

## 7 EVALUATION METRICS

$$\text{ACCURACY} = \frac{\text{TRUE +VE} + \text{TRUE -VE}}{\text{ALL}}$$

 MEDICAL MODEL

$$\frac{P_{\text{TRUE}}}{P_{\text{TRUE}} + P_{\text{FALSE}}}$$

- FALSE +VE OK
- FALSE -VE NOT OK

“FIND ALL THE SICK PEOPLE”

HIGH RECALL

 SPAM DETECTOR

$$\frac{P_{\text{TRUE}}}{P_{\text{TRUE}} + P_{\text{FALSE}}}$$

“FALSE +VE NOT OK  
FALSE -VE OK”

“MARKED SPAM?  
BETTER BE SPAM!?”

HIGH PRECISION

F1 SCORE - COMBINES PRECISION (P) AND RECALL (R)

$$F_1 = \frac{2PR}{P+R}$$

HARMONIC MEAN

$F_B$  - WEIGHTS PRECISION OR RECALL MORE.

$\beta < 1$  WEIGHS PRECISION MORE

$\beta > 1$  WEIGHS RECALL MORE

$$F_\beta = \frac{(1+\beta^2)PR}{\beta^2P+R}$$

## 8 MODEL SELECTION

### • TYPES OF ERRORS

GODZILLA w/ FLYSWATTER  
OVERSIMPLIFYING

UNDERFITTING

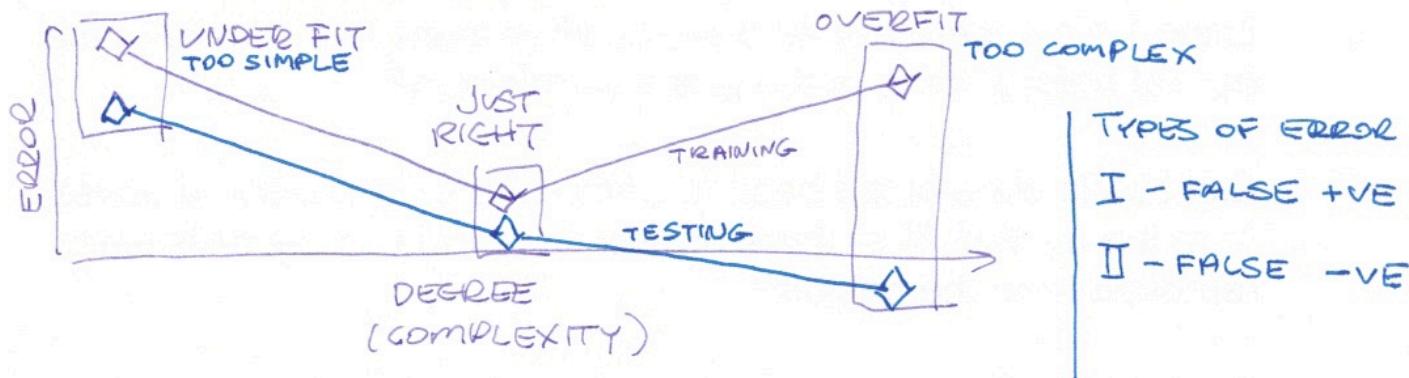
DOES NOT DO WELL IN TRAINING  
ERROR DUE TO BIAS

FLY w/ BAZOOKA  
OVER COMPLICATING

OVERFITTING

GOOD IN TRAINING, BAD TESTING  
ERROR DUE TO VARIANCE

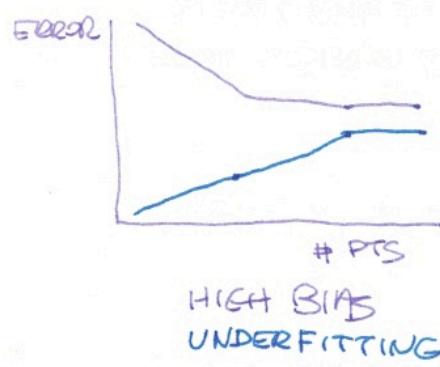
### • MODEL COMPLEXITY GRAPH



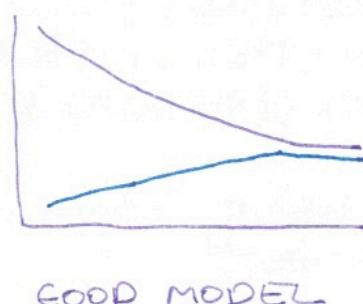
### • CROSS VALIDATION - SET APART DATA FOR FINDING MODEL PARAMS

• K-FOLD CROSS VALIDATION - TRAIN K TIMES w/ DIFFERENT  $1/K$  PORTION OF DATA FOR TESTING. HELPS PREVENT OVERFITTING  
`sklearn.model_selection.KFold`

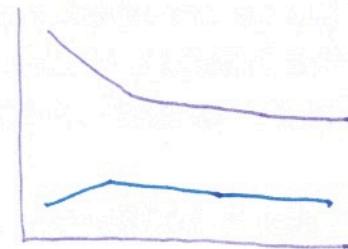
### • LEARNING CURVES



HIGH BIAS  
UNDERFITTING



GOOD MODEL



HIGH VARIANCE  
OVERFITTING

• GRID SEARCH - USED TO FIND BEST HYPERPARAMETERS  
`SKLEARN: sklearn.model_selection.GridSearchCV`

# 11 LINEAR REGRESSION $y = w_1 x + w_2$

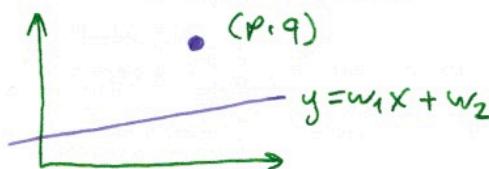
## • SUPERVISED MACHINE LEARNING

CLASSIFICATION  
YES/NO

REGRESSION  
How much?

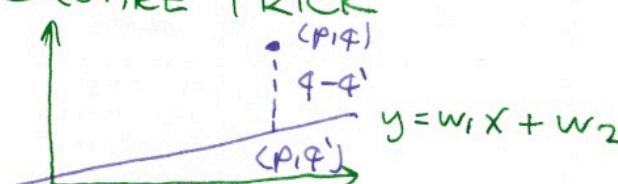
- LEARNING RATE:  $\alpha < 1$  AVOIDS TAKING TOO-LARGE STEPS

## • ABSOLUTE TRICK



$$y = (w_1 + \alpha p)x + (w_2 + \alpha q)$$

## • SQUARE TRICK

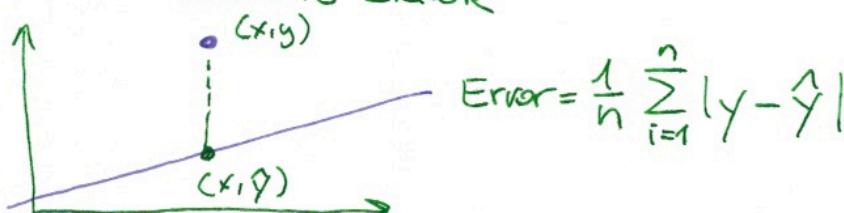


$$y = (w_1 + \alpha p(q - q'))x + (w_2 + \alpha(q - q'))$$

- GRADIENT DESCENT  $w_i \rightarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Error}$

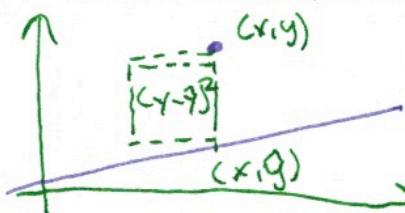
## • ERROR FUNCTIONS

### MEAN ABSOLUTE ERROR



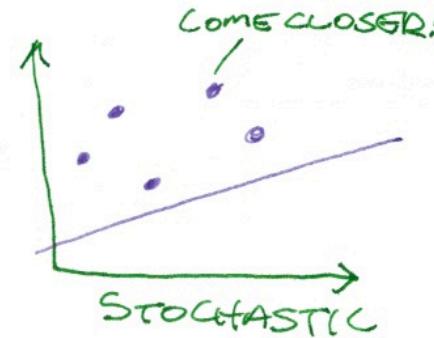
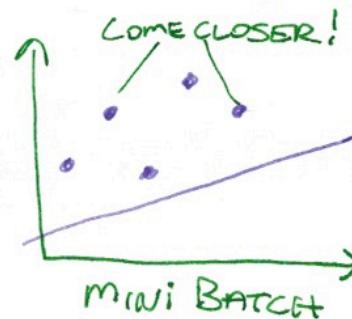
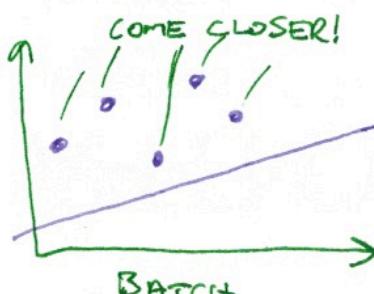
$$\text{Error} = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

### MEAN SQUARED ERROR



$$\text{Error}_{\text{MSE}} = \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2$$

## • TYPES OF GRADIENT DESCENT



# 11 LINEAR REGRESSION (CONTINUED)

## • LINEAR REGRESSION IN SCIKIT LEARN

```
from sklearn.linear_model import LinearRegression  
model = LinearRegression()  
model.fit(x-values, y-values)
```

## • HIGHER DIMENSIONS

n DIMENSIONAL SPACE → PREDICTION IS n-1 DIMENSIONAL HYPERPLANE

$$\hat{y} = w_1x_1 + w_2x_2 + \dots + w_{n-1}x_{n-1} + w_n$$

## • LINEAR REGRESSION HAS A CLOSED-FORM SOLUTION

FOR WEIGHTS  $w$ , VARIABLES  $X$  AND LABELS  $y$

$$w = (X^T X)^{-1} X^T y$$

EXPENSIVE IN PRACTICE - STILL USE GRADIENT DESCENT

## • LINEAR REGRESSION WARNINGS

1. WORKS BEST WHEN DATA IS LINEAR.
2. SENSITIVE TO OUTLIERS.

## • POLYNOMIAL REGRESSION

### • NON-LINEAR DATA

### • SAME GRADIENT DESCENT METHOD

## • REGULARIZATION

### • ADD TERMS TO ERROR TO PENALIZE COMPLEXITY

### • L<sub>1</sub> - ADD ABSOLUTE VALUE OF FIT COEFFICIENTS

### • L<sub>2</sub> - ADD SQUARES OF FIT COEFFICIENTS

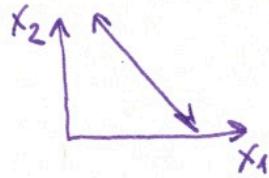
### • ADJUST $\lambda$ TO ADJUST PUNISHMENT FOR COMPLEXITY LARGE $\lambda$ - PUNISH COMPLEXITY

### • WHICH

L <sub>1</sub>	L <sub>2</sub>
COMPUTATIONALLY INEFFICIENT SPARSE OUTPUTS FEATURE SELECTION	EFFICIENT NON-SPARSE OUTPUTS

## 12 PERCEPTRONS (CLASSIFICATION)

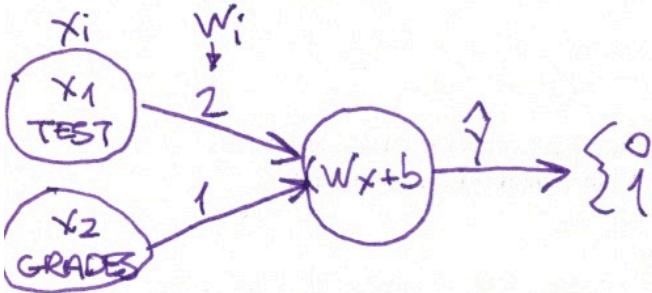
### • CLASSIFICATION



LINE BOUNDARY  $w_1x_1 + w_2x_2 = 0 \quad Wx + b = 0$   
 PREDICTION  $\hat{y} = \begin{cases} 1 & \text{if } Wx + b > 0 \\ 0 & \text{if } Wx + b < 0 \end{cases}$

$n$ -DIMENSIONS: BOUNDARY IS  $n-1$  DIMENSIONAL HYPERPLANE.

### • PERCEPTRON



### • PERCEPTRON TRICK

- GOAL: SPLIT DATA CLASSES w/ ONE LINE
- POINTS ON THE WRONG SIDE WANT THE LINE TO "COME CLOSER".

### • TRICK

- EX:  $3x_1 + 4x_2 - 10$  BAD POINT AT  $(4, 5)$

$$\frac{3 \quad 4 \quad -10}{(4 \quad 5 \quad 1)} \times 0.1 \quad \text{LEARNING RATE}$$

- - IF MISCLASSIFIED IN +VE AREA
- + IF MISCLASSIFIED IN -VE AREA

### • PERCEPTRON ALGORITHM

1. START w/ RANDOM WTS  $w_1, \dots, w_n, b$
2. FOR EVERY MISCLASSIFIED PT  $(x_1, \dots, x_n)$ 
  - $x_i = w_i = w_i + \alpha x_i \quad \forall x \in 1, \dots, n \quad \left\{ \begin{array}{l} \text{PREDICTION} = 0 \\ \text{TUE PT IN -VE AREA} \end{array} \right.$
  - $b = b + \alpha$
  - PREDICTION = 1 - VE PT IN +VE AREA  
 $w_i = w_i - \alpha x_i$   
 $b = b - \alpha$

## 13 DECISION TREES

- INTRO: CLASSIFY BY ASKING QUESTIONS OF THE DATA  
NEED TO KNOW WHAT QUESTION (DATA) TELLS US  
THE MOST

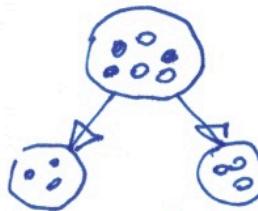
### ENTROPY

- $m$  RED BALLS,  $n$  BLUE BALLS

$$\text{Entropy} = -\frac{m}{m+n} \log_2 \left( \frac{m}{m+n} \right) - \frac{n}{m+n} \log_2 \frac{n}{m+n}$$

$$\text{LET } P_1 = \frac{m}{m+n}, P_2 = \frac{n}{m+n} \Rightarrow S = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$
$$S = -\sum_{i=1}^2 P_i \log_2 P_i$$

### INFORMATION GAIN



$$\text{Information Gain} = S_{\text{parent}} - \left( \frac{m}{m+n} S_{\text{child}_1} + \frac{n}{m+n} S_{\text{child}_2} \right)$$

BUILD DECISION TREES BY MAXIMIZING INFORMATION GAIN.

### RANDOM FOREST

- DECISION TREES PRONE TO OVERFITTING
- BUILD MULTIPLE TREES FROM RANDOMLY SELECTED FEATURES AND HAVE THEM VOTE

### HYPERPARAMETERS

- MAX DEPTH
- MIN SAMPLES PER LEAF

- MIN SAMPLES PER SPLIT
- MAX NUMBER OF FEATURES

- DECISION TREES IN SKLEARN: DECISIONTREECLASSIFIER

## 14 NAIVE BAYES

PROBABILISTIC ALGORITHM BASED ON CONDITIONAL PROBABILITY

EASY TO IMPLEMENT, FAST

### • BAYES THEOREM

- GUESS THE PERSON

• w/o INFO  $P = 0.5$  FOR EACH (PRIOR)

• WE SEE A RED SWEATER  $P_{1,rs} = 0.4$   $P_{2,rs} = 0.6$  (POSTERIOR)

- KNOWN  $\rightarrow$  INFERRRED

$$P(A) \Rightarrow P(A|R)$$

$$P(R|A)$$

- ADD MORE INFO SCHEDULE:  $P_{1,in} = 0.75$   $P_{2,in} = 0.25$

$$P_{1,rs} = 0.4$$

$P_{1,in}$	$A$	$R$	$0.75 \cdot 0.4$
	$\bar{A}$	$R$	$0.75 \cdot 0.6$
$P_{2,in}$	$B$	$R$	$0.25 \cdot 0.6$
	$\bar{B}$	$R$	$0.25 \cdot 0.4$

PERSON IN RED SO  $\bar{R}$  NOT POSSIBLE

$$P(A|R) = \frac{0.75 \cdot 0.4}{0.75 \cdot 0.4 + 0.25 \cdot 0.6} = 0.67$$

- MORE FORMALLY  $P(R|A) = P(A)P(R|A)$

$$\begin{array}{ll} P(R|A) & P(R \cap A) \\ \xrightarrow{\text{EVENT}} \begin{cases} P(A) \\ P(B) \end{cases} & \begin{cases} \xrightarrow{A} R \\ \xrightarrow{B} R^c \end{cases} \quad \begin{cases} P(R \cap A) \\ P(R^c \cap A) \end{cases} \\ P(R|B) & P(R \cap B) = P(B)P(R|B) \\ \xrightarrow{\text{EVENT}} \begin{cases} P(A) \\ P(B) \end{cases} & \begin{cases} \xrightarrow{B} R \\ \xrightarrow{B} R^c \end{cases} \quad \begin{cases} P(R \cap B) \\ P(R^c \cap B) \end{cases} \end{array}$$

$$\Rightarrow P(A|R) = \frac{P(A)P(R|A)}{P(A)P(R|A) + P(B)P(R|B)}$$

- BAYESIAN LEARNING - SPAM EXAMPLE

- NAIVE - ASSUME EVENTS ARE INDEPENDENT

$$P(\text{spam} | \text{'easy'}, \text{'money'}) \propto P(\text{'easy'} | \text{spam}) P(\text{'money'} | \text{spam}) P(\text{spam})$$

$$P(\text{ham} | \text{'easy'}, \text{'money'}) \propto P(\text{'easy'} | \text{ham}) P(\text{'money'} | \text{ham}) P(\text{ham})$$

COMPUTE THE TWO ABOVE AND NORMALIZE TO ONE.

## USING GridSearchCV

- 1 DEFINE DICTIONARY W/ PARAMETERS TO TEST:  
params = {"param\_name": [val1, val2, ...], ...}  
\* BE CAREFUL WHEN USING FRACTIONS IN PARAMETER GRID
- 2 CREATE MODEL/CLASSIFIER  
model = DecisionTreeClassifier()  
USE FIXED random\_state FOR REPEATABILITY
- 3 CREATE A SCORER. USE make\_scorer TO CONVERT A METRIC TO A SCORER.  
scorer = make\_scorer(accuracy\_score)
- 4 SET UP THE GRID SEARCH.  
clf = GridSearchCV(model, param\_grid, cv=4, n\_jobs=3,  
scoring=scorer)
- 5 RUN THE GRID SEARCH.  
clf.fit(X\_train, y\_train)
- 6 PREDICT AND CALCULATE ACCURACY.  
 $y\text{-train-pred} = \text{clf. predict}(X\text{-train})$  } USES BEST FIT.  
 $y\text{-test-pred} = \text{clf. predict}(X\text{-test})$  }  
train\_accuracy = accuracy\_score(y\_train, y\_train\_pred)  
test\_accuracy = accuracy\_score(y\_test, y\_test\_pred)

# BAYESIAN INFERENCE TUTORIAL

- DATA SET: SMS SPAM COLLECTION
  - SPAM/HAM, MESSAGE
  - USE pd.Series.map to map "spam"  $\mapsto 1$ , "ham"  $\mapsto 0$
- BASE OF WORDS
  - sklearn.feature\_extraction.text.CountVectorizer
  - MANUALLY
    1. str.lower()
    2. str.translate(str.maketrans("", "", string.punctuation))
    3. str.split()
    4. collections.Counter  $\rightarrow$  dict of word occurrences
  - Counter(s)  $\mapsto$  dict of word frequencies
- sklearn CountVectorizer
  1. cv = CountVectorizer() LOWER CASE, STRIP PUNCT, NO STOPS
  2. cv.fit(docs)
  3. cv.transform(docs).toarray()  $\mapsto$  CONVERT FREQS TO DOC MATRIX
- BAYES THEOREM IMPLEMENTATION

$P(D) = 0.01$  PROB OF DIABETES

$P(+)$  = PROB +VE TEST RESULT

$P(-)$  = " -VE " "

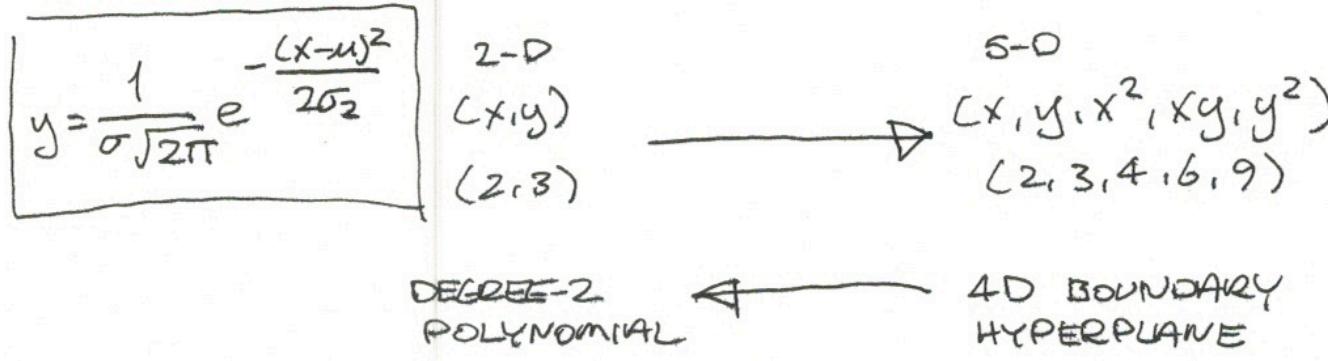
$P(+|D) = 0.9$  SENSITIVITY = TRUE +VE RATE

$P(-|\sim D) = 0.9$  SPECIFICITY = TRUE -VE RATE
- EVALUATION
  - ACCURACY = CORRECT PREDICTIONS / TOTAL PREDICTIONS
  - PRECISION = WHAT PORTION OF MSGS MARKED SPAM ARE SPAM  
 $P = \text{TRUE +VE} / (\text{TRUE +VE} + \text{FALSE +VE})$
  - RECALL (SENSITIVITY) = WHAT PORTION OF SPAM WAS MARKED AS SPAM  
 $R = \text{TRUE +VE} / (\text{TRUE +VE} + \text{FALSE NEGATIVE})$

## 15 SUPPORT VECTOR MACHINES (SVMs)

- TRIES TO MAXIMIZE DISTANCE BETWEEN BOUNDARY LINE AND DATA POINTS
- ERROR = CLASSIFICATION ERROR + MARGIN ERROR
  - MISCLASSIFICATIONS PUNISHED & DISTANCE FROM DIVIDING LINE
  - CLASSIFICATION ERROR:  $\sum |y_i - \hat{y}_i|$  → ERROR TERM  
SUM DISTANCES OF POINT FROM MARGIN
  - MARGIN ERROR:  $\text{ERR} \propto 1/\text{SIZE MARGIN} \propto \frac{2}{\|w\|} = \text{MARGIN}$   
SO ERROR =  $\|w\|^2$
- ERROR FUNCTION = MARGIN ERROR + CLASSIFICATION ERROR
- C PARAMETER - COEFFICIENT FOR CLASSIFICATION ERROR
  - SMALL C → LARGE MARGIN, POSSIBLE CLASSIFICATION ERRORS
  - LARGE C → SMALL MARGIN, CLASSIFIES PTS WELL
  - HYPERPARAMETER
- POLYNOMIAL KERNEL
  - A LINE MAY NOT BE ABLE TO DIVIDE CLASSES
  - KERNEL TRICK - NONLINEAR FNS OR ADDITIONAL DIMENSIONS

KERNEL TRICK



DEGREE OF KERNEL IS A HYPERPARAMETER

- RADIAL BASIS FUNCTION (RBF) KERNELS

- USE LINEAR COMBINATIONS OF RBFs TO SEPARATE CLASSES
- HYPERPARAMETER  $\gamma$  DETERMINES WIDTH OF RBF

$$\gamma = \frac{1}{2\sigma^2}$$

linear, poly, rbf

- SVMs in sklearn

• `sklearn.svm.SVC()`

• HYPERPARAMETERS: C, KERNEL, DEGREE, GAMMA

MAX DEGREE FOR POLY

RBF  $w^{1074}$

# 16 ENSEMBLE METHODS - USING DIFFERENT MODELS TOGETHER

## • METHODS

- BAGGING (BOOTSTRAP AGGREGATING)
- BOOSTING

COMBINE WEAK LEARNERS  
 ↓  
 STRONG LEARNER

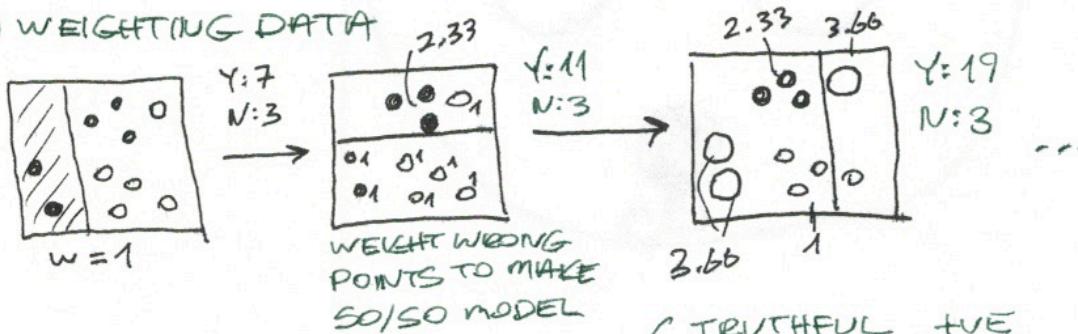
## • BAGGING

- WEAK LEARNER - ONLY SLIGHTLY BETTER THAN RANDOM
- TRAIN MANY WEAK LEARNERS ON DIFFERENT SUBSETS OF DATA
- COMBINE WEAK LEARNERS BY VOTING/AVERAGING

## • ADABOOST

1. TRAIN WEAK LEARNER TO SOME MINIMUM ACCURACY
2. TRAIN 2ND LEARNER - PENALIZING POINTS MISCLASSIFIED BY FIRST LEARNER MORE
3. CONTINUE w/ MORE LEARNERS
4. COMBINE MODELS

## • WEIGHTING DATA



{ TRUTHFUL +VE  
 RANDOM (SO/SO) O  
 LIAR -VE

## • WEIGHTING THE MODELS

$$\text{WEIGHT} = \ln \left( \frac{\text{accuracy}}{1 - \text{accuracy}} \right) = \ln \left( \frac{\text{correct}}{\text{wrong}} \right)$$

For ACCURACY = 1 or 0 - JUST USE THOSE MODELS

## • COMBINING THE MODELS

CALCULATE WEIGHTS USING WEIGHTED DATA RESULTS

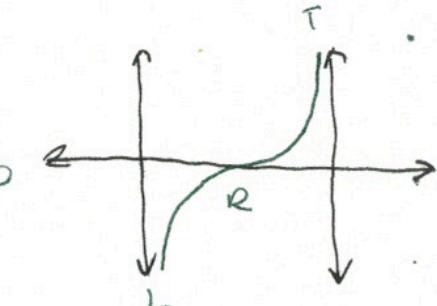
## • ADABOOST IN sklearn

- `sklearn.ensemble.AdaBoostClassifier()`

### • HYPERPARAMETERS

- `base_estimator`: WEAK LEARNER MODEL

- `n_estimators`: MAX NUMBER OF WEAK LEARNERS



# PROJECT: FINDING DONORS FOR CHARITY ML

## • OVERVIEW

- EXPLORE DATA
- TRANSFORM DATA TO WORKABLE FORMAT
- EVALUATE SEVERAL LEARNERS
- OPTIMIZE ONE OF THE LEARNERS
- TEST

• GOAL: MAXIMIZE DONATIONS, MINIMIZE NUMBER OF LETTERS SENT.

## • SUBMISSION

- COMPARE ALL SECTIONS AGAINST RUBRIC
- SUBMIT IPYNB NOTEBOOK AND EXPORTED report.html
- SUBMIT IN Finding-donors FOLDER IN GITHUB REPO.

## • MEASURING PERFORMANCE

$$\text{ACCURACY} = \frac{TP + TN}{N_{\text{PRED}}}$$

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

$$\text{RECALL } R = \frac{TP}{TP + FN}$$

$$F_B = (1 + \beta^2) \frac{P \cdot R}{\beta^2 P + R}$$

## [19] CLUSTERING

- Unsupervised learning - no labels but can still group things
  - clustering
  - dimensionality reduction
- K-means: how many clusters?
  - Steps: 1. Assign 2. Optimize
    1. assign points to the nearest class center
    2. Optimize cluster center to minimize sum of square of distances to points in the cluster
      - ▷ Initial positions of centroids is important!
- sklearn.cluster.KMeans
  - parameters
    - n-clusters (8) usually need to set this
    - max\_iter (300)
    - n\_init (10) number of times to run w/ different centroid seeds

## 21 HIERARCHICAL AND DENSITY-BASED CLUSTERING

### R-means considerations

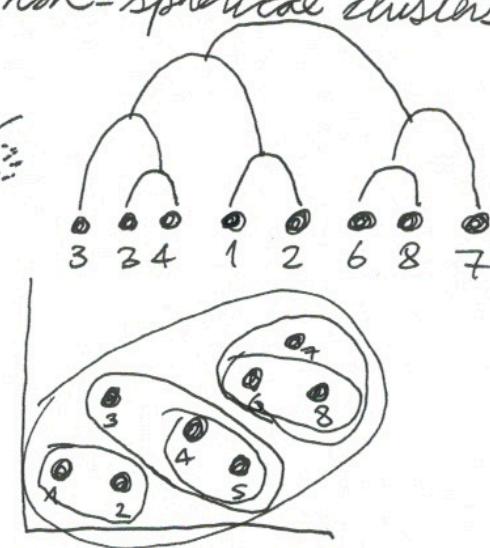
- need to know number of clusters
- clusters are hyperspherical - can't do non-spherical clusters
- data sets

uniform



### Overview of other clustering methods

- Hierarchical clustering
- Density clustering (DB Scan)
- Hierarchical clustering: single-link
- distance: smallest distance between any two points in two clusters
- handles two arcs, two rings, distinct clusters as well or better than R-means



### Hierarchical clustering: three other methods (agglomerative)

- complete link: farthest distance between any two points in two clusters
- average link: average distance between all points in two groups

• Ward's method  $A(A_1, B) = C_1^2 + C_2^2 + C_3^2 + C_4^2 - A_1^2 - A_2^2 - B_1^2 + B_2^2$



### Hierarchical clustering: implementation

#### • sklearn.cluster.AgglomerativeClustering

`clust = AgglomerativeClustering(n_clusters=3, linkage = "ward")  
labels = clust.fit_predict(X)`

#### • Dendograms - use `scipy.cluster.hierarchy`

#### • Advantages

- Resulting hierarchical representation is very informative
- Visualization
- Good at modeling real hierarchical relations

#### • Disadvantage

- Sensitive to noise, outliers
- expensive

- Density-based clustering: DBSCAN

1. Select a pt

2. Look for pts w/in  $\epsilon$

- no pts  $\rightarrow$  noise

- less than MinPts  $\rightarrow$  noise

- $\geq$  MinPts  $\rightarrow$  cluster w/ core, border pts

- Implementation

```
from sklearn.cluster import DBSCAN
```

```
db = DBSCAN(eps=0.5, min_samples=5)
```

```
db.fit(X)
```

```
db.labels_ # clusters, -1 = noise
```

- DBSCAN examples and applications

- Advantages

- Don't need to specify number of clusters

- Flexibility in shapes, sizes of clusters

- Handles noise and outliers

- Disadvantages

- Border pts reachable from 2 clusters

- Takes difficulty w/ clusters of varying densities

## 22 Gaussian Mixture Models

- soft clustering - all points in all clusters - membership
- cluster analysis
- cluster validation
- somon Clustering

- based on Gaussian (normal) dist
- parameters (1D):  $\mu \equiv \text{mean}$   $\sigma \equiv \text{std dev.}$
- fit pt density w/ Gaussians

w/in	
10	68%
20	95%
30	99%

- Expectation-Maximization for Gaussian mixtures

1. Init  $K$  Gaussian dists (can use k-means)
2. soft-cluster data - « expectation » \*
3. Re-estimate Gaussians - « maximization » \*\*
4. Eval log-likelihood to check for convergence \*\*\*

$$* E[Z_{iA}] = \frac{N(X_i | \mu_A, \sigma_A^2)}{N(X_i | \mu_A, \sigma_A^2) + N(X_i | \mu_B, \sigma_B^2)}$$

$$** \mu_A = \frac{\sum_{i=1}^N E[Z_{ij}] X_i}{\sum_{i=1}^N E[Z_{ij}]}$$

$$\sigma_A^2 = \frac{\sum_{i=1}^N E[Z_{iA}] (X_i - \mu_A^{\text{new}})(X_i - \mu_A^{\text{new}})^T}{\sum_{i=1}^N E[Z_{iA}]}$$

$$*** \ln p(X | \mu, \sigma^2) = \sum_{i=1}^N \ln \left( \sum_{k=1}^K \pi_k N(X_i | \mu_k, \sigma_k^2) \right)$$

- Implementation

gmm = sklearn.mixture.GaussianMixture(n\_components=3)  
gmm.fit(X)

clustering = gmm.predict(X)

↳ array of cluster memberships

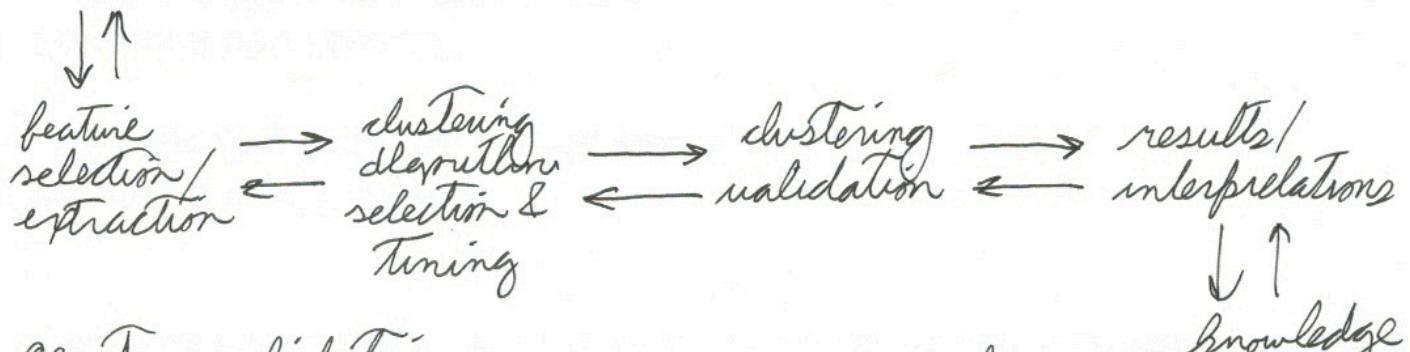
- Examples / Advantages

- soft clustering
- cluster shape flexibility } pros
- sensitive to init values } cons
- local optima
- slow convergence

22

## Gaussian Mixture Models

- Cluster analysis process
- data



- Cluster validation

1. External indices (labels)
2. internal indices (no labels)
3. relative indices

compactness

separability

- External indices: adjusted rand index

- Internal indices: silhouette coefficient [-1, 1]

$$S_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

$a_i = \text{avg dist in cluster}$   
 $b_i = \text{avg distance closest cluster}$

maximize  $S_i$  to find  $K$ .

never use this to measure DBSCAN

## 23 Feature Scaling

- sometimes must rescale so all features have equal weight
  - min-max scaler
- $$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad 0 \leq x \leq 1$$
- pro: fixed interval
  - con: outliers ruin scaling
- Implementation
    - `sklearn.preprocessing.MinMaxScaler`

- 24 Principal Component Analysis (PCA)
- axes shifts and rotations
  - measurable vs latent  
square footage      size  
no. of rooms
  - school ranking      neighborhood
  - neighborhood safety
  - principal component - composite feature that more directly probes the underlying phenomenon  
Used for: dimensionality reduction, unsupervised learning
  - How to determine the principal component?  
principal component of a data set is the direction that has the maximum variance
  - Maximal variance and information loss
  - Review / definition of PCA
    - transform input features to principal components
    - use principal components as new features
    - PCs are directions in data that maximize variance (minimize info loss)
    - more variance along PC  $\rightarrow$  higher ranked PC
    - max. no. of PCs: no. of input features
  - Implementation
 

```
pca = sklearn.decomposition.PCA(n_components=2)
pca.fit(data)
```

pca.explained\_variance\_ratio\_  
pca.components\_ [
  - When to use PCA
    - latent features
    - dimensionality reduction
    - reduce noise
    - preprocessing for other algorithms
-

- 26 Random Projection and ICA - Dimensionality Reduction
- Random Projection
    - good for high-dimensional problems
    - computationally efficient
    - $X_{k \times N}^{RP} = R_{k \times d} X_{d \times N}$   $d$  dims to  $k$  dims
    - Johnson-Lindenstrauss lemma - can project to lower dim space largely preserving distance between pts
    - Implementation
 

```
from sklearn.random_projection import SparseRandomProjection
rp = SparseRandomProjection()
new_X = rp.fit_transform(X)
```
  - Independent Component Analysis
    - Extract independent sources from a mixed sample  
↳ blind source separation >>
    - FastICA
      1. center, whiten  $X$
      2. choose initial random weight matrix  $W_1, W_2, \dots, W_n$
      3. estimate  $W_i$  containing vectors
      4. decorrelate  $W$
      5. repeat from 3 until convergence

▷ Need as many observations as signals we are trying to separate.
    - Implementation
 

```
from sklearn.decomposition import FastICA
X = list(zip(signal_1, signal_2, signal_3))
ica = FastICA(n_components=3)
components = ica.fit_transform(X)
    ↳ independent components.
```