

Lecture 1 1/7/202

Lecture 1

Data Prep 3: Preprocessing

Missing Data - several methods

Smoothing noisy data

smoothing out fluctuations

1. Binning

partition using equidepth

a. use the mean

b. use the boundaries

2. Clustering

should the outliers well

3. Create regression model

that function you create

is fed data and you replace
data points w/ results

Data Integration

from multiple data sources

Redundant - look at correlation

between attributes to see if
could be eliminated / combined

Data Transformation:

Normalization techniques

1. Min-Max

$$v' = \left[(v - \text{min}) / (\text{max} - \text{min}) \right] \times (\text{newmax} - \text{newmin}) + \text{newmin}$$

ex. \$30,000 b/w [10,000.. 45,000] to [0..1]
[30000 - 10000 / 35000] = 0.514

2. Z-score

$$v' = (v - \bar{x}) / SD$$

3. Decimal scaling

$$v' = v / 10^i$$

Good when one attribute dominates magnitude of another (like age vs. salary \rightarrow normalizing both)

Discretization

3 types of attributes:

1. nominal - categorical

2. ordinal

3. numeric

Descretization reduces # of values
for a continuous attribute
usually done by dividing range
into intervals

Can be used to turn it into category
Can also be used to do data reduction
Some algorithms require categorical
data only

1. Binning - see slides
2. Histogram analysis
3. Cluster analysis - unsupervised
4. Decision tree - supervised
5. Correlations

for example χ^2 analysis
Unsupervised

Binning -

Equal Width - N bins of equal
size - outliers cause havoc
so do data smoothing or
noise reduction first

$$\text{width} = W = (\text{High-Low}) / N$$

Equal Depth

each partition contains
same # of samples

Classification

e.g. decision tree

use entropy to determine

split points - top down, recursive splits

Correlation Analysis

supervised

Bottom-up merge

e.g. χ^2 -based discretization

If you need to convert categorical
to numerical attribute

standard spreadsheet format

Ex

<u>outlook</u>	<u>leaves</u>	<u>outlook</u>	<u>outlook</u>	<u>sunny</u>	<u>rain</u>
<u>sunny</u>		<u>overcast</u>		1	0
<u>overcast</u>			1	0	
<u>rain</u>			0	0	1

Data Reduction

Techniques

Data cube aggregation

Dimensionality reduction

Discretization

Numerosity reduction - reduce # rows

regression

histograms

clustering

sampling

Data cube aggregation

reduce data to the concept
level needed in the analysis

Dimension reduction

Curse of dimensionality - as data dimensions increase, data starts to become sparse

Dimensionality Reduction

avoids the curse
reduces noise

reduces time/space requirements
easier to see patterns

Techniques for Dimensionality
reduction

PCA

attribute subset selection
attribute or feature generation

PCA - find projections of data
that captures largest amount
of variation in data

Original data projected onto
smaller space, resulting
in dimensionality reduction

This is done by finding eigenvectors
of the covariance matrix, and
the e-vectors define the new space

Steps for PCA

N rows

n dimensions (attributes)

find K ($\leq n$) orthogonal vectors
(ie - principal components) that
can best be used to represent
your data

1. Normalize the input data
2. Compute K orthonormal (unit)
vectors - these are the
principal components
3. Each input data (vector) is
a linear combo of the K
PCVs
4. sort the PCs in decreasing
order of strength
gives you order of variance
they represent
5. Eliminate those w/ weak
components

attribute generation

create new attributes

that capture the important info more effectively than original

1. Attribute extraction

2. mapping to new space (like in data reduction)

3. attribute construction

a. combine features

b. data discretization

Numerosity reduction - reduce rows

Parametric

regression is derived original data

using least squares method
replace data points w/ the model

Histogram - keep buckets

Cluster / consider items in same cluster as same item

Sampling - if normally distributed
not appropriate if skewed
Stratified - data is grouped into
strata is taken a sample from
each stratum

can we clustering for
strata

Correlation Analysis

for nominal data use χ^2

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{Expected}}$$

larger the χ^2 , more likely they
are related

cells that contribute most to
 χ^2 those where count is
very different from expected

for numeric data use correlation coeff
(Pearson's)

r = correlation
 $+$, $-$ or 0 (not correlated)

correlation is a measure of the linear relationship b/w objects

so you standardize the objects
and take their dot product

$$a'_k = (a_k - \text{mean}(A)) / \text{STD}(A)$$

$$b'_k = (b_k - \text{mean}(B)) / \text{STD}(B)$$

$$\text{corr}(A, B) = A' \bullet B'$$

Cross Tab to look at
visual pattern of 2 categorical
variables
Can plot w/ Bar charts

2 pt → distribution where, for ex,
look at frequency of word, and
sort in descending order

2 - more normalization

your margin std dev away
from mean

Discretization: turns numerical
data into categorical

inc_bins = pd.qcut(vstable.Income, 3)

↑

Hq bins

can specify labels, too
and specify the quantiles

groupby - against categorical vars
to get aggregated state

Plotting categorical - Bar
using value - count

vstable[["genre"]].value_counts().plot
(kind="bar")

scatter plot

alpha is transparency

$S \rightarrow$ size of circles

$C \rightarrow$ color according to
color map (cmad)

Crosstab

can we group by w / multiple
variables or

$gg = pd.crosstab(cat.\text{ variables})$

then you can plot it

Opposite of discretization (going from
numerical attribute to categorical).

can convert categorical to numeric
using pd.get_dummies

need this to compute similarities
or distances, or correlation
analysis

When data is in all numeric form
said to be in "standard spreadsheet"
format

You also need to normalize
so it is in same scale
- correlation analysis
- distance

Any Function

vsTable [vsTable.isnull(),

any (axis=1)]

gives me rows where there
is a null

vsTable [vsTable.Gender, isnull()]
values where just Gender
is null

Fill in missing values

vsTable, Age, fillna (age - mean)
axis=0, inplace=True)

look across rows

Distances & Similarity Measures

Proximity measures

distance used in context
of categorical variables

Similarity
often normalized to fall
b/w [0,1]

Euclidean - not normalized
so no upper bound

Grouping similar items requires
measuring distance b/w objects
requires representation of objects
as "feature vectors".

Distance Matrix - distance one row
to another $O(3,1)$ distances b/w row
3 and row 1

Proximity measure of
categorical

based on matching
how many attributes match

$m = \text{# matches}$ (color match, shape match)

$P = \text{total # of variables}$

$$\text{method 1} \rightarrow d(i, j) = \frac{P - m}{P}$$

$\uparrow \quad \uparrow$
2 objects % of attributes where they match

% of attributes where they match

method 2 - convert to std spreadsheet
then use vector-based similarity
(numerical based)

Binary distance Proximity measure
2 rows, object i & object j

		Obj j	
		1	0
Obj i	1	q	r
	0	s	t

so q is # of times i & j match

Distance for symmetric binary vars
 $\frac{r+s}{q+r+s}$ → where different
 $q+r+s+t$

symmetric - care about when both match

asymmetric $\frac{r+s}{q+r+s}$

Jaccard

$$\frac{q}{q+r+s}$$

Symmetric - care about matches
0z and 1z

Asymmetric - maybe zeros don't count

Document 3 Terms

only care when Term appears
in several docs, don't care if it doesn't appear
or does not include 3

Jaccard - look at similarity
only

if 2 items have maximum
similarity, $J=1$

Have to standardize or normalize
data to compare

Common Distance Measures for Numeric Data

Consider 2 rows

$$x = (x_1, x_2 \dots x_n)$$

$$y = (y_1, y_2 \dots y_n)$$

Manhattan

$$\text{dist}(x, y) = |x_1 - y_1| + |x_2 - y_2| \dots$$

Euclidean

$$\text{dist}(x, y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$$

Dual (similarity measure)

$$\text{sim}(x, y) = \frac{\sum x_i \cdot y_i}{\sqrt{\sum x_i^2} \times \sqrt{\sum y_i^2}}$$

$$\text{dist}(x, y) = 1 - \text{sim}(x, y)$$

Manhattan - location 1 block this way, 2 blocks that way

Euclidean - straight line distance

Minkowski Distance

The generalization

$$d(i,j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

Manhattan is when $h=1$ (AKA L_1 Norm)

Euclidean is when $h=2$ (AKA L_2 Norm)

$$\sqrt{|x_{i1} - x_{j1}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

Both are Distance

Vector-Based similarity - take dot product

$$\text{sim}(X, Y) = X \cdot Y = \sum x_i \cdot y_i$$

scalar value represents how similar they are

Cosine similarity is the normalized dot product

$$\text{sim}(X, Y) = \frac{X \cdot Y}{\|X\| \times \|Y\|} =$$

$$\frac{\sum_i (x_i \cdot y_i)}{\sqrt{\sum_i x_i^2 \times \sum_i y_i^2}}$$

) norm, magnitude
of vector

$$\text{distance} = 1 - \text{sim}(X, Y) \quad \text{SEE SLIDES (4)}$$

you are not concerned w/magnitude
of vector, just the distance

$\cos 90^\circ = 0$ orthogonal

$\cos 0^\circ = 1$ on top of one another

Example SL:dp 17

Document similarity example

71, 72 73...

Doc1 0 4 0 ...

Doc2 3 1 4 ...

Doc3 0 1 0 ...

: : : : Similarity of Doc1, Doc3

1) get dot product (Doc1, Doc3)

2) get Norm doc 1

get Norm doc 3

3) cosine Doc1, Doc2 = 10 / Norm1 * Norm2

you could also use correlation
as a similarity measure.
Good where high mean values
across objects

Pearson

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\text{stdev}(x) \cdot \text{stdev}(y)}$$

used in recommender systems
used to compare rows, objects
(before we looked at cov between
variables (columns)).

How are distances used in Classification
KNN Classify new instance based
on similarity to or distance
from instances we have
seen before

Compare object you want to classify
to all other objects

KNN - slide 21

given object X , find k most similar objects to X

find class label for those neighbors & vote for what class X should be in

KNN can also be used for numeric prediction

find value of attribute from neighbors

then return a weighted average of neighbors as predicted value

WEEK 3 1/28/20

Classification Techniques

Cos Similarity is

$$\text{sim}(x, y) = \frac{x \cdot y}{\|\text{norm } x\| \cdot \|\text{norm } y\|}$$

forwards direction, not magnitude

$$\text{norm} = \sqrt{\sum x_i^2 + \sum y_i^2} \rightarrow \text{unit vector}$$

Euclidean and Manhattan - Vector Distance

$$\text{distance also} = d(x, y) = 1 - \text{sim}(x, y)$$

X and Y are 2 vectors

vector is a record, a row

Cosine similarity ignores magnitude
and measures angle between 2 vectors

Distance comes up in KNN
given X, find K most similar objects

If you have high dimensional sparse dataset,

KNN for classification - video store

example .51:00

1. split into test / train
Evaluate error rate w/ test partition
KNN w/ training partition
good to randomize rows / st to reduce bias
2. create var w/ column names
create var for target column
 $\text{VS_Target} = \text{VS_Incidentals}$
Convert categorical data into SSF
excluding Target
and getdummies
3. Split
choose training / test 90
split data frame
split target

4. Preprocessing Normalize Data

from sklearn import preprocessing
does normalizing for you !!

min-max-Scale = preprocessing.MinMaxScaler
Create normalizations object
min-max-Scale.fit(X_train)
Apply to training data
Create normalized DF's for both
train & test

Scaler returns numpy array

5. Run KNN function 1:07

X = new instance (test)

D = training data, normalized
for K, try different values and
see how it's error rate
too small - overfit
too large - underfit

the distance return will be
normalized distance

6. Classify - voting 1.11

there is a method called Counter
from collections import Counter
print(Counter(neigh_labels))

most_common() gives you majority
Counter(neigh_labels).most_common(1)

↑

both parts of
tuple

Numpy - takes input and creates that
many copies of it

7. Calculate error rate

error count / test instances

Classification

when prediction is discrete
or nominal values

(e.g. class/category labels)
fixed set of categories

Regression -

predicting continuous variable

Consistency in ML means

hypothesis consistent
w/ training data

What you are trying to get to

Hypothesis?

2:01

shape = circle \rightarrow category = positive

red \rightarrow positive slide 5

true? consistent

May have to get more & more specific
to get to consistency

shape = circle & color = red \rightarrow category
positive

Classification accuracy
% of instances classified
correctly

Lazy - unlike decision tree, KNN
doesn't build model beforehand
Happens when you go to
classify a new instance

Classification is a 3 step process

2.05 slides

1. Model construction (learning)
2. Model evaluation (accuracy)
3. Model use (classification)

Classification methods

DT

SVM

KNN

Genetic Algorithms

Bayesian

Random Forest

→ Logistic Regression

ensemble

Neural Net

Linear Discriminant Analysis

Evaluating Models

Training Set

Test Set

Evaluation Set \rightarrow used to test model performance

Cross validation (slide 13)

divide data (randomly) into n folds
each w/ same # of records

Each model is trained w/
 $n-1$ folds & tested w/ remaining
fold

Used to find best algorithm &
optimal training parameters
steps - slide 13

assess final model using
evaluation set

For example - 5 folds
train on 4/5ths, test on 1/5th
do it again w/ another 4/5ths
until all data used
average results

Then Δ parameters ? do it again

Build model on full training data
not just folds

See slide 14 for graphic

Another approach - Bootstrap
Validation
sampling w/ replacement
n instances in data set
sampled n times

Measuring Effectiveness

in Classification, most common measure is accuracy using confusion matrix (S1: 10 16 2.22)

		Predicted	Total
Actual	T	F	20
	18	2	
F	3	15	18

Total	21	17	38
-------	----	----	----

error rate is inverse of accuracy

		Predicted	
Actual	T	TP	FN
	F	FP	TN

$$\text{accuracy} = \frac{TP + TN}{\text{Total}}$$

$$\text{error rate} = 1 - \text{accuracy}$$

$$\text{Sensitivity} = \frac{TP}{\text{all P}} \quad (\text{class label } \begin{matrix} \text{apparuse} \\ \text{rarely} \end{matrix})$$

$$\text{Specificity} = \frac{TN}{\text{all N}} \quad (\text{class label } \begin{matrix} \text{cancer} \end{matrix})$$

other measures

Precision

% instances classifier predicted
as + that were + $\frac{TP}{TP+FP}$
perfect = 1

Recall

% + classifier predicted correctly
as + $\frac{TP}{TP+FN}$
perfect = 1

F score

mean of precision & recall

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

Numeric Prediction / Estimation
typically uses some form
of Regression

Linear \rightarrow response variable y
predicted by single variable x

Multiple Linear \rightarrow $> 1 x$

use least squares

Non - Linear regression

$$y = w + w_1 x + w_2 x^2 + w_3 x^3$$

can convert to linear w/
new variables

$$x_2 = x^2, x_3 = x^3$$

Regression tree

CART - used both for class
and regression

Model Tree

generalization of Regression Tree

There do better than linear regression

Evaluation Metrics of Prediction
(Numerical Predictions)

how close you stand to actual

$$MAE = |P_i - a_i| + |P_n - a_n| / n$$

$$RMSE = \sqrt{\frac{(P_1 - a_1)^2 + \dots + (P_n - a_n)^2}{n}}$$

std deviation of error

High Bias - model underfits
Low Var data

Low Bias - over fit
High Var

High Bias (underfitting) ?

add features

use more complex model

Low Bias - overfitting

more training instances

↓ features

simpler model

regularization coefficient - used in
regression

Text Categorization Learning

Bayesian (naive)

NN

relevance Feedback

Nearest Neighbor

SVM

similarity - Cosine most effective
of composite weights for terms
based on how frequently
term appears across docs
 $TF-IDF$

TF = term frequency

IDF = inverse doc frequency

goal: assign a $TF \times IDF$ to each
term in each document

IDF - penalize words that appear
in every doc, like stop words

Relevance Feedback (Rocchio)

get initial results of query
get feedback

modify original query
re-issue search query

Read TF-IDF notebook

Bayer class video

Decision of see video

General ML Procedure

1. Get Data

2. Clean data

3. Separate Class into its own
dataframe → target

4. Get Dummies for cat data

5. Train-Test-Split w/ sklearn which

① gets train, test, target train,
target test

giving it training data and
target and test size (0.2)

(random state value will give you same
random data vs. new random each
time)

6. Normalize Train & test data

a. min-max-Scaler = preprocessing. MinMaxScaler
().fit(X-train)

b. apply to train & test datasets

7. Import & Fit classifier

8. Call Predict function on a TEST instance

and print the test

exampl

Create classifier

knnclf = neighbors. KNeighborsClassifier(
5, weights='distance')

Fit it

knnclf. fit (vs-train-Norm, vs-target-Norm)

Predict

knnclass-test = knnclf. predict (vs-test-Norm)

10. Run Classification Report, pyc. to Evaluate Accuracy

from sklearn. metrics

print (classification_report (vs-target-test, knnclass-test))
↑
labels

11. Get Avg Accuracy across test Instances

print (knnclf.score (vs-test-norm, vs-target-test))

12. Compare to Performance on training

data to ✓ for over/under fitting

print (knnclf.score (vs-train-Norm, vs-target-train))

Don't need to normalize for Decision Tree
Overfitting - Train Data Score high
TEST Data Score low

① regression - get X_{train} y_{train}
 X_{test} y_{test}

Set aside Test set

model selection (looking at
parameters, diff model(s) done

on Training

→ Cross-validation done on training
Once you get params of model
Then train model on full Training
Then test on test.

Typical model use

create model object

Trains
model

Fit using that mode on data

Evaluate

$\text{predict}(\text{y_predict} = \text{clf.predict}(X))$

Show Accuracy

Show Classification Report

Show Confusion matrix

Feature Selection

from sklearn import feature_selection

most informative attributes

percentile = 30 \rightarrow top 30%

measure correlation between

response & variables

classification uses cat vars or user χ^2

\rightarrow regression user + values - independently

does not look at interactions among
variables or comb of variables

$X_{train}^{fs} = fs.fit_transform(X_{train}, y_{train})$

smaller subset of attributes

2.24 2/11/20

fs.get_support tells you how feature did (χ^2)

fs.scores $\rightarrow \chi^2$ values \rightarrow higher better

Now need to evaluate performance of
new smaller feature set

- do training again

1. Create classifier (like dtree...)
2. fit it (X_{train} , y_{train})
3. Transform on TEST $\xrightarrow{\text{FS}} \text{transform}(X_{\text{test}})$
4. Evaluate performance

Can use cross-validation to do feature selection to determine percentile (line 24)

finds optimal percentile
 Then translate % to column #
 there are my optimal # of features

Evaluate the best # of features on test set
 Feature-selection. Select Best (..., optimal_num_features)

$X_{\text{train}}_{\text{FS}} = \text{FS. Fit_Transform}(X_{\text{train}}, y_{\text{train}})$

$dT = \text{classifier}$
 $dT. \text{Fit}$

$X_{\text{TEST}}_{\text{FS}} = \text{FS. transform}(X_{\text{TEST}})$

measure performance

Model Selection

determine best param for model

You can see model params w

```
print(dt.get_params())
```

(one way)

1. build classifier object w/param

2. scores = cross_val_scores(dt, Xtrain, ytrain,
 $CV=5$)

3. scores.mean()

4. try another value for param 3
run again

like criteria = 'entropy'

criteria = 'gini'

for example

5. try best on full training data
vs. cross validation

6. measure performance

You fit to train (X_{train}, y_{train})
Measure performance on test

For more systematic way
to test params, look
at function calc_params

line 32

User X validate to test
parameters

Classification \rightarrow score returns
classification accuracy

Regression \rightarrow score needs to return
MSE or RMSE

* Grid Search CV gives more systematic
way to look at different number of
multiple params

from sklearn.model_selection import GridSearchCV
so much easier

Regression steps in general

See Ridge Regression in

Regression - scikit-learn.ipynb

1. Create regression object

$\text{ridge} = \text{Ridge}(\alpha=0.8)$

2. Train model using training data set

$\text{ridge}.fit(X, y)$

3. Compute RMSE on Training data

$P = \text{ridge}.predict(X)$

$err = P - y$

Total Error = $\text{np.dot}(err, err)$

$\text{rmse_train} = \sqrt{\frac{\text{Total_Error}}{\text{len}(P)}}$

4. Compute RMSE using 10-fold cross validation

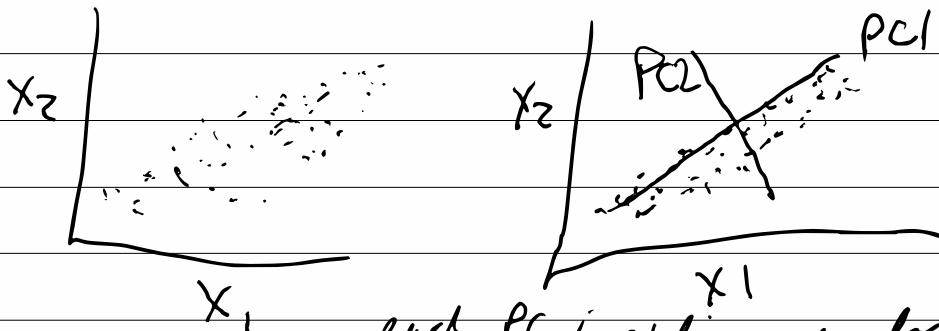
$\text{rmse_10cv} = \text{cross_validate}(\text{ridge}, X, y, 10)$

see Notebook for k -folds

compare RMSE and RMSE-CV

PCA

orthogonal means they are not correlated - reduces noise



each PC is a linear combination of original 2 variables

PCA

data matrix A

n objects

P variables (which may be correlated)

summarizes it by creating uncorrelated axes (PCs)

that are linear combinations of the original P values

Computed using e-vectors 3 e-values of the cov-variance matrix

so each component is a linear combo of original variables

co-variance of 2 variables is their tendency to vary together

1st PC captures the most variance in the data

2nd next most variance
first few give you most of the info in the data

you will ultimately have P PCs

Computed using COV matrix and
 λ -values and λ -vectors
PCs are the λ -vectors

the day!

You have P original variables

x_1, x_2, \dots, x_p

Produce new variables

$y_1, y_2, y_3, \dots, y_p$

$$y_1 = a_{11}x_1 + a_{12}x_2 + a_{1p}x_p$$

$$y_2 = a_{21}x_1 + a_{22}x_2 + \dots + a_{2p}x_p$$

:

y_i 's are the Principal Components,

The a 's are the values in
the eigenvectors

y_i 's are uncorrelated w/ each
other (orthogonal)

y_1 explains the most variance
in data set

So you pick the top k y 's,
drop the rest

1st component is y_1 . - 2nd will be orthogonal to that y_2

Eigenvalues represent the variance of the PC



Eigen vector \rightarrow the coefficients in the linear combination

$$z_1 = [a_{11}, a_{12} \dots a_{1P}] \rightarrow \lambda_1$$

eigen vector of the covariance matrix
is coefficient of 1st PC

$$z_2 = [a_{21}, a_{22}, a_{23} \dots a_{2P}] \rightarrow \lambda_2$$

2nd eigen vector of cov matrix and coefficient

Dimensionality reduction - you don't need all P variables. You try to figure out how much variance is captured in the top PCs starting from 1st and going down (the lambdas)

For independent variables, $\text{cov}(x,y) = 0$
 $\text{cov}(x,x) = \text{var}(x)$
Cov matrix is square & symmetric
diag for independent variables are just the variables themselves
and cov are the non-diag elements

To calculate cov matrix from dataset
1 - center data by subtracting the mean of each variable
2 - compute $\frac{1}{n-1} (A^T A)$
example slide 10
use this table to find PCs

Eigenvalues are the variance
of the coordinates on each PC axis

Do you recognize the matrix
and the subset of eigenvalues? 3
vectors is compared to original
matrix

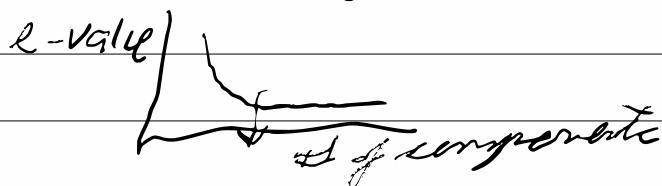
Higher $\lambda \rightarrow$ contains more info

Sum of λ 's called the trace
score for a PC -

$$y_{ki} = z_{1k}x_{1i} + z_{2k}x_{2i} + \dots$$

y is $n \times k$ matrix of PC scores

Screep plot - the knee is
the # of components



Single Value Decomposition

$$\text{Matrix } D = m \times n$$

m customers n items

use where

- D is sparse (less than 1% of entries have rating, for example)
- n is large
so finding matches to less popular items difficult

Main idea:

compress the columns (items)
into lower dimensional representation
decompose

$$D = U \Sigma V^T$$

$m \times n \quad m \times n \quad n \times n \quad n \times n$

$V^T \rightarrow$ e.vectors of $D^T D$

Σ is diagonal - singular values -
like PC's

U rows are coefficients \downarrow

show variance

Numpy SVD - pass to it matrix
you want to factorize
gives out V , S , VT matrices

V - $m \times m$ - LEFT singular values of original matrix

S - matrix of singular values $m \times n$

VT - v transpose $n \times n$ - right Singular Values

$$\text{original matrix} = np.dot(V, np.dot(np.\text{diag}(S), VT))$$

SVD gives you top \rightarrow largest singular
values in S from
select the columns and the
rows from VT

goal is smaller # of dimensions
pick small # of singular values (S)
corresponding # of columns in
 V matrix and rows in VT matrix

$\text{low_rank} = \text{U}\text{.dot}(\text{V}\text{.T}, \text{nP}\text{.dot}(\text{sigma}, \text{V}\text{.T}\text{.T}))$

V^T is the lower dimension version
of original data

TD in example

3 dimension (vs 10 dimension)
version of each of the documents
3 "terms" vs. 10

Doct for info retrieval

line 65 - give you lower dimension
version of query

You also have lower dimension
version of document

(7) - measure similarity between
query and documents

Combined terms are tele-themes

like clustering terms!

Recommender Systems:

You have users ratings (r)
items

measure utility of items for
user v

recommendation to each user made
by offering items w/ highest
predicted rating

some function R : user \times items \rightarrow
ratings

Collaborative Recommender system
Neighborhood based models
User based

K most similar users

Knn

Predict rating on item as
weighted avg of ratings by neighbor

model-based

item-based (items vector of user ratings)

- item-item similarity

measured by if rated similarly
by users

- prediction for target item

based on user's own rating on
similar items

Content-based

based on content of items

movies - based on ratings of
director

pairwise similarity of items

find $P \in Q$ that gives you lowest
sum of square error

using some gradient descent
converges

SVD for Matrix Factorization

$$M_K = U_K \times \Sigma_K \times V_K^T$$

U_K = user-feature matrix

V_K^T = feature item matrix

Ratings prediction computed using
dot products

where K is lower dimension

Evaluating Recommender System

normally - separate train/test
of ratings for each user

Use Knn or whatever on train data

Predict set aside ratings

Look at difference between predict
and test (MAE or RMSE)
 \uparrow
regression

Prob 2 - Difference is your set aside a
portion of each user's ratings

Stand Est

data matrix is user \times item matrix
user - index for one of the users
item - index
goal - generate predicted
rating for this user, this item

Finds all items similar to item

measures similarity of all items

Then looks at user ratings
of those items

greater weighted avg of all
ratings is the rating

user = user

item = j

if rating is 0, do nothing
if rating not 0 look at
overlap of item j and
item class looking at
in terms of overlap

if overlap is not 0, find
similarity
predicted rating will be weighted
avg of all ratings
 $\frac{\text{user rating} \times \text{similarity}}{\text{all similarity}}$

return rating

recommend

Neg of recommendation you want
for each user, predict rating
of all items 3 returns top N

SVD Est - only differences:

estimation vs. SVD based

1st does singular value

decomposition of data

Creates transformed representation

of items based on transformed
feature space of items

Visualized it as dimension space

do not use this for assignment 4
(listing 14.3)

Test of item based recommendation system
all similarity measures

[12] reduced dimension space for
items reduced to 4 dimensions
now all items
recommend function
pass data matrix
a user ($\#4$)
similarity
estimate method (old items -
based)
returns index of item 3 predicted
rating

Jesse Raunay

use gradient descent optimization
approach for learning the P and Q

models the rating data

rows = users

columns = movies

values = rating

matrix-factorization functions

uses gradient descent to

learn matrix $P \times Q$

slow - full GD, not stochastic

3000 steps (epochs)

Prediction for user 975, movie 9

dot product of $P[975]$, $Q[9].T$

if you want all predictions, $\text{dot}(P, Q.T)$

$\text{preds} = \text{np.dot}(P, Q.T)$

includes ratings that were
already known

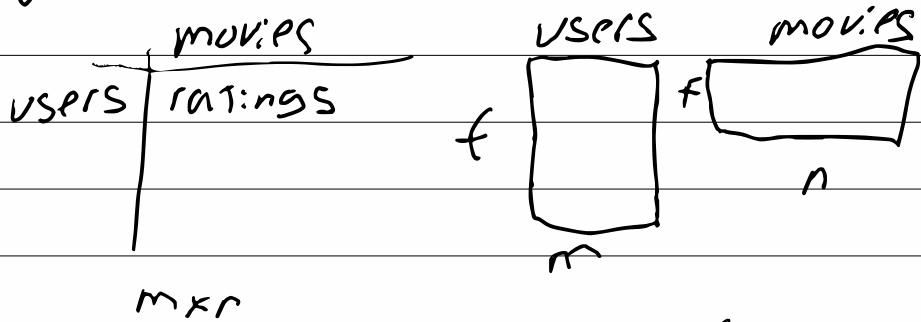
Fast - goes thru rating data,

generate predictions

compute, find error, avg for MAE

Matrix factorization is good
for sparse data like user
ratings in Netflix
we factorization of a sparse
matrix to create two
dense matrices

e.g.



multiply
to improve values
also for reducing large #
of features to fewer

Support Vector Machines

another supervised method

primarily used for classification
particularly for image classification
and text classification

Many class. models use distribution
to determine labels for new
points (like Naive Bayes)

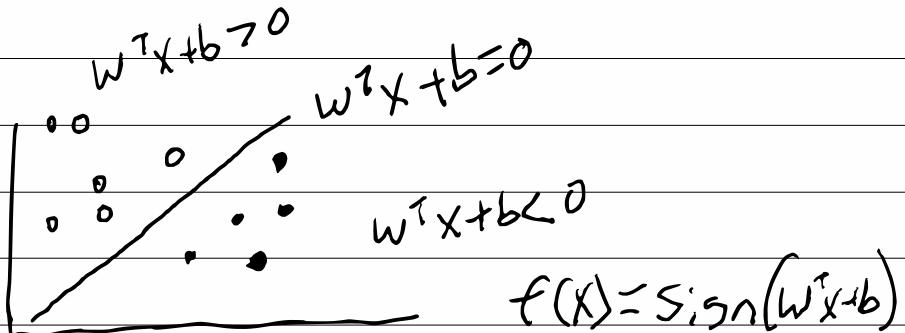
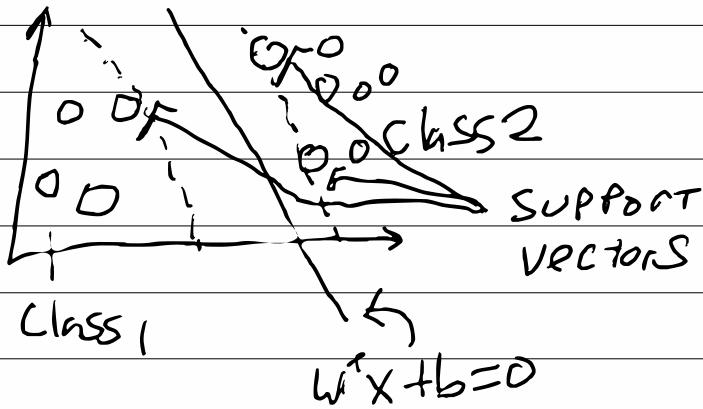
SVM: discriminative classification:
instead of modeling each class,
find a line or curve in 2d
that divides the classes from
each other

LOA is another discriminative
classification

Want a decision boundary that is max
distance from both classes

SVM tries to find that boundary - called maximizing the margin

The points that push up against the margin are the support vectors



$$\text{norm} = \sqrt{\text{sum of squares}}$$

Nice because only points close about one the support vectors

Size of dataset doesn't matter

Soft Margin Classification
when training set not clearly linearly separable

add slack variable to manage noise resulting in margin called "soft"

Conjugacy and priors

Subjective Prior

expert opinion used to determine probability

Objective Prior

I don't know much about problem

so I put in a non informative prior

no prior info available
use a wide distribution

like Uniform - all probabilities equally likely
or diffuse distributions

Conjugacy

when posterior distribution
is the same family as prior,
but w/ new parameter values

So, for example

Beta prior - Beta posterior

↳ Binomial

Advantage - do not need to
calculate integral

Week 8

Use MAP if you need a point estimate
for the posterior (maximum value)

Use MLE if you need the likelihood

EAP - expected A Posteriori - expected
value of posterior distribution

Trace plot - estimation from sampling
pm.find.map - gives you a point estimate

PM, Summary (trace)

mean

interval

Markov chain

for each one of your unknown
parameters

Bayesian sampling methods

Used to find posterior density
distributions

You sample from posterior
distribution to approximate
the posterior

Become more exact as more
samplers are used

PMC3 - need hundreds of thousands of
samplers

Markov Chain Monte Carlo sampling (mcmc)
main method to sample
many variations

Monte Carlo is just a sampling method

Markov Chain - all info from past captured
in my present and can be used
to predict future
Used to create prob distributions over
time

Metropolis Hastings - type of MCMC
a variation
weighted random walks
accepts or rejects proposal from MC