VALUE

## OF A MATRIX

$$A\vec{v} = \lambda \vec{v}$$

where $\vec{v}$ = eigen vector
$\lambda$ = eigen value

Matrix-vector multiplication <u>is same as</u> scaling the vector

LHS = Matrix · Vector
RHS = Scalar · vector

Writing scalar in matrix format = diagonal matrix with diagonal element

$$= \lambda$$

$$= \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$A\vec{v} = \lambda I \vec{v}$$

$$(A - \lambda I)\vec{v} = 0$$

Now we don't want $\vec{v} = 0$ since we want non-zero eigen vector

$$\therefore (A - \underbrace{\lambda I)\vec{v}}_{\text{matrix}} = 0$$
vector

Now a non-zero vector when transformed via a matrix is zero only when the matrix transformation "squishes" the vector into lower dimension

$$\Rightarrow \det(A - \lambda I) = 0$$

Note: det. of matrix = 0 when the space transformation done by matrix squished into lower dimension

get the value of $\lambda$ by solving $\det(A-\lambda I)=0$
and substitute in $(A-\lambda I)\vec{v} = 0$ to get value of $\vec{v}$

E.g. $A = \begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \\ 2 & 6 & 5 \end{bmatrix}$

$A-\lambda I = \begin{bmatrix} 3-\lambda & 1 & 4 \\ 1 & 5-\lambda & 9 \\ 2 & 6 & 5-\lambda \end{bmatrix}$

$\det(A-\lambda I) = \det\left( \begin{bmatrix} 3-\lambda & 1 & 4 \\ 1 & 5-\lambda & 9 \\ 2 & 6 & 5-\lambda \end{bmatrix} \right) = 0$

Simpler example: $A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$

$A-\lambda I = \begin{bmatrix} 2-\lambda & 2 \\ 1 & 3-\lambda \end{bmatrix}$

$\det(A-\lambda I) = 0 \Rightarrow (2-\lambda)(3-\lambda) - (2)(1) = 0$
$\therefore \lambda = 1$

$(A-\lambda I)\vec{v} = 0$

$\begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix}\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow x\begin{bmatrix} 1 \\ 1 \end{bmatrix} + y\begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$= \begin{bmatrix} x+2y \\ x+2y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

here x & y = 0 but this is just an example

# DIAGONAL MATRIX INTERPRETATION IN TERMS OF EIGEN VECTORS & EIGEN VALUES

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \qquad \begin{bmatrix} -5 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$



Original space     . Transformation

$$\begin{bmatrix} -5 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad \begin{bmatrix} -5 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$$

Eigen values — factor by which eigen vectors expand or contract in the transformed space

Eigen vectors as they stay on their span, may be expand or contract

# EIGEN DECOMPOSITION (MATRIX FACTORIZATION TECHNIQUE)

Eigen decomposition of a matrix is a type of decomposition that involves decomposing a square symmetric matrix into a set of eigenvectors and eigenvalues.

If the eigenvectors of a matrix $A$ are $v_1, v_2, \ldots v_n$ and $v_1, v_2, \ldots v_n$ are stacked column-wise to create a matrix/eigen $V$, and $\lambda_1, \lambda_2, \ldots \lambda_n$ are eigen values, stacked diagonally to create matrix $\lambda$; then

$$\boxed{A = V \, diag(\lambda) \, V^{-1}}$$

eigen decomposition of matrix $A$
($A$ needs to be square & symmetric)

Eigen decomposition is used to simplify calc. of other complex matrix operations. Also used in PCA.

The decomposition can be derived from fundamental property of eigen vectors:

$$A\vec{v} = \lambda\vec{v}$$
↳ if you pack all vectors into a matrix $V$ column-wise

$$\Rightarrow AV = \lambda V$$
$$\Rightarrow AV = V\lambda$$
$$\Rightarrow A = V\lambda V^{-1}$$

# UNITARY MATRIX:

Square matrix where conjugate transpose = Inverse

For real matrix, conjugate transpose = transpose

$$U^* = U^T = U^{-1} \text{ , if } U \text{ is real matrix}$$

Unitary matrix symbol

Since $V$ in eigen decomposition is unitary, eigen decomposition can also be written as:

$$\boxed{A = V \Lambda V^T}$$

# VECTOR SPACE:

A space where the rules of vector algebra apply.

# LINEAR REGRESSION

(No assumption abt normal dist. of Independent or dependent vars)

## ASSUMPTIONS:

1. Linear Relationship b/w X & Y
   - Use Scatter Plot Matrix to check
   - if not linear relationship, transform variables
     - e.g. right skew — log transform
       - left skew — reciprocal — then log transform

2. No multi-collinearity (if interpretation is the objective) between predictors
   (if else not needed)
   - Use correlation matrix to check multi-collinearity
   - if VIF > 10, handle multi-collinearity
     1) identify var which are co-related & remove one of them & remove them
     2) highest VIF var → remove
   (+ other methods provided in multi collinearity notes)

3. No- Auto correlation b/w residuals
   - Use Residual plot to check (residuals should be i.i.d. → random var.) → no pattern
   - Auto-correlation: residuals should not be on each other / residuals
     no pattern → homoscedastic fit
     
     Residual vs Fitted value (ŷ)

4. Homoscedasticity: - Equal var of residuals
   - Use Residual plot to check
   
   Residual vs Fitted value (ŷ)
   
   ← Heteroscedastic (cone shape)

5. Normality Assumption of Residuals:
   - ~~Use~~ Residual plot to check
   - Residual histogram → Normal

6a) Regression is sensitive to outliers : check it from residual plot

# Linear Regression Eq^n :

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n \quad ①$$

or

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots \ldots + \theta_n x_n$$

or

$$\hat{y} = \theta^T \cdot X$$
$\rightarrow$ contains $\theta_0$, bias term

✓ Loss Fn:

$$SSE = \sum (y_i - \hat{y}_i)^2$$

or more generally $MSE = \frac{1}{m} \sum_{i=1}^{m} (y_i - \hat{y})^2$

(1) Solving for $\beta_0, \beta_1 \ldots \beta_n$ using Var-Cov. method :

✓ (OLS or Line of best fit method)

$$\beta = \frac{Cov(x,y)}{Var(x)}$$

where $Cov(x,y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$

Where N = # of data values/ examples

$$Var(x) = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

After calculating $\beta$'s , substitute in eq ① to get $\beta_0$

$\rightarrow$ Cov. is a measure of how much two variables change together linearly

$\rightarrow$ If these $\beta$ no linear relationship between two variables, their covariance will be equal to zero; the variables are linearly uncorrelated but not necessarily independent.

In vector form: Normal Equation

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot y$$

(2) Gradient Descent Way of solving

$SSE = \frac{1}{2} \sum (y - \hat{y})^2$

→ for mathematical convenience as it helps in calculating gradient easily

Update rules:

$$\theta^{(next\ step)} = \theta^{(prev.)} - \text{learning rate} \times \frac{\partial SSE}{\partial \theta}$$

e.g. $\hat{y} = a + bx$

New a = a - learning rate $\times \frac{\partial SSE}{\partial a}$

New b = b - learning rate $\times \frac{\partial SSE}{\partial b}$

POLYNOMIAL REGRESSION: — Linear Model ( we assume the form of function as linear in y

$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \dots$ (eg.: )
interaction terms

→ For given predictors e.g. $x_1$ & $x_2$ how to generate polynomial features

from sklearn.preprocessing import PolynomialFeatures

poly = PolynomialFeatures ( degree = 2 ) → generate features for $x_1, x_2$ if so ))
all comb. of $x_1$ & $x_2$

$x\_poly = poly.\ fit\_transform\ (X)$

$x_1, x_2, x_1 x_2, x_1^2, x_2^2$ = all comb indefy
original $x_1, x_2$

→ Now fit linear Regression Model

# ASSESSING REGRESSION MODELS:

## 1) $R^2$:

$$R^2 = \frac{\text{Variance Explained}}{\text{Total Variance}}$$

→ Interpretation
→ Always between 0 & 1

$$R^2 = 1 - \frac{SS_{res}}{SS_{total}}$$

where $SS_{res} = \sum (y_i - \hat{y}_i)^2$

$SS_{total} = \sum (y_i - \bar{y}_i)^2$

→ Mean of observations values of response variable (not predicted)

## 2) RMSE or MSE → Prediction

$$RMSE = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n}}$$

$$RMSE = \sqrt{MSE}$$

→ easier to calc. gradients (compared to MAE)

## 3) MAE : Mean Absolute Error
(In case of outliers)
→ reduces the effect of outliers

$$= \sum |y_i - \hat{y}_i|$$

# FLOWCHART FOR REGRESSION

1) Start with Linear Regression: 1) Linear relationship b/w X & Y Scatterplot + correlation matrix
   2) No multi-collinearity → Residual plot

2) Once we get the Residual plot look for two things:
   1) Pattern - There shouldn't be [pattern, should be random]
   2) Outliers :- Few if more

   → 1) RANSAC or tree based model
   2) Evaluation metric
   - MAE (reduces the effect of outliers)

3) If there are patterns in residual plot, it means if we are unable to capture some explanatory information

   ↓

   Polynomial Regression : to account for interaction terms
   or
   Non-Linear Model (e.g. RF regressor) : to capture non-linear relationship b/w predictors

# LOGISTIC REGRESSION

## Assumptions:

1) Binary outcome

2) linear relationship between independent variables and log odds of outcome

3) Observations need to be independent

4) No or little multi collinearity between independent variables (to obtain reliable estimates)

5) Large sample size

6) No influential outliers that excessively influence the estimation of model parameters

---

Why not use linear regression for predicting probabilities

1. Response/target is binary, not continuous in LR

2. Error terms are not normally distributed in LR
   (if error terms normally distributed
   MLE = LS
   max. likelihood    least squares
   estimation)

---

ODDS, LOG ODDS, LOGIT FN & LOGISTIC (SIGMOID) FN:

ODDS = ratio of something happening to something not happening

$$odds = \frac{p}{1-p}$$

where $p$ = prob. of something happening

Probability is between 0 & 1 but when transforming prob. to odds, it removes the upper bound (upper bound = $+\infty$) but lower bound is still 0.

Taking log makes the lower bound $-\infty$    $\left(\log 1 = 0,\ \log 0 = -\infty\right)$

So $\log(\text{odds})$ range $[-\infty, +\infty] \Rightarrow$ range of values linear combination of independent variables will output

Also $\log(\text{odds}) = \log\left(\frac{p}{1-p}\right) = \text{logit}(p)$
  i.e. logit fn converts prob. to log of odds

Logistic/Sigmoid fn is the inverse of logit fn and converts $\log(\text{odds})$ into probability
i.e. squashes values from $[-\infty, +\infty]$ to $[0,1]$

Logistic/Sigmoid $(\text{prob}(t)) = \dfrac{1}{1+e^{-\log(\text{odds})}} = \dfrac{1}{1+e^{-\log\left(\frac{p}{1-p}\right)}}$

Prob (sigmoid(x))

Depending on range of x
(linear comb. of independent var)
the shape can be shifted
towards left or right

Sigmoid fn$(x) = \dfrac{1}{1+e^{-x}}$

Derivation of Above Formula:

$\log\left(\frac{p}{1-p}\right) = \log(\text{odds})$

Exponentiate both sides.

$\dfrac{p}{1-p} = e^{\log(\text{odds})}$

$\dfrac{1-p}{p} = e^{\log(\text{odds})}$

$p = (1-p)\, e^{\log(\text{odds})}$

$p = e^{\log(\text{odds})} - p\, e^{\log(\text{odds})}$

$p + p\, e^{\log(\text{odds})} = e^{\log(\text{odds})}$

$p = \dfrac{e^{\log(\text{odds})}}{1+e^{\log(\text{odds})}}$

∴ Logistic Regression eqn.

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = w_0x_0 + w_1x_1 + \ldots + w_kx_k$$
$$\text{OR}$$
$$\beta_1x_1 + \beta_2x_2 + \ldots \beta_0$$

Exponentiating

$$\frac{p}{1-p} = e^{w_0x_0 + w_1x_1 + \ldots}$$
$$= e^{w^T \cdot x}$$

$$p = \frac{e^{w^T x}}{1+e^{w^T x}}$$

$$p = \frac{1}{1+e^{-(w^T x)}}$$

↗ p ← prob. of outcome 1

## INTERPRETATION:

$$\log\frac{p}{1-p} = \beta_0 + \beta_1x_1 + \beta_2x_2 + \ldots$$

$$\therefore \frac{p}{1-p} = e^{\beta_0+\beta_1x_1+\beta_2x_2}$$

$$\frac{p}{1-p} \rightarrow \text{odds}$$

Quantitative Var: For every unit change in predictor, the odds of event (with outcome =1) change by a factor of $e^{<coeff. value>}$

If $e^{coeff \, value} > 1$ , change → increase

$e^{coeff \, value} < 1$ , change → dec.

e.g. Let's assume only one prediction $p_1$

⇒ odds of event=1 change for 1 unit change in predictor (in terms of factor)

$$= e^{\beta_1(x_1+1)}$$
$$= e^{\beta_1 x_1} \cdot e^{\beta_1}$$
$$= e^{\beta_1}$$

# Qualitative Var:

E.g. $x_2 = \begin{cases} 1 & \text{for gold} \\ 0 & \text{for silver} \end{cases}$ for dimmy coding

If $x_2 = $ gold when compared to silver $(x_2 = 0)$,
i.e $x_2 = 1$

the odds of event $(=1)$ change by the factor
of $e^{\text{coeff value}}$

when $e^{\text{coeff value}}$
when $e^{\text{coeff val}} \begin{cases} >1 & \text{, change} \rightarrow \text{increase} \\ <1 & \text{, change} \rightarrow \text{decrease} \end{cases}$

# Cost Function:  (Logloss function or negative log likelihood)

✓ We know that

$$L(\underset{\underset{\text{parameter}}{\downarrow}}{\Theta} | \underset{\underset{\text{outcome of }i^{th}\text{ instance}}{\downarrow}}{y_i}) = P(y_i | \Theta)$$

i.e. The likelihood of true parameters being a certain value given data
(i.e. outcome)

== Prob. of observing data given some true parameter values
(i.e. outcome)

MLE: maximizes LHS, but since LHS = RHS above, we maximize
$P(y_i | \Theta)$

Prob. of $i^{th}$ sample being $1$ is given by logistic fn (with some params)
$$P(x_i) = \frac{1}{1 + e^{-z_i}}$$

Mathematically, for samples labeled as '1', we try to estimate $\theta$ such that prob of probability $p(x)$ is as close to $1$ and for samples labeled as '0', we try to estimate $\theta$ such that product of all prob. is as close to 0. ie. $1-p(x)$ should be close to 1

(product of prob is taken since all samples are independent ie.
$P(A \text{ and } B) = P(A) \cdot P(B)$)

∴ For samples labelled 1 : $\prod p(x_i)$
s.t. $y_i = 1$

For samples labelled 0 : $\prod (1-p(x_i))$
s.t. $y_i = 0$

$L(\theta) = \prod\limits_{s.t.\ y_i=1} p(x_i) \quad \prod\limits_{s.t.\ y_i=0} (1-p(x_i))$

$L(\theta) = p(x_i)^{y_i} \quad (1-p(x_i))^{1-y_i}$

→ likelihood fn that needs to be maximized

4f 1 sample prob = $p(x_i)$
2 sample prob = $(p(x_i))^2$

∴ $y_i$ sample prob = $(p(y_i))^{y_i}$

Taking log
$\log L(\theta) = y_i \log p(x_i) + (1-y_i) \log (1-p(x_i))$ ⓘ

Maximize above ⓘ or minimize the $-ve$ of above eq ⓘ

$= - \left[ y_i \log p(x_i) + (1-y_i) \log (1-p(x_i)) \right]$

$y_i =$ true label
(either 0 or 1)

Average over $n$ samples
Log loss fn $= \dfrac{-1}{n} \left[ y_i \log p(x_i) + (1-y_i) \log (1-p(x_i)) \right]$
Log loss heavily penalizes classifiers that are confident about incorrect predictions.

Nodes: $y_1$, $y_2$, ... $y_n$ with weights $w_1$, $w_2$, ... $w_n$ → $M$ ($w^T x$) → Sigmoid fn → prob fn → $y$

## UNDERSTANDING MAXIMUM LIKELIHOOD ESTIMATION

Goal in logistic regression is to find the best fitting S curve for given data points
& In logistic regression, we transform the y-axis from the probabilities to $\log(\text{odds})$



Label 1 — Sigmoid(x)
Label 0
X

$\log(\text{odds})$
$+\infty$
X
$-\infty$

If $p = 1$ (relabeled)

$$\log\left(\frac{p}{1-p}\right) = \log\left(\frac{1}{1-1}\right) = \log\left(\frac{1}{0}\right)$$

$$= \log 1 - \log 0$$

$$= 0 - (-\infty)$$

$$= +\infty$$

Similarly if $p = 0$, $\log\frac{0}{1-0} = \log\frac{0}{1} = \log 0 = -\infty$

The log(odds) transformation pushes the data points to +ve & -ve infinity

Side Note: That's why we can't use least squares to find best fitting line as the residuals are also equal to +ve and -ve infinity

MLE → A likelihood fn is defined that calculates the probability of observing the outcome given input data & model (maximized.)

log(odds) of data pts is taken, now project these data pts onto the log(odds) line which gives each data pt, a log(odds) value



Then transform this log(odds) value to probabilities using

$$p = \frac{e^{log(odds)}}{1 + e^{log(odds)}}$$

After calculating prob, plot them on S-curve

Now keep rotating the log(odds) line projecting data pts onto it transforming it to probabilities calculating the log likelihood (until maximized) (e.g. using logit fn)

Note: The algorithm that find the line with max. likelihood does so in a way that it increases log likelihood each time it rotates the line. So in few rotations, we get optimal fit

Methods that maximize log likelihood:

- Newton Raphson method
- Fixed point iteration
- Bisection method
- Muller's method
- Gradient Descent (minimizes negative log likelihood)

In general, to maximize log likelihood, take derivative of likelihood fn & equate it to zero.

MLE does not tell "how" to find the optimal value of $\theta$, it just tells how one value of "$\theta$" is more likely than others (optimization algos like GD tell "how" to find this optimal value)

# DECISION TREES

- AIM: The aim of decision trees → reduce impurity
  i.e. splitting the nodes using features which lead
  to maximum improvement in impurity i.e. moving
  towards "pure" nodes
  (performing a split when building a decision tree ⇒ dividing up the feature space)

- IMPURITY MEASURES: prob. of incorrectly classifying an element, if randomly chosen
  (SPLITTING CRITERIA)  (node value from start of tree)

1) GINI SCORE: or GINI INDEX

$$G_i = 1 - \sum_{K=1}^{n} P_{i,K}^2$$

$P_{i,K}$ = ratio of class K instances to total instances at i'th node

& n = total no. of classes

- (does binary splits only)
- Higher Gini = ~~higher~~ lower purity / higher impurity

e.g. $1 - \left(\frac{0}{54}\right)^2 - \left(\frac{49}{54}\right)^2 - \left(\frac{5}{54}\right)^2 \approx 0.168$

measure of disorder → range [0,1]

2) Entropy: $H_i = - \sum_{K=1}^{n} P_{i,K} \log(P_{i,K})$   → range [0,1]

$P_{i,K}$ = ratio of class K instances at i'th node

higher entropy "lower homogeneity"

e.g. $-\frac{49}{54}\log\left(\frac{49}{54}\right) - \frac{5}{54}\log\left(\frac{5}{54}\right)$

Information gain = reduction in impurity = Entropy_parent - May_impurity

Information gain = reduction in impurity = Entropy_parent - Entropy_children
(wt.avg of all subnodes that & parent node splits into)

i.e.  $\dfrac{n_{child_i} \cdot Entropy_{child_i}}{n_{parent} \cdot Entropy_{parent}}$   weights

→ want max gain → amongst features

impurity ⌐0

∴ AIM of DT:

NOTE: Start with the node (feature) having impurity (e.g. Gini, entropy) which reduces impurity maximally
& then select nodes which reduces impurity maximally

Stopping Cond.
1) max_depth is reached
2) impurity does not decrease
3) Min num of examples required to split an internal node is reached / Min num of samples required to be at leaf node is reached

# DT Algos:

1. ID3 — Can produce decision trees with nodes that have more than 2 children
2. C4.5 ⟩→ use entropy/inf. gain
3. CART (Classification & Regression Tree)
   — greedy algorithm: Optimum split at each level without
     checking whether the split will lead to
     lowest possible impurity several levels down
   — searches for lowest

— produces only binary trees
— uses gini

## Decision Boundary of DT:

Depth=1

Depth=2

Depth=3

— Decision boundary is parallel to some axes OR perpendicular
  to other axes

## REGULARIZATION:

— If depth of decision tree is increased, DT is
  prone to overfitting

— To regularize decrease max* params or
  increase min* params

# PRUNING:

- Pruning is done to avoid overfitting if not pruned, the tree will come very large and will overfit the data

- Pruning deletes unnecessary nodes Remove the nodes whose child nodes provide no purity improvement i.e. if purity improvement is not statistically significant

$\chi^2$ test: Null hypothesis: Improvement is due to chance

Hence if p-value > 0.05, prune the node

# ESTIMATING CLASS PROBABILITIES: (INFERENCE/PREDICTION)

- Traverse the tree to find leaf node for the instance
- ratio of instances of class $k$ to total instance at this node

# ESTIMATING TARGET CLASS: Majority vote at leaf node
(PREDICTION)

# NON-PARAMETRIC:

- DT are non-parametric; as they do not make an assumption on the distribution of data
⇒ generally do not require (centering or scaling)

# NON-LINEAR:

- DT are non-linear: non-linear relationship b/w dependent & independent variables

# REGRESSION TREES: (outcome is continuous not feature)

- Minimize MSE, instead of impurity
- Prediction is simply the average target value of instance associated with the leaf node

# ADVANTAGES

1. Easy to understand & interpret
2. Standardizing the data is not required
3. "Feature Selection" inbuilt as less informative features are not part of decision tree
4. Can handle missing values

# DISADVANTAGES:

1. Prone to over-fitting
2. Sensitive to small changes in data : a slight change in data can result in a very different tree

# DIFFERENCE B/W STOPPING COND & PRUNING:

Stopping cond. = criteria for stopping the growth of tree as soon as cond. is met
pruning
  1) look few steps ahead & decide whether we want to stop
  2) make DT to a large depth giving -ve returns
  3) Suppose start at the bottom & start splitting a split is giving -ve gain (i.e loss) and next DT will stop at -10 but removing leaves which are compared to the top gives overall gain of +10 supports it
     Suppose a split of 20. A simple gain of +10 DT and keep both leaves but pruning we will see overall gain of +10 supports it
     pruning but different nodes at the top of tree - based models, the nodes at the top of

NOTE: sklearn DT does not provide max gain → this fact can be utilized.
NOTE: In tree - based models, the
tree provide max gain
to get feature importance

Consider the dataset

| Loves Popcorn | Loves Soda | Age | Loves (Cool) As Ice |
|---|---|---|---|
| F1 | F2 | F3 | Target |
| Yes | Yes | 7 | No |
| Yes | No | 12 | No |
| No | Yes | 18 | Yes |
| No | Yes | 35 | Yes |
| Yes | Yes | 38 | Yes |
| Yes | No | 50 | No |
| No | No | 83 | No |

1. Find the feature (along with its split) with lowest impurity (Gini or Entropy)
Note that when we are looking at a feature, we have to take into account it's all possible values.



Loves Popcorn
— Yes → Target F 3 T 1
— No → Target F 1 T 2

Loves Soda
— Yes → Target F 1 T 3
— No → Target F 3 T 0

Looking at the diagram, we can see "Loves Soda" (No branch) results in pure node whereas none of the branches of "Loves Popcorn" results in pure node. Hence "Loves Soda" does a better job in predicting the target ("Loves Cool As Ice")

Quantifying it via Gini Impurity:

$$\text{Gini Impurity} = 1 - (\text{prob. of "N" in target var})^2 - (\text{prob. of "N" in target var})^2$$

Gini Impurity for "yes" branch of "Loves Popcorn"

$$= 1 - (\text{prob. of "y" in target var})^2 - (\text{prob. of "N" in target var})^2$$

$$= 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2$$

$$= 0.375$$

Gini Impurity for "No" branch of "Loves Popcorn"

$$= 1 - (\text{prob of "yes" in target var})^2 - (\text{prob of "No" in target var})^2$$
$$= 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2$$
$$= 0.444$$

∴ Total Gini Impurity for "Loves Popcorn" (includes both branches)

= Weighted avg. of both branches

$$= \left(\frac{1+3}{1+3+2+1}\right) \times 0.375 + \left(\frac{2+1}{1+3+2+1}\right) \times 0.444$$
$$= 0.405$$

Similarly, Total Gini Impurity for "Loves Soda" (both branches)
$$= 0.214$$

Hence, we choose "Loves Soda" as the root node

The final tree looks like this: (upto this pt.)

[ Loves Soda ]   Yes ↙   ↘ No

"Yes" branch will now only see the dataset where "Loves Soda" value is "Yes".

Similarly "No" branch will now only see the controlled dataset where "Loves Soda" value is "No".

| Loves Popcorn | Loves Soda | Age | Loves ? |
|---|---|---|---|
| Yes | Yes | 7 | No |
| No | Yes | 12 | Yes |
| Yes | No | 18 | Yes |
| Yes | Yes | 35 | Yes |
| No | Yes | 38 | Yes |

2. Now on this controlled dataset, use the other features (e.g. Age or Loves Popcorn) to find next best feature for splitting.

NOTE: We exclude "Age" in our discussion for simplicity. In reality, we would compute impurity for "Age" as well. Since it is a numerical feature, the splits creation is different.

Creating splits for "Age" feature

i) Sort the "Age" feature

| Age | Target Loves Soda Alice |
|-----|-------------------------|
| 7   | No                      |
| 12  | No                      |
| 18  | Yes                     |
| 35  | Yes                     |
| 38  | Yes                     |
| 50  | No                      |
| 83  | No                      |

→ 9.5
→ 15
→ 26.5

ii) Take avg of consecutive "age" feature value as "age" feature split point

e.g. Age
9.5 ← Age → 79.5

iii) Calc. gini impurity of all split points & at avg. for "Age" feature

---

Let's say the next best feature to split is "Age" and we end up with a tree like below:

Loves Soda
Yes ↙        ↘ No
Age<12        Target
              N | 2
              0 | 0

Yes ↙   ↘ No
Target     Target
N | 4      N | 0
0 | 0      0 | 1

3. Since all nodes are pure, there is no need to continue building tree. (In practice, stopping condition is used to end building tree)

4. Now, we need to assign output values to each leaf:
Classification : Majority vote
Regression : Mean of values

So, the final tree is:

Loves Soda
— yes → Age<12.5
  — yes → Loves Cool As Ice = NO
  — No → Loves Cool As Ice = Yes
— No → Loves Cool As Ice = No

---

INFERENCE/PREDICTION:
If a new data point comes along with the following values:

| Loves Popcorn F1 | Loves Soda F2 | Age F3 | Target Loves Cool As Ice |
|---|---|---|---|
| Yes | Yes | 5 | ? |

Traversing the tree:

Loves Soda → Age<12.5 — yes → Loves Cool As Ice = Yes

# BAGGING & BOOSTING

**Bagging:** Sampling with replacement i.e. duplicate records/examples in diff. samples (bootstrap aggregation)

**Pasting:** Sampling without replacement i.e. no duplicate records/examples in different samples

**Boosting:**
- Combining several weak learners into strong learner
- Train models sequentially, each trying to correct the mistakes of its predecessor

Note: row sampling & column sampling facility is provided by most implementation (before each boosting iteration)

# BAGGING Vs BOOSTING: DIFFERENCE

## BAGGING
(Not RF since RF = bagging + random subset/features for each tree)

1. Parallel ensemble: Each model is built independently on diff. subset of data

2. AIM: To decrease variance, not bias

3. Trees are built independently in a parallel fashion hence easily parallelizable / distributed

## BOOSTING
(All kinds of boosting eg. (adaboost, GBM, Xgboost) at nothing specific to one boosting algo

1. Sequential ensemble: next model to be trained corrects previous model's errors

2. AIM: Primarily to decrease bias. However, variance is also reduced due to combining multiple trees

3. Although trees are built sequentially, construction of diff branches & split finding procedure can be parallelized / distributed (available in Xgboost but not in standard GBM implementation)

WHY BAGGING MODELS like RF reduce variance & not bias whereas
BOOSTING MODELS like GBDT reduce primarily bias and also variance?

## RF

RF uses "fully grown" DTs (low bias, high variance)
The trees are made uncorrelated to maximize
decrease in variance but the algo cannot
reduce bias : hence need for large, unpruned
trees so that bias is initially as low as
possible (No pruning happens in RF)

## GBDT / Xgboost

## GBDT / Xgboost

Boosting is based on "weak
learners" i.e. "shallow trees"
(high bias, low variance),
sometimes even as small as
decision stumps (trees with
2 leaves). By sequentially
training on residuals i.e.
fixing mistakes, bias is
reduced.

However, combining multiple
trees reduces variance too.

# RANDOM FOREST

RF = Bagging + "Feature selection"

— Build n parallel & independent trees by:

for each tree
- ① — selecting bootstrap sample (samples with replacement) (random subset of features for each tree)
- ② — selecting "random" subset of features → random in Random Forest

① ② → (mechanism to reduce over fit)

— Combine predictions Using
- → majority voting : classification
- → Avg. : regression

— SPLITTING CRITERIA: Same as decision trees (e.g. Gini, Entropy → classification, MSE → Regression (based on target var. not predictor))

— WHY SUBSET OF FEATURES & NOT SIMPLE BAGGING

1) If subset of features not different for each tree
   → learners (individual trees) will be highly correlated
   → "Diversity" in ensemble decreases

   → All learners will point in the same dir
      ie. not a great way to combine

2) Using subset of features helps in → over-fitting (reducing)

— WHY NO PRUNING:
— RF is robust to noise from each individual tree

- HYPERPARAMS
1) Num. of trees        (usually sufficient to tune)
2) Size of bootstrap sample
3) Num. of features for each tree

- EASILY PARALLELIZABLE:
  - Since trees are built independently in parallel

# GBDT / GBM (GRADIENT BOOSTED DECISION TREES / GRADIENT BOOSTING MACHINES)

- Sequentially adds model/tree to an ensemble, each one correcting the errors made by the predecessor model

- In particular, the new model is fit to the "residual" errors i.e. difference between true label and predictions from prev. model

- This "residual" is called "pseudo-residual" since it is the same as taking $-ve$ gradient of loss fn (e.g. MSE) w.r.t. predictions
  $$\underline{\text{"GRADIENT" in GBDT/GBM}}$$

$$L_{MSE}(Y, \hat{Y}) = \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

Now, let's take partial derivative of loss fn with respect to specific $\hat{y}_j$'s

$$\frac{\partial L_{MSE}(Y, \hat{Y})}{\partial \hat{y}_j} = \frac{\partial}{\partial \hat{y}_j} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

$$= \frac{\partial}{\partial \hat{y}_j} (y_j - \hat{y}_j)^2$$

$$= 2 (y_j - \hat{y}_j) \frac{\partial}{\partial \hat{y}_j} (y_j - \hat{y}_j)$$

$$= 2 (y_j - \hat{y}_j)(-1)$$

$$= -2 (y_j - \hat{y}_j)$$

Dropping constant $-2$
$$= (y_j - \hat{y}_j) \rightarrow \text{"residual"}$$

Note: we can now remove Summation since partial derivative of loss $f_n$ for $i \neq j$ is 0

i.e. Chasing Residual vector in GBM = Chasing (-ve) gradient of loss fn via gradient descent

(magnitude + dir.)

[-ve sign is imp since we want to move in opp dir of gradient for minimization
(in gradient descent)]

- Each new tree corresponds to another step of gradient descent

- GBM/GBDT uses Loss fn and gradient of loss fn to chase the "residual" (ie pseudo-residual)

- GBM/GBDT have both (split criteria and loss fn to build next tree)
  (splitting criteria and loss fn to build next tree)

- GBM/GBDT does not use Gini or Entropy as splitting criteria
  instead it usually is based on loss fn e.g. MSE
  (reason: trees in gbm predict gradient of loss (even for classification) which is a numerical value hence like MSE (regression))

- Loss fn like mse, mae (regression) or logloss (classification) is used to calculate the residual on which the next tree is fit
- How GBM is trained:

1) Usually the first tree takes avg. of true label as the initial prediction for all examples/instances ($P_0$)

2) Next, residual or (true label - initial prediction) is calculated ($R_0$) (by taking gradient of loss fn)

3) Next, 2nd tree is trained to predict the residuals (calc.)
   (i.e. residuals become the target) in step 2)
   (splitting criteria: usually based on loss fn)

4) To make new prediction: ($P_1$)

   $$P_1 = P_0 + \text{Learning-Rate} \times (\text{Predictions from 2nd tree})$$
   
   (initial model pred eg. avg.)

5) Again calc. residuals and repeat from step #2

| $x_1, x_2$ | $x_3$ | Label/True Value(s) | $P_0$ | $R_0$ | $P_1 = P_0 + LR \times (\text{pred. from next tree})$ | $R_1$ | ...... |
|---|---|---|---|---|---|---|---|

- PREDICTION USING GBM:

Prediction = Sum of predictions of all models weighted by learning rate

i.e. $1^{st}$ Model Pred + LR x $2^{nd}$ Model Pred.

$+ LR * 3^{rd}$ Model Pred

⇒ "Sum" since models are trained on residuals

⇒ Learning Rate: How quickly the error is corrected from one tree to next.

$0 < LR < 1$

If LR is low: more trees are needed to fit the training set but usually results in better generalization
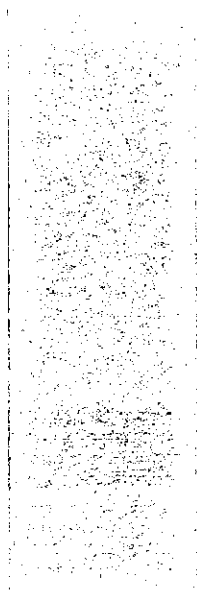
— Over-fitting problem is usually seen in GBMs

— Boosting can be considered as "Gradient Descent" in Prediction Space

Since we are optimizing the combined model predictions and not model parameters (as in neural networks)

— HYPERPARAMS:

— Num. of trees
— Max depth of tree
— Learning Rate
— Row sampling & Column sampling facility provided by most implementations beside each boosting round

# XGBoost

## Loss Function:

Xgboost, rather than explicitly fitting the pseudo-residuals (as in GBM), aims to minimize the following objective at each iteration

$$L^{(t)} = \sum_{i=1}^{n} l\left(y_i, \hat{y}_i^{(t-1)}\right) + f_t(x_i) + \Omega(f_t)$$

- label for instance i
- → prediction from ensemble of trees until prev time step $f^{(t-1)}$ for instance i
- → prediction from the new tree at time t for instance i
- → regularization term for the new tree

The above eq. is a fn of fns ie. includes fn as a parameter to another fn and cannot be optimized using traditional optimization methods in euclidean space [as mentioned in xgboost paper]

Hence, 2nd order approximation is used to optimize the above loss fn utilizing Taylor Series Expansion

Note: Taylor series Expansion is used to calc. the value of entire fn at every point if the value of the fn and all of its derivatives are known at a single point

Since $y_i$ & $\hat{y}_i^{(t-1)}$ are constants in above eq.

$$L^{(t)} = \sum_{i=1}^{n} l\left(y_i, \hat{y}_t(x_i)\right) + \Omega(f_t)$$

Now applying Taylor Series Expansion

$$= \sum_{i=1}^{n} \left[ g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t)$$

- → gradient of loss fn w.r.t prediction for instance i
- → Hessian of loss fn w.r.t prediction for instance i

# GRADIENT & HESSIAN CALCULATION:

To minimize the objective fn, XGBoost uses gradient descent.
In each iteration, the first and second order derivatives (gradient and Hessian) of the loss fn are calculated w.r.t. predicted output of each instance in the dataset, giving vectors $g$ and $h$ with gradient & hessian values for each instance

Note: Hessian provides a more precise estimate of the dir. of highest decrease of loss fn which allows the model to converge faster.

# CONSTRUCTION OF TREES:

In the $t$th iteration, the best tree $f_t$ that minimizes the objective fn is added using the calculated gradients and Hessians in a greedy fashion (or approx. greedy for large datasets)

- Start with single node (containing all instances) $\rightarrow$ residual of instances
  - i) Exact greedy - split pt. obtained from all features and values for greedy feature $\land$  (obtained from initial pred.)
  - ii) (approx. greedy / weighted quantile sketch)
  - For larger datasets, XGBoost partitions features into quantiles instead of scanning all feature values
- Iterate over i) all features and values ii) all possible split
  - evaluate loss reduction

# GAIN

$$GAIN = loss_{parent} - (loss_{left\ branch/child} + loss_{right\ branch/child})$$

[Gain must be $\gamma$ min-split-gain hyperparam (Gamma) else stop growing the branch]

Gain is fn of gradient & hessian of left and right branches + (lambda) (L2 regularization term)
$\gamma$ min. loss reduction required for a split (controls model complexity) $\frac{\lambda}{2}\frac{(G_L/H_L+\lambda)}{...}$

Gain eqn combines both loss reduction & regularization term (helping prevent over-fitting & making optimal trade-off b/w complexity & predictive power)

## ILLUSTRATION

**Step 1:** Let's say initial Prediction = avg. of target values
(same as GBM)

$$= \frac{100 + 90 + 110 + 75}{4}$$

$$= \frac{375}{4} = 293.75$$

| Feature | Target |
|---------|--------|
| X | Y |
| 10 | 100 |
| 15 | 90 |
| 20 | 110 |
| 25 | 75 |

**Step 2:** Calc. residuals for each instance
(same as GBM)

| Feature | Target | |
|---------|--------|---------|
| X | Y | Residual |
| 10 | 100 | 6.25 |
| 15 | 90 | -3.75 |
| 20 | 110 | 16.25 |
| 25 | 75 | -18.75 |

**Step 3:** All the residuals go to the root node
(Diff. than GBM)
from now on

$$\boxed{6.25, -3.75, 16.25, -18.75} \rightarrow \text{Root Node}$$

**Step 4:** Now the first cand. split pt is taken (either from greedy or approx. greedy / weighted quantile sketch method)
Let's say it is 12.5

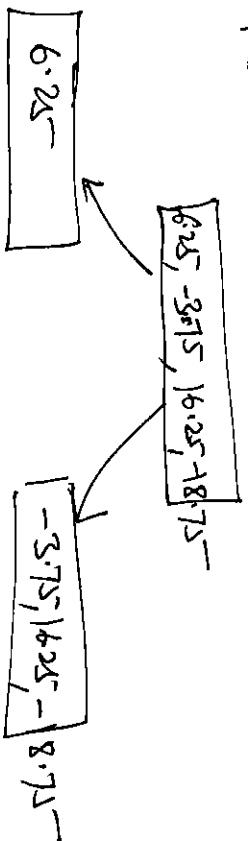then all residuals whose instance feature value is < 12.5 go in left subtree
& remaining i.e. >12.5 go in right subtree

$$\boxed{X < 12.5}$$

left subtree → 6.25

right subtree → -3.75, 16.25, -18.75

**Step 5:** Now loss is calc. for parent (with all residuals) and left & right subtree.

ie.

6.25

6.25, -3.75, 16.25, 18.75

-3.75, 16.25, -18.75

**Step 6:** loss is a fn of gradient & hessian of instances in respective nodes

& finally, gain is calculated

$$GAIN = loss_{parent} - (loss_{left\ subtree} + loss_{right\ subtree})$$

+ reg. param $\lambda$

Gain is calc. for all candidate split pts and the split pt. with max gain is chosen as the actual/true split pt.

Let's say optimum split pt is $X < 12.5$ then the tree looks like this:

$X < 12.5$

6.25

-3.75, 16.25, 18.75

Only 1 obs, so cannot split further

can split further

## HOW TO PREDICT / OBTAIN OUTPUT VALUE FROM LEAF NODES:

Let's say this is the final tree

$X < 15$

-10.5

O/P: -10.5

$X < 30$

6.75

O/P: 7

-7.5

O/P: -7.5

O/P value from leaf node = Sum of values (ie. residuals) / No. of residuals + $\lambda$ (reg. param)

($\lambda$ = 0 in calculation of O/P values for this tree)

(Leaf nodes are residuals)

# HOW IS HESSIAN CALCULATED, SINCE IT IS COMPUTATIONALLY EXPENSIVE

## TO CALCULATE 2nd Order Derivatives

Ans: Xgboost uses a diagonal approximation to the Hessian.

A diagonal 'n×n' matrix has atmost 'n' non zero elements

The diagonal approximation scales nicely bcz it grows linearly in 'n', as opposed to dense Hessian which grows quadratically

# HANDLING OF MISSING VALUES IN XGBoost → SPARSITY-AWARE SPLIT FINDING ALGO

At the time of training the instances with missing values of a feature are placed on both left and right branch of split and Gain is calculated. The side where gain is maximized, becomes the default dir. whenever this feature value is missing (e.g during prediction if the example has this feature value missing)

Split into 2 tables

(w/o missing values & only with missing values)

Illustration:

Feature →
→Target

| X | y | Residual |
|---|---|---|
| 10 | -7 | -6 |
| 20 | 8 | -8 |
| 25 | 10 | +12 |
| 30 | 12 | +16 |
| ? | 18 | -10 |
| ? | 10 | +8 |

| X | y | Residual |
|---|---|---|
| 10 | -7 | -6 |
| 20 | 8 | -8 |
| 25 | 10 | 12 |
| 30 | 12 | 19 |

| X | y | Residual |
|---|---|---|
| ? | 18 | -10 |
| ? | 10 | 8 |

At each split point for X

e.g is 22.5, 27.5

all the missing value instances are first placed on left, then right and default placed whenever gain is maximized

Note: Residuals are used in finding best split etc.

# SPLIT FINDING

— Exact Greedy

— Approx. Greedy (→ Weighted Quantile → Sparsity Aware Split Finding)

Xgboost provides a
→ makes Xgboost parallelizable/distributed
→ hyperparam to choose either: i) exact greedy
                                              ii) Approx. greedy

## EXACT GREEDY

→ Enumerates 'all' possible splits on 'all' features
→ In order to do so for a feature, all values of a feature must be stored and sorted → becomes an issue when dealing with continuous features in large datasets

→ Most single machine boosting implementations e.g. GBM in sklearn, R' gbm or single machine version of Xgboost support exact greedy

→ Though exact greedy can be parallelized (one feature on each node, if node can accommodate all values of one feature, usually not possible in large datasets), it is computationally expensive and parallelization gains are much lower than approximate greedy

## (PARALLEL LEARNING)

APPROX GREEDY: → Enumerating 'all' possible splits on 'all' features is impossible to efficiently compute when the data does not fit in memory or split needs to be found in distributed setting
→ Approx. Greedy

& (WEIGHTED QUANTILE SKETCH)
→ Cand. split points are proposed based on the quantiles of feature distribution. E.g. if a continuous feature has values sorted from 1...100 then the split points are 10, 20, ... 90.

→ If the dataset is huge, sorting and finding quantile cannot be done on a single machine. So the dataset is chopped and put on diff. machines (distributed)
→ Quantile Sketch Algorithm combines the quantiles value from each machine to make an approx. histogram, which in turn is used to calc. approx. quantiles of the full dataset

→ When every instance in the dataset has equal wt → Quantile sketch
  When instances in the dataset have unequal wt → Weighted Quantile
  Sketch

E.g. In regression, all instances have equal weight
  however, in Classification, instances could have wts in imbalanced
  datasets
  (e.g. Scale_pos_weight hyperparam)

→ The 'weight' in weighted Quantile sketch is a fn of
  hessian for the instance

# DIFFERENCE BETWEEN Xgboost VS GBM

## Unique Features of Xgboost (not found in GBM) [Also Advantages]

— Regularization term in Loss fn + Gamma : min loss needed for split

— Use of Hessian
    Xgboost → Converges faster
    → fbs longer to train

— Approximate (+ Exact greedy for similar datasets) greedy Algo (GBM uses only exact greedy algo)

— Parallel Learning (Xgboost parallelizes the construction of 'one' the tree by building several nodes at a given depth simultaneously, standard GBM does not have this parallelizing implementation)

— Weighted Quantile Sketch

— Sparsity Aware Split Finding

— Cache-aware access (Xgboost caches gradient & hessian in CPU cache for quick access)

— Blocks for out-of-core Computation
(Xgboost — if dataset is too large to fit in memory, it must be stored on hard drive but reading & writing to hard drive is super slow, so Xgboost compresses the data while storing on machine, decompressing the data while reading (reading is faster than reading the uncompressed data))

— Pruning difference: Xgboost split upto max-depth and then start pruning the tree backwards & remove split beyond which there is no -ve gain.
GBM uses greedy approach and stops as soon as -ve gain is seen at a split.

# CONS OF XGBOOST:

- Doesn't support Categorical features natively i.e. categorical features need to be encoded before feeding to xgboost algo
- Somewhat sensitive to outliers as new trees are fit to fix errors made by prev. trees (though regularization helps)

EXTREME in Xgboost: Refers to the engineering goal to push the limit of computational resources for boosted tree algo

# POPULAR HYPERPARAMS TO TUNE:

1. scale_pos_weight: Usely for imbalanced datasets should be set to $\frac{\# majority\ class\ instances}{\# minority\ class\ instances}$

2. Eta (shrinkage or Learning Rate)
3. L1 (alpha), L2 (lambda)
4. Gamma: Min. loss reduction to create a split
5. subsample: Fraction of data to train on
6. colsample: Fraction of features to train on

## K-MEANS.

Step1: Pick centroids

Step2: Calc. dist of every pt to each centroid and assign pts to nearest centroid (cluster)

Step3: Now for every cluster, calculate the mean of all data points in that cluster ---> New centroid

Repeat step2: Calc. dist of every pt (irrespective of current cluster) to the new centroid and re-assign pts to nearest centroid (cluster)

Repeat step3: Now for every cluster, calc. the mean of all data points in that cluster → New centroid

— Continue until:

1) Threshold for the difference b/w cost fn value b/w iterations is reached

2) Threshold for the difference b/w positions of centroid b/w subsequent iterations is reached.

COST FN OF K-MEANS: $\sum_{i=1}^{K} |x_i - \mu_k|^2 \rightarrow$ within cluster SSE

$\mu_k \rightarrow$ centroid of the cluster where pt $i$ belongs

ELBOW METHOD: value of cost fn

HOW TO FIND K of K-MEANS:

— The elbow method plots different cost fn produced by different K.

- Intuitively : If K increases, distortion or within cluster SSE will decrease.

Cost fn value



k=optimum value

HYPERPARAMETER TO TUNE : K

Advantages
+ Scalable
$O(n K i)$
  n: no. of different pts
  K: no. of clusters
  i: no. of iterations

Disadv
- K- beforehand
- Non deterministic, depending on initial choice of centroids
- Centroids can be non-existent sample

# AFFINITY PROPAGATION

EXEMPLARS: Data points representative of clusters

SIMILARITY FN: — Euclidean Distance (in general, can be other fn too)

Pairwise Similarities between data points are taken as input and clusters are found by maximizing the Similarity between data points and Exemplars.

## IMPORTANT HYPERPARAMETERS:

Preference: The similarity value b/w Exemplars and other data points in the cluster

· Low preference = cluster size will be large but fewer clusters

· High preference = cluster size will be smaller (compact) but large no. of clusters

Advantages:

+ No K beforehand

+ Deterministic

+ Exemplars are real data points

DisAdvantages:

- Not scalable

$O(n^2 i)$

$n$ = no. of data points
$i$ = iteration

# DBSCAN

1) Core point : if Min. no. of neighbour points (MinPts) fall within a specific radius $\varepsilon$

2) Border point : falls within radius $\varepsilon$ of a core point but has fewer neighbours than MinPts

3) Noise point : Except core & border point

Notion of Density : no. of points within radius $\varepsilon$

- Form a separate cluster for each core point or connected grp of core points
- Assign border point to the cluster of its corresponding core point

HYPERPARAMETERS TO TUNE

1) MinPts   2) $\varepsilon$ - radius

Adv:
+ Arbitrary shaped clusters
+ No K beforehand
+ Notion of Noise - assigned
-1 - useful in anomaly detection

DisAdv:
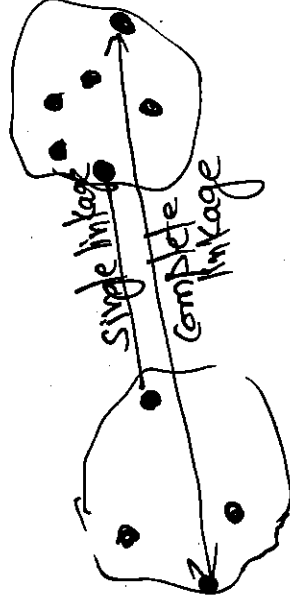- Hyperparameters $\varepsilon$ & MinPts to tune
- Not entirely deterministic, pts assignment of border pts can be to either neighbouring clusters on order of data being processed/update
- Not scalable

# HIERARCHICAL CLUSTERING

Agglomerative     Divisive

Every pt is a cluster & then merge.
    All pts belong to one cluster and then break

## CRITERIA FOR MERGING:

1) Single Linkage - minimum
2) Complete Linkage - maximum
3) Average Linkage - mean/mean
4) Ward's Linkage: Those Clusters are merged which leads to minimum increase of total within-cluster SSE



single linkage / complete linkage

1) Compute dist matrix of all data pt
2) Merge according to criteria (min, max, avg, ward)
3) Update dist. matrix
4) Repeat 2 & 3 until one cluster remains.

## HYPERPARAMETER:    LINKAGE

- In sklearn implementation of Agglomerative clustering
   K can be pre-specified

## Adv

+ No K beforehand
+ Deterministic
+ Dendrograms to Understand Structure in data

## DisAdv

- Members/Data pts when clustered, do not unmerge
- Not scalable
  Time complexity of Agglomerative with single linkage

✓ Worst case: $O(n^3)$          Best case: $O(n^2 \log$
  i) dist. from every pt to every other pt needs to be calc. $(n^2)$
  ii) search thru this distance matrix to find closest pt $(n$ or $\log$
     linear search  binary search

## GMM:

— Probabilistic Clustering

— Deterministic

— K - No. of Gaussian Mixture (Gaussian Distr.)
need to be pre-specified.

---

## ASSESSING THE QUALITY OF CLUSTERING:

1) Visualization, if possible (if no. of dim > 3, difficult)

2) Silhouette plots & scores

3) If labels present → all classification metrics

---

## SILHOUETTE PLOTS:

Silhouette coeff is calculated for every data point

$$s = \frac{b-a}{\max\{b, a\}} = [-1, 1]$$

$a$ = cluster cohesion, avg dist b/w a data point and all other data points in same cluster

$b$ = cluster separation, avg dist b/w a data point and all other data points in nearest cluster

if $s = 0$   if cluster separation and cohesion are equal

silhouette score : avg. silhouette coeff of all data point in a cluster
(for a cluster)

Cluster

1

2

3

0.1 0.2 0.3 0.4 0.5 0.6 0.7 → Silhouette coeff

Silhouette score

outliers

Shape or thickness of individual plot

Size = size of cluster

# KNN

1. Choose the number K and a distance metric
2. Find K nearest neighbors of the sample that we want to classify.
3. Assign the class label by majority vote. (avg in case of regression)

In case of a tie (if K = even) the answer is implementation specific. In sklearn, the algorithm will prefer the neighbours with closer distance to the sample.

## CHOICE OF K :

1) Odd (so that majority can be calculated) in case of 2 class problem

2) $K = \sqrt{n}$   where $n$ = no. of samples

3) If K is very high → under-fitting
   K is small → over-fitting
   Perform cross-validation

## WEIGHTED KNN

Wts could be assigned to K nearest neighbors to make the final prediction.

## COMPUTATIONAL COMPLEXITY OF KNN (test) prediction time

If $d$ = no. of dimensions
   $O(dn)$

However, efficient data structures like KD-tree can reduce time complexity at the cost of increased training time

# LAZY LEARNER & INSTANCE BASED LEARNING

Models based on instance based learning are characterized by memorizing the training dataset and lazy learning is a special case of instance based learning that is associated with ref(zero) cost during the learning process.

## NON-PARAMETRIC:

- Does not make explicit assumption about the form of target fn.
- It cannot be characterized by a fixed set of params and the number of params grow with training data

## CURSE OF DIMENSIONALITY:

- Prone to over-fitting when no. of features increase compared to size of training set. This is because even closest neighbors are too far away in a high-dimensional space to give a good estimate

Pros:

1. Immediately adapts to the new training data
2. Somewhat insensitive to outliers
3. No assumption about data

Cons:

1. Computational complexity is high for classifying new samples
2. Storage req are high as all data is needed
   Prediction time.
3. Sensitive to scale in data i.e standarization must be done before applying it.

# DISTANCE / SIMILARITY METRIC

1. Euclidean (real-valued)

   L2 norm: $d(x,y) = \sqrt{\sum (x_k - y_k)^2}$

2. Manhattan (real valued)

   L1 norm: $d(x,y) = \sum |x_k - y_k|$

3. Cosine Similarity (real valued)

   $\cos\theta = \dfrac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|}$   $[-1, 1]$

   $= 1$   similar
   $= 0$   dissimilar
   $= -1$  in opp. dir

   judgement of orientation rather than mag.

   $[1,2]$ $[100, 200]$
   → cosine similar
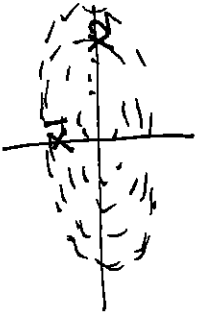   high but
   euclidean distance
   → far off points

4. Pearson Correlation (real valued)

   $r = \dfrac{\sum\limits_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$   (trend)

   $[-1, 1]$

5. Mahalanobis distance (real-valued)

- Useful if features/attributes are correlated

- Accounts for the fact that variances in each dir is different



An euclidean space $x_1$ is near to origin
but variance in $x_1$ dir is less to origin

∴ $x_2$ is near if mahalanobis distance is considered
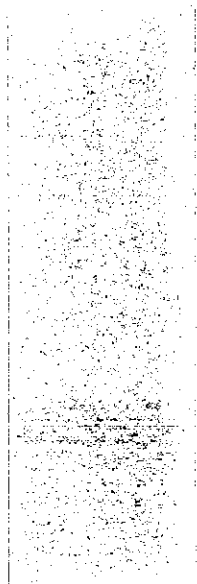since normalized by variance in that direction

(categorical)

6. Jaccard similarity

$$\frac{A \cap B}{A \cup B}$$

NOTE: If observation contains both real & cat. values

- real dist. metric for real valued + dist. metric for categorical

∴ compare observations by real & cat. dist metric separately.

## NAIVE BAYES:

$P(A \text{ and } B) = P(A) \times P(B)$ if events $A$ & $B$ are independent

Conditional probability:

$$P(A|B) = \frac{P(A \text{ and } B)}{P(A \text{ and } B)}$$

if events $A$ & $B$ are not independent

### Bayes Theorem:

$$P(O|E) = \frac{P(E|O) \; P(O)}{P(E)}$$

Likelihood ↗

$P(O)$ ← Class prior probability

$P(E)$ ← Predictor prior probability

Posterior prob.

It is easier to calculate conditional probability in one direction; often we know freq. of some evidence, given known outcome, we use it to calculate the reverse

E.g. Prob. of Disease) given Test-positive

$=$ Prob (Test is +ve | Disease) $\times$ Prob (Disease)
(scaled by) Prob (Testing +ve, with or without disease)

### NAIVE BAYES:
In reality, there are multiple evidence and in naive bayes we calculate class probabilities of each outcome & select the outcome with higher prob.

$$P(C|X) = \frac{P(X|C) \; P(C)}{P(X)}$$

$P(X) \rightarrow$ need not calc. since common to all classes. → "intuition behind → Naive Bayes"

predictors/attributes/features assumed to be independent,

If all predictors are : $P(X|C) = P(x_1|C) P(x_2|C) \times P(x_n|C)$

$P(C|X) = P(X|C)P(X|C) \times P(C)$

(add-1 variate) $\underbrace{}_{\text{if o then everything}}$ =0 (zero freq problem) Intuition behind multiplying by prob setting) is to give higher prob to more common outcome → base rate

**Ex1:** Swine flu is suspected to affect 1 in 10,000 people. Test is 99% accurate with FP=1%, FN=0. What is the probab. that you have swine flu you test +ve. What is the probab. that you have swine flu

$P(D)$ = prob. of swine flu    $P(T)$ = prob. of +ve test

$\therefore P(D|T) = \dfrac{P(T|D)P(D)}{P(T)}$

which can be rewritten as:

$$P(D|T) = \dfrac{P(T|D)\,P(D)}{P(T|D)P(D) + P(T|ND)\,P(ND)}$$

Where $P(ND)$ = prob. of not having swine flu

$P(D) = \dfrac{1}{10000} = .0001 \rightarrow$ prior

$P(T|D) = 1$

$P(T|ND) = .01$ (1% chance of false +ve)

$P(ND) = 1 - P(D) = .9999$

$$P(D|T) = \dfrac{1 \times .0001}{1 \times 0.0001 + 0.01 \times 0.9999} \approx 0.01$$

Pred $\begin{array}{c|cc} & T & F \\\hline T & 1 & \\ F & 0 & 1 \end{array}$ Actual

$\therefore$ even though your test is +ve your chance of having swine flu is 1% (compared to population relative to population)

**Ex2:**

| Weather | Play |
|---------|------|
| Sunny | No |
| Overcast | No |
| Rainy | No |
| Sunny | Yes |
| Sunny | Yes |
| Overcast | Yes |
| Rainy | No |

Prob. of Play=Yes if weather is sunny

$$P(Yes|Sunny) = \frac{P(Sunny|Yes) \cdot P(Yes)}{P(Sunny)}$$

$$= \frac{(2/3) \times 3/7}{3/7} = \frac{2}{3} = 0.667$$

$$= 66.67\%.$$

Adv. of NB: → Fast

Disadv of NB → a) Independence of features may not be true

b) Zero frequency problem:

$$P(c_i|x_i) = P(x_1|c_i) \times P(x_2|c_i) \ldots$$

if = 0 then Overall = 0

→ MAP is the foundation for Naïve Bayes

→ In case of two classes, we are assuming our
data is drawn from two distributions.

ie. we have a bunch of data where we know
the class, and want to be able to predict
$P(\text{class} \mid \text{data point})$

→ $P(\text{spam} \mid w) \propto \prod_i P(w_i \mid \text{spam}) \, P(\text{spam})$

$P(\neg\text{spam} \mid w) \propto \prod_i P(w_i \mid \neg\text{spam}) \, P(\neg\text{spam})$

Whichever is bigger of two, wins
[Maximizing posteriors]

→ Since $P(C \mid x_i) = P(x_1 \mid C) \times P(x_2 \mid C) \times P(x_3 \mid C) \times \dots \times P(x_i \mid C)$

Can lead to numerical underflow
as $[0,1] * [0,1] * \dots = $ very small value

Take log:
$= \log P(C) + \sum \log P(x_i \mid C)$

$\Rightarrow \log P(C)$

→ If variables are continuous = Use Gaussian NB
which calculates prob. using a continuous function
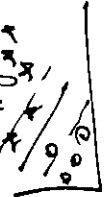
# SVM

Objective Fn: $\dfrac{1}{2}||w||^2 + C\sum\xi$  → also called as margin constant

penalty on misclassification

→ Large C → higher penalty on misclassification
∴ the left & right margins shrink
= low bias, high variance

Small C → lower penalty on misclassification, ∴ the left & right margins equal
= High bias, low variance

maximize the margin



support vectors (closest to separating hyperplane)

separating hyperplane

margin

**Hard Margin classifier**: no point are allowed to cross the margin

**Soft Margin classifier**: some point are allowed to cross the margin

A Linear classifier is of the form (essentially an eq. of a line)

$$f(x) = w^T x + b$$

dot product of two vectors : wt vector and x

$w^T x + b \geq 1$ → one class
$w^T x + b \leq -1$ → other class  } Hard margin

$w^T x + b > 0 \rightarrow$ class 1
$w^T x + b < 0 \rightarrow$ class 2  } Soft margin

only if data is linearly separable without misclassification

The decision boundary defined by a hyperplane
is said to be linear because it is linear in input feature space.

Linear SVM works well if decision boundary is linear in input space

What if decision boundary is non-linear in input space? — Kernel trick

Transform the training data into higher dimensional space. Via a mapping function φ and train a 'linear' SVM model to classify data in this new feature space. Then how unseen data can be transformed using the same mapping function φ to classify using linear SVM.

✗ Computationally very expensive to classify soon → use kernel trick

Note: The hyperplane always remains linear in the input feature dimensional space. However, in the higher space, this hyperplane takes the form of a curve

Kernel trick: A method to map features to new producing feature space, computing dot product essentially producing a scalar) without explicitly doing it thru a
⇒ achieve this without ↗ series of operations

A kernel is a function that returns the same value as the dot product of its corresponding feature vectors in the original feature space (mapped given the classifier RBF (Gaussian) kernel)

In the new feature space, the classifier becomes

E.g. polynomial kernel, kernel can be interpreted as ⇒ similarity b/w a pair of samples
$v_m$similar=1, 0 very dissimilar such that margins are maximized & misclassification is low

↓

Quadratic Programming Problem
Solved using SMO (sequential Minimal Optimization)
sub gradient descent

$f(x) = \omega^T \phi(x) + b$

## Adv of SVM:

1. ✓ Can classify data point when decision boundary is non-linear

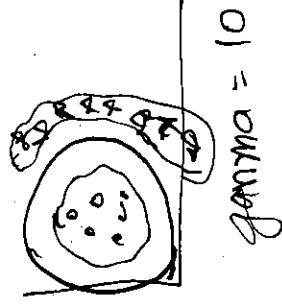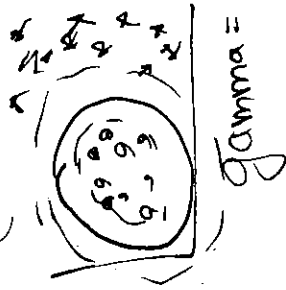2. Less overfitting as obj fn has regularization parameters

## Disadv: of SVM:

1) Can be computationally very expensive - both memory wise & training time wise

2) Choosing the right kernel & tuning hyperparameters (gamma) can be tricky

3) Not very robust to outliers

4) Fundamentally a binary classifier, in multi class setting - One vs Rest paradigm.

## PRACTICAL CONSIDERATIONS WHEN APPLYING SVM:

1. Choosing C

2. Choosing Kernel



Linear    Polynomial    Gaussian

A smaller gamma in higher dimensions = pointed bump & vice-versa for larger gamma
↑ low bias & high variance ↓ & vice-versa for larger gamma

3. Choosing Gamma (in case of Gaussian kernel)



gamma = 0.1          gamma = 1          gamma = 10

Note: If multiple of the above has to be chosen : Grid Search

Kernel Trick:

Kernel methods represent the data only thru a set of "pairwise similarity" comparisons between the original data observations $X$ (with the original co-ordinates in the lower dimensional space), instead of explicitly applying the transformation $\phi(x)$ and representing the data by these transformed co-ordinates q in the higher dimensional feature space.

In kernel methods, the dataset $X$ is represented by $n \times n$ kernel matrix of pairwise similarity comparisons where the entries $(i,j)$ are defined by kernel function $K[x_i, x_j]$

↗ ↖ Two observations in the dataset

The kernel function acts as a modified dot product.

The ultimate benefit of the kernel trick is that the obj. function we are optimizing to fit the higher dimensional decision boundary only includes the dot product of transformed feature vectors... we just substitute these dot product terms with the kernel fn and we don't even use $\phi(x)$

→ Let $w$ be the minimizer of the SVM problem for some dataset with $m$ examples $\{x_i, y_i\}$

Then for $i = 1 \cdots m$ there exits $\alpha_i \geq 0$ such that optimum $w$ can be written as $w = \sum_{i=1}^{m} \alpha_i x_i y_i$

Prediction : $\text{sign}(w^T x)$

&

Wt. matrix: $w = \sum_{i=1}^{n} \alpha_i y_i x_i$

$\therefore w^T x = \sum_{i=1}^{n} \alpha_i y_i \boxed{x_i^T x} \Rightarrow$ we only need to compute dot products b/w training examples and new example

This is true even if we map examples to high dimensional space

$w^T \phi(x) = \sum_{i} \alpha_i y_i \underbrace{\phi(x_i)^T \phi(x)}_{\downarrow}$

This is provided by Kernel trick

# HINGE LOSS (incorporates dist. from classification boundary into cost calc)

(soft margin SVM) *accounts for classification errors

$$H(y) = \max(0, 1 - t \cdot y)$$

where $t$ = true label
$(1 \text{ or } -1)$

$y$ = pred. prob (output of classifier)

**Ex1:** outcome = 1, pred = 0.5

$$H(y) = \max(0, 1 - 1 \times 0.5) = 0.5$$

**Ex2:** outcome = -1, pred = 0.5

$$H(y) = \max(0, 1 - (-1) \times 0.5) = 1.5 \quad \text{(Higher compared to above)}$$

# COST FUNCTION OF SVMs:

- Maximize margin (inv. prop. to margin)
- Minimize $w$ which is a vector orthogonal to the data-separating hyperplane onto which we project our data pts
- Primal form of Hard Margin SVM

For soft margin SVM, we combine this minimization obj with a loss fn such as hinge loss

# DIMENSIONALITY REDUCTION

→ Dimensions in ML = Features

Reducing dimension = reducing feature space

→ Why dimensionality reduction?

- If the feature space is very large, the distance b/w data pts increases hence models may not generalize well

- Computation time for algos increase in large feature space hence to reduce the computation time

→ Main Approaches for dimensionality Reduction

├ Projection (projecting onto lower d-dimensions)

├ Manifold Learning (learning the lower dimension manifold on which training instances lie e.g. a 2D manifold is a 2D shape that can be twisted & bent in higher dimensional space)

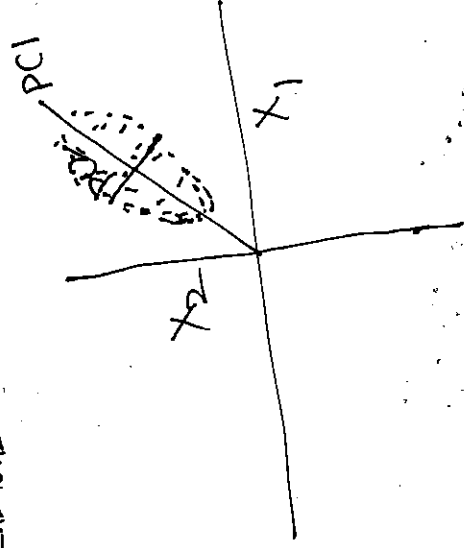→ All dimensionality red^n techniques suffer from "reconstruction error"

NOTE:

HYPERPLANE
(In context of
classification)

1) n-1 dimensions
2) Flat
3) Divides the space (n dimension feature space)
into two

# Principal Component Analysis (linear, unsupervised)
## (PCA)

IDEA: The basic idea behind PCA is to rotate the axis of dataset towards directions that maximize the variance along the new axis

Ex: in 2-dimensions:
Scatter plot



PC1 is the axis = max. variance

PC2 is the axis = second max variance

Note:-1) All these axes are orthogonal (linearly independent)

2) $PC_i$ = linear combination of features

e.g. $PC_i = a x_1 + b x_2$

  real numbers

You take top $m$ PCs where $m$ < no. of features And do a dot product with Original dataset to get the projection of data on these $m$ axes

# STEPS TO CALCULATE PCA

1) Standardize the data

— meaning subtract the mean of a dimension/feature
from each of the value in that dimension/feature
and divide the result with std. dev of the
dimension/feature

i.e. $\dfrac{x_i - \bar{x}}{\sigma_x}$    for each feature value $x_i$

       where $\bar{x}$ = mean of feature

             $\sigma_x$ = std dev of feature

2) Calculate Covariance Matrix of all features

$$Cov \left( \underset{\underset{features}{\nearrow \quad \nwarrow}}{x_1, x_2} \right) = \frac{\sum\limits_{i=1}^{n} (x_{1i} - \bar{x_1}) \cdot (x_{2i} - \bar{x_2})}{n-1}$$

Where $\bar{x_1}$ = mean value of feature $x_1$

       $\bar{x_2}$ = mean value of feature $x_2$

       $x_{1i}$ = $i$th feature value of feature $x_1$

       $x_{2i}$ = $i$th feature value of feature $x_2$

       $n$ = no. of data points / examples

Sign of covariance = +ve : both features → together
                   —ve : one feature ↑ another feature ↓
                     0 : no relation (independent features)

$$
C = \begin{bmatrix} Cov(x_1, x_1) & Cov(x_1, x_2) & \dots & \dots & Cov(x_1, x_n) \\ \vdots & & & & \vdots \\ Cov(x_n, x_1) & \dots & \dots & & Cov(x_n, x_n) \end{bmatrix}
$$

↑
Covariance matrix

IMP: For a matrix X, co-var. matrix in code can be computed as
$C = np.dot(X.T, X)/n-1$

IMP: The reason we create covariance matrix is bcz eigen decomposition requires square & symmetric matrix & co-var. matrix has this property

→ Square matrix & symmetric along diagonal

Note $Cov(x_1, x_1) = $ variance of $x_1$

3) Calculate Eigen vectors and Eigen values

Recall Eigen vectors of a matrix are the vectors that stay on their span and eigen values are the factors by which eigen vectors stretch or squish after the matrix transformation. Also eigen vectors = axis of rotation

Finding Eigen vectors of covariance matrix amounts to finding axes that remain unchanged during the matrix transformation of space according to co-variance of features

Now eigen values "stretch or squish" the eigen vectors during this transformation based on co-variance of features,

Eigen values = amt. of variance corresponding to axes

**Selecting eigen vectors based on their eigen values**

= finding axes based on amt. of variance

Recall eigenvectors and eigen values are found using eqn

$$C\vec{v} = \lambda \vec{v}$$

or In eigen decomposition form:

$$C = V \, diag(\lambda) \, V^{-1}$$

We take top-d eigen vectors based on eigen values and stack them columnwise to build eigen vector matrix

4) Project I/P data on the eigen vector matrix

Input $\cdot$ top-d eigen vector matrix = projection of input on the hyperplane defined by d-axis

$\uparrow$ dot product

Note: eigen vector = principal component of covariane matrix (usually both are unit size)

eigen vector with highest eigen value = 1$^{st}$ PC
eigen vector with 2$^{nd}$ highest eigen value = 2$^{nd}$ PC

# USING SVD TO CALCULATE EIGEN VECTORS

— Recall eigen decomposition of a matrix is only possible when the matrix is square & symmetric

$$A = X \Lambda X^T$$

↗ eigen vectors
stacked (column-wise)
needs to be
symmetric (& square)

→ diagonal matrix
where the diagonal elements = eigen values

— However, observations/examples for a ML task may not be square & symmetric

I/P
$$\begin{bmatrix} 2 & 1 & 4 \\ 3 & 9 & 12 \\ 4 & 17 & 19 \end{bmatrix}$$

wts
$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$$

labels
$$= \begin{bmatrix} 2 \\ 5 \\ 13 \end{bmatrix}$$

→ SVD is a way to factorize *any* matrix (even if not symmetric +square)

$$A = U \Sigma V^T$$
$$m \times n$$
$$m \times m \quad m \times n \quad n \times n$$

SVD decomposition $A = U \Sigma V^T$ ①

Now since eigen decomposition is only defined for symmetric, square matrices let's convert $A$ into symmetric square matrix so that it can be compared to eigen decomposition

$A^T A$ yields symmetric, square matrix for any matrix $A$

So,

$A^T A = (U \Sigma V^T)^T \cdot U \Sigma V^T$

$\qquad = (V^T)^T \Sigma^T U^T \cdot U \Sigma V^T$

$\qquad = V \Sigma^T U^T U \Sigma V^T$

$\qquad = V \Sigma^T \Sigma V^T$

$\qquad = V \Sigma^2 V^T \rightarrow$ This resembles eigen decomposition

$\qquad = \underset{X}{\underbrace{V}} \underset{\lambda}{\underbrace{\Sigma^2}} \underset{X^T}{\underbrace{V^T}}$

$A^T A$ ← Covariance matrix of $A$

$(AB)^T = B^T A^T$

$(A^T)^T = A$

$U^T U = I$
since $U$ is orthonormal (orthogonal & unit vectors)

$\Sigma^T \cdot \Sigma = \Sigma^2$
since $\Sigma$ is diagonal matrix

Hence $V =$ Matrix of Eigen Vectors as columns

$\Sigma =$ Diagonal matrix of $\sqrt{eigen\ values}$ (or singular values)

Hence from SVD decomposition① using $V$ & $\Sigma$ we can

- - form SVD decomposition① using $V$ & $\Sigma$ we can
- - Compute eigen vectors and eigen values

# WHY SVD METHOD IS USED TO CALCULATE EIGENVECTORS & EIGEN VALUES IN PCA?

1. Calc Covariance matrix is computationally expensive and can result in rounding errors. SVD way is numerically stable. Another way to say it is — singular values are more stable than eigen values

2. SVD can work on sparse matrices

# SVD (SINGULAR VALUE DECOMPOSITION)

→ Matrix Factorization technique

## SINGULAR VECTORS & SINGULAR VALUES :

For a given matrix A :

$A^T \cdot A$    $\underline{\underline{\text{results in}}}$    Symmetric, square matrix
and
$A \cdot A^T$      and both matrices have same eigenvalues

positive ($\geq 0$)

Also   $A^T \cdot A = A \cdot A^T$   (dot product is commutative
ie $a \cdot b = b \cdot a$ )
$\Rightarrow |a||b| \cos\theta = |b||a| \cos\theta$

These are infact,
covariance matrix
of A

## Fundamental PROPERTY of Symmetric Matrices :

For symmetric matrices, we can choose its eigenvectors to be orthonormal

Two vectors are said to be orthonormal, if

i) they are of unit length

ii) perpendicular to each other (orthogonal)

let's name   eigenvectors of $A A^T = u_i$  ⎫
                              ⎬ Singular vectors
        eigenvectors of $A^T A = v_i$  ⎭ of A

Both matrices $A A^T$ & $A^T A$ have same +ve eigenvalues
(or zero)

$\sqrt{\text{eigen values}}$ = Singular Values

# SVD DECOMPOSITION:

*Any* matrix $A$ can be decomposed as

$$A = U \Sigma V^T$$

where $U =$ orthonormal eigen vectors of $AA^T$

$V =$ orthonormal eigen vectors of $A^TA$

$\Sigma =$ diagonal matrix with $\sqrt{\text{eigen values}}$ in decreasing order

Dimension-wise:

$$A_{m \times n} = U_{m \times m} \ \Sigma_{m \times n} \ V^T_{n \times n}$$

$$= \left[ \begin{array}{ccc} \uparrow & & \uparrow \\ u_1 & \cdots & u_m \\ \downarrow & & \downarrow \end{array} \right] \left[ \begin{array}{cccc} a & & & \\ & a_2 & & \\ & & \ddots & \\ & & & a_n \\ & & & 0 \end{array} \right] \left[ \begin{array}{ccc} \uparrow & & \uparrow \\ v_1 & \cdots & v_n \\ \downarrow & & \downarrow \end{array} \right]$$

eigen vectors of $AA^T$

$a_1 \geq a_2 \geq \cdots a_n$

eigen vectors of $A^TA$

Since $a_1 \geq a_2 \geq \cdots \geq a_n$, if we want to ignore some very low $\sqrt{\text{eigen values}}$, we can take first $\gamma$ $\gamma$ values. In that case

$$A_{m \times n} = U_{m \times \gamma} \ \Sigma_{\gamma \times \gamma} \ V^T_{\gamma \times n} \Rightarrow \text{dimensionality reduction}$$

Note: Orthogonal matrices are useful as their inverse = transpose which means calc. their inverse is computationally cheap & stable.
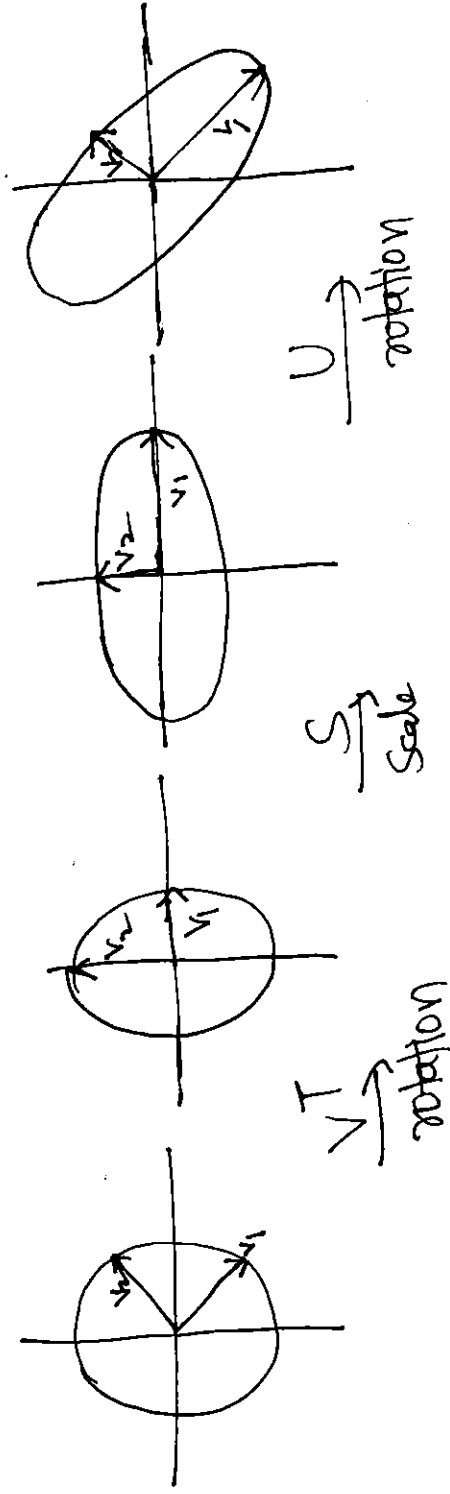
# GEOMETRICAL INTERPRETATION OF SVD

$$A = U \Sigma V^T$$

Since $A$ is a matrix = Transformation of space

This transformation can be expressed as three transformations in sequence: rotation $(V^T)$, scaling $(\Sigma)$ and again rotation $(U)$

Since $U$ & $V$ are orthonormal matrix and orthonormal matrix just change the coordinate axes via same rotation but no scaling

$\Sigma$ is a diagonal matrix and hence just scales the dimensions

# SVD in Recommendation Systems

$$A_{m \times n} = U_{m \times r} \; \Sigma_{r \times r} \; \left( V \right)^t_{n \times r}$$

$A_{m \times n}$ ↓ user-item rating matrix

$U$ → (columns are unit vectors & orthogonal)

$\Sigma$ → diagonal matrix, diagonal entries = singular values

$V$ → Columns are unit vectors & orthogonal

then

U → user to concept

Σ → strength of concept

V → item to concept

Ques) Given a data pt, how can we find which concept it relates to?

$$d = \begin{bmatrix} 0 & 4.5 & 0 & 0 \end{bmatrix}_{1 \times 4}$$

$\underbrace{V}_{n \times r}$ projecting into concept space

$$= \begin{bmatrix} 5.2 & 0.4 \end{bmatrix}_{1 \times r}$$

Higher affinity to this concept →

Note: SVD is defined for sparse matrices but you get missing values