

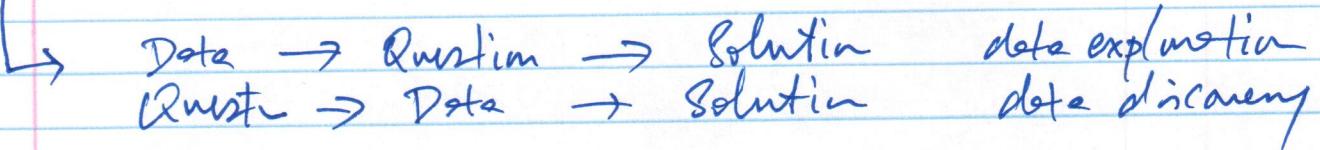
Data Scientist Nanodegree Program

2. Introduction to Data Science

Lesson 1: Data Science Process

CRISP-DM - Cross-Industry Standard Process for Data Mining

- business understanding
- data understanding - data collection, sorting, cleaning
- prepare data
- data modeling
- evaluate results
- deploy



Descriptive analytics - basic statistics to answer basic questions

- e.g. percent,
histograms
mean, max, min
summary ...

Data modeling - predictions.

e.g. machine learning

- Points of concern for passing data to ML models in sklearn:
 1. missing values
 2. categorical values

3 strategies:

1. remove or drop rows with missing values
2. impute missing values
3. build models around them

- Removing data can lead to bias in model.
Track missing values using indicator values or counts.

OK to drop missing data

- errors
- not relevant
- no variability, or constant values

Drop column

- significant % of rows missing values
- replace with indicator of missing or present values (e.g. 1 for present, 0 for missing)

Importing

- most common way to fill missing values

- mean, median, mode

- ML and Bayesian statistical approaches

- k-nearest neighbors

- AMELIA library

- introducing Bias

- lead to overgeneralize

- dilute importance of feature, by reducing variability of feature and thereby power of predict

- common -

- impute mean of numerical

- impute mode of categorical

- use KNN

- impute either a very small value, or very large number to differentiate missing values.

- Categorical variables - convert to numbers
 - easy
 - flexible
 - scaling issues
 - added columns may lead to inverting solutions

} advantages
} disadvantages

pandas.get_dummies(df, prefix = ' ',
 prefix_sep = ' ',
 drop_first = True,
 dummy_na = True)

Overfitting - model does not generalize to data not seen before.

Cold start problem - recommend engine with no prior data

steps

1. how many mining rules
2. remove data
3. import data
4. categorical to numerical
5. different cutoffs for number of non-zero values / mining rules
 vary and find optimum outcome per cutoff.
 determine how much variability is needed.

git

- git clone
- git add *
- git commit -m 'commit here'
- git push

b. separate predictors and response

c. split into train + test dataset

d. fit to model

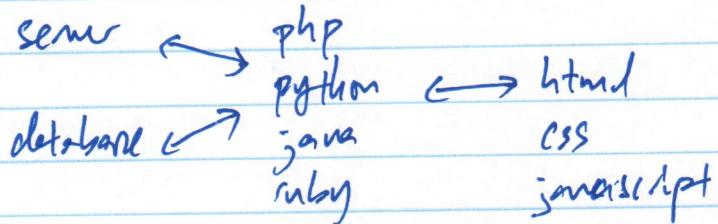
e. score and evaluate.

3. Software Engineering

Lecture 1 - 5: Introduction to Software Engineering

Lecture 6: Web Development

- Basics of web app
 - html - content
 - css - style, color, layout
 - javascript - interactivity, animations
- Front-end libraries
 - bootstrap
 - plotly - visualization
- Back-end libraries
 - flask
- Deployment



Example HTML

```
<html>
<head>
  <title>    </title>
</head>
```

```
<body>
```

```
  <h1> Example HTML </h1>
  <a href = "https://www.google.com" > Home Page </a>
  <p>    </p>
  <img src = "photo.jpg" alt = "photo" >
</body>
</html>
```

```
<div>  
  <p> This is a <span> test </span> </p>  
</div>
```

ID - name for each tag, used one time only
- e.g. `<p id='first-id'> </p>`
`<div id='first-div-id'> </div>`

class - class name reusable for multiple tags for same styling.
- multiple classes together
- e.g. `<p class='class-one' 'class-two'> </p>`

Writing CSS

- inline

```
<p style="font-size: 20px;"> Test </p>
```

- style sheet within HTML

```
<head>  
  <style>  
    p { font-size: 20px; }  
  </style>  
</head>
```

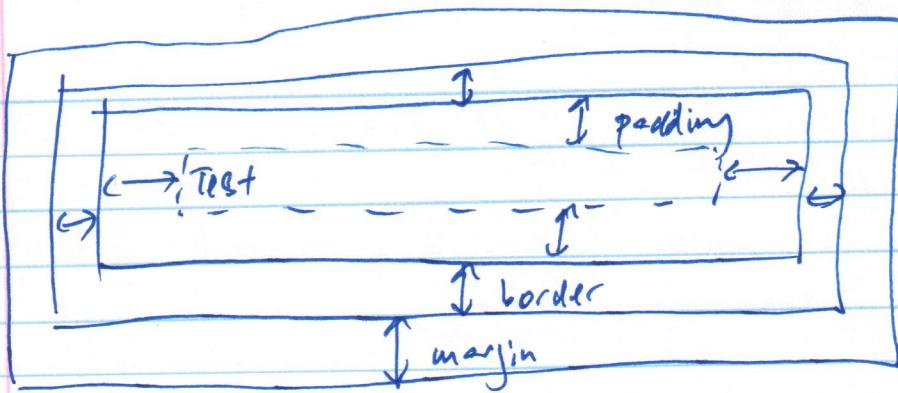
- separate .css file

```
<head>  
  <link rel="stylesheet" type="text/css" href="style.css">  
</head>
```

```
p {  
  color: red;  
}
```

```
a {  
  text-decoration: none;  
  color: blue;  
  font-weight: bold;  
}
```

```
.class_name {  
}  
#id_name {  
}
```



Specifying sizes :

1. px
2. percentage
3. em (default font size)

Bootstrap library

- class = "img-fluid"
- div class = "col-x"
- div class = "mt-x" "mb-x"
"ml-x" "mr-x"
- h4 class = "text-muted"

JavaScript

```
- jQuery - JavaScript library
<script src="https://ajax.googleapis.com/.../jquery>
</script>
<script>
$(document).ready(function () {
    $("img").click(function () {
        $("h1").text("Replace text");
    });
});
</script>
```

- high level language - python, PHP, Ruby, C++
- newer JavaScript tools - React, Angular for front end
- node.js for backend

ready(function ()) - wait for html document to load
\$("img") - element
\$("p#first-id") - paragraph element with id = first-id

Plotly:

```
<script> src = "https://plotly ... "</script>
```

```
<body>
<div id="plot1"></div>
<script src="plot1.js"></script>
</body>
```

plot1.js -
var xdata
var ydata
var trace1 = {
 x: xdata,
 y: ydata,
 mode: 'lines',
 type: 'scatter'
};

```
var data = [trace1];
```

```
var layout = {  
    title: ' ',  
    xaxis: {  
        title: ' ',  
        tick0: ' ',  
        dtick: ' '},  
    yaxis: {  
        title: ' '},  
};
```

```
plotly.newPlot('plot1', data, layout);
```

Chart library

- d3.js
- chart.js
- google charts
- nvd.js - based on d3.js
- plotly. - Javascript, Python libraries
- syntax similar between 2 languages.

Flask

HTTP

- use in Backend - route files to render website from server to computer via HTTP
- written in Python - scikit-learn, Pandas, numpy

e.g.

```
worldbank.py - from worldbankapp import app  
app.run(host='0.0.0.0', port=3001,  
'debug = True')
```

```
worldbankapp/__init__.py - from flask import Flask  
app = Flask(__name__)
```

```
from worldbankapp import routes
```

/routes.py -

```
from worldbankapp import app  
from flask import render_template  
@app.route('/')  
@app.route('/index')  
def index():  
    return render_template('index.html')
```

```
@app.route('/project-one')  
def project-one():
```

```
    return render_template('project-one.html')
```

template/index.html, project-one.html

Prom information back + forth from backend to frontend.
render-template ('index.html', data-set = data)

index.html:
aken with {{ data-set }}

Flask comes with template engine - Jinja.
{% for tuple in data-set %}
 <p> {{ tuple }} </p>
{% end-for %}

With Plotly: rental.py
import plotly.graph_objs as go
import plotly.json

```
graph-one = [go.Scatter(  
    x = data[0][1],  
    y = data[0][2],  
    mode = 'lines',  
    name = 'country'  
)]
```

```
layout-one = dict( title = ' ',  
    xaxis = dict( 'title = ' ,  
        autotick = False,  
        tick0 = 1988,  
        dtick = 25 ),  
    yaxis = dict( 'title = ' , ) )
```

```
figures = []  
figures.append( dict( data = graph-one,  
    layoutone = layout-one ))
```

```
ids = [f'figure-{i}'.format(i) for i, _ in enumerate(figures)]
```

```
figuresJSON = json.dumps(figures, cls=plotly.utils.PlotlyJSONEncoder)
```

```
@app.route('/')
```

```
@app.route('/index')
```

```
def index():
```

```
    return render_template('index.html', ids=ids)
```

```
    figuresJSON = figuresJSON)
```

```
index.html:
```

```
<script type="text/javascript">
    var figures = {{figuresJSON | safe }};
    var ids = {{ids | safe }};
    for (var i in figures) {
        Plotly.plot(ids[i],
            figures[i].data,
            figures[i].layout || {});
    }
</script>
```

Connect Flask to SQL database.

Heroku - PaaS
- free tier

```
web-app / worldbank.py
         / worldBankApp
         / scripts
         / data
```

```
python3 -m venv worldbankenv - create virtual env.
source worldbankenv/bin/activate
```

pip install flask pandas plotly gunicorn

cd web-app

curl https://heroku.com/....

heroku login

-install Heroku

worldbank.py - delete ^{line} app.run()

touch Procfile

web gunicorn worldbank:app

pip freeze > requirements.txt

git init

git add .

git commit -m 'commit'

git config user.email,
email
git config user.name
'name'

heroku create name-of-webapp --buildpack heroku/python

git remote -v

git push heroku master

→ https://name-of-webapp.herokuapp.com

set in environmental variable

heroku config: set SLUGIFY_USES_TEXT_UNIDECODE=yes

heroku config: set AIRFLOW_FPL_UNIDECODE=yes

heroku config

4. Data Engineering

Lesson 2: ETL Pipelines

ETL - extract, transform, load.

ELT - extract, load, transform

→ Extract :

1. CSV
2. JSON
3. XML
4. SQL
5. text files

1. CSV:

```
pd.read_csv(' ', dtype=str, skiprows=4)
```

2. JSON:

```
pd.read_json(' ', orient='records')
```

or, import json

```
with open(' ') as f:  
    json_data = json.load(f)
```

3. XML:

```
from bs4 import BeautifulSoup  
with open(' ') as f:  
    data = BeautifulSoup(f, "xml")
```

```
for record in data.find_all('record'):
```

```
    for field in record.find_all('field'):
```

```
        fields_dict[record['name']].append(record.text)
```

4. SQL

```
conn = sqlite3.connect(' ')
```

```
pd.read_sql(' query ', conn)
```

or, engine = create_engine('sqlite:///path.db')

```
pd.read_sql(' query ', engine)
```

API:

import requests

import pandas as pd

url = 'http://...'

r = requests.get(url)

r.json()

pd.DataFrame(r.json()['data'])

→ Transform:

- Combining data

pd.concat([df1, df2, df3])

pd.melt(df1,

id_vars=[],

var_name=' ',

value_name=' ')

df.merge(df2, on=[' ', ' '], how=' ')

- Cleaning Data

• missing values

• inconsistencies

• duplicate data

• incorrect encoding

df[[' ']].drop_duplicates()

• unique()

df[' '].str.replace(' ', ' ')

pd.numeric(df[' '])

pd.to_datetime('January 1st, 2017')

('01/01/2017 5:30', format='%d/%m/%Y %H:%M')

- encoding:

```
import chardet  
with open('mystery.csv', 'rb') as f:  
    print(chardet.detect(f.read()))
```

- missing data:

- for ML, can't have missing values.
but decision trees, can.

- impute missing data

- mean substitution, median, mode
 - forward fill, backward fill - best if ordered by time.

- $\rightarrow df.\text{transform}(\lambda x: x.\text{fillna}(x.\text{mean}))$

- $\rightarrow df.\text{set_values}('year').\text{fillna}(\text{method} = \text{'ffill'})$

- order of fill forward/backward will have different results.

- duplicate data:

- check for similar names, abbreviation, ...

- dummy variables

```
df.replace("", "", "", "")
```

```
df.replace("^\(1| )", "", "", regex=True)
```

- outliers

- PCA

- clusters

- quantile = $df.\text{quantile}(0.25) = Q1$

$$Q3 = df.\text{quantile}(0.75)$$
$$\text{Interquartile range, IQR} = Q3 - Q1$$
$$\text{max_value} = Q3 + 1.5 \text{ IQR}$$
$$\text{min_value} = Q1 - 1.5 \text{ IQR}$$
$$\bullet df[(df < \text{max_value}) \& (df > \text{min_value})]$$

Tukey-rule

find extremes

2D analysis:

```
model = LinearRegression()
```

```
model.fit(df.values.reshape(-1, 1),  
          df.values.reshape(-1, 1))
```

```
inputs = np.linspace(1, 2 * 10^6, num=50)
```

```
predicts = model.predict(inputs.reshape(-1, 1))
```

```
df.plot(x, y, kind='scatter')
```

```
plt.plot(inputs, predicts)
```

- Feature Scaling

- Rescaling = scale range from 0 to 1
 - keep distribution.

- Standardize - mean = 0 standard deviation = 1
 - also keep distribution.

$$x = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- Feature Engineering

- manipulation of existing features to create new features
 - mean, max, min
 - power, add, subtract.

→ Good!

```
df.to_json(' ', orient='records')
```

```
df.to_csv(' ', index=False)
```

```
df.to_sql(tablename, con=conn, if_exists='replace')
```

```
pd.read_sql('SELECT * FROM tablename', con=conn)
```

con.commit()

con.close()

```
conn = sqlite3.connect('db')
cur = conn.cursor()
cur.execute('sql query')
conn.commit()
```

```
[ cur.execute("SELECT * FROM table")
  cur.fetchall()
  conn.close()
```

```
cur.execute("DROP TABLE IF EXISTS test")
```

```
cur.execute("CREATE TABLE test
            (colA, colB, colC REAL, PRIMARY KEY
             TEXT, INTEGER, (colA, colB));")
```

```
for index, values in df.iterrows():
```

```
    colA, colB, colC = values
```

```
    query = 'INSERT INTO test (colA, colB, colC)
```

```
        VALUES ("{}","{}","{}");'.format(colA, colB, colC)
```

ETL process

1. read csv, create tables
2. each row → data frame → processing → row
3. row → insert into sql table

Lesson 3: NLP Pipeline

NLP Pipeline

1. Text processing
 - cleaning
 - normalize
 - tokenizat
 - stop word removal
 - part of speech tagging
 - named entity recognition
 - stemming & lemmatization

2. Feature extraction - extract + produce features
- bag of words
 - TF-IDF
 - word embeddings

3. Modeling - design a statistical or ML model

1. Text Processing

- text from web, PDFs, word, books etc.
- remove punctuation, common words (and, a, the, of, etc)

• Cleaning:

```
from bs4 import BeautifulSoup
soup = BeautifulSoup(r.text, 'html5lib')
```

```
summaries = soup.findAll("div", class_="course-summary-card")
```

{
import requests
r = requests.get('url')}

```
summaries[0].select_one("h3 a").get_text().strip()
```

```
summaries[0].select_one("div [data-course-short-summary]").get_text()  
.strip()
```

• Normalization:

text = text.lower() - convert to lowercase

```
import re
text = re.sub(r"[^a-zA-Z0-9]", " ", text)
- remove punctuation
```

```
import nltk  
nltk.download('punkt')  
nltk.download('stopwords')
```

- Tokenization:

- splitting each sentence into individual words.
words = text.split()

```
from nltk.tokenize import word_tokenize, sent_tokenize  
words = word_tokenize(text)  
sentences = sent_tokenize(text)
```

- Stop Word Removal:

```
from nltk.corpus import stopwords  
print(stopwords.words("english"))
```

corpus = collect of text

words = [w for w in words if w not in stopwords.words("english")]

- Part of Speech Tagging

```
from nltk import pos_tag
```

sentence = word_tokenize("I always lie down to tell a lie.")
pos_tag(sentence)

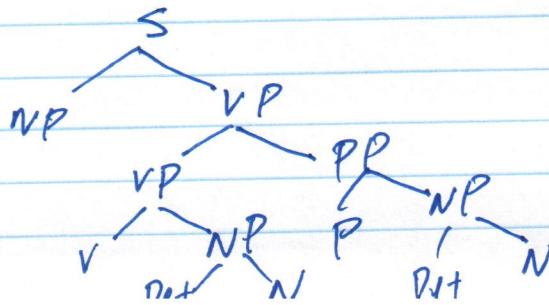
- noun, verb, pronoun ...

my_grammar = nltk.CFG.fromstring("S->NP VP
NP->NP VP
VP->V NP
VP->VP P NP
P->not
NP->NP N
NP->D Vt N")

parser = nltk.ChartParser(my_grammar)

sentence = word_tokenize("I always lie down to tell a lie.")

for tree in parser.parse(sentence):
tree.draw()

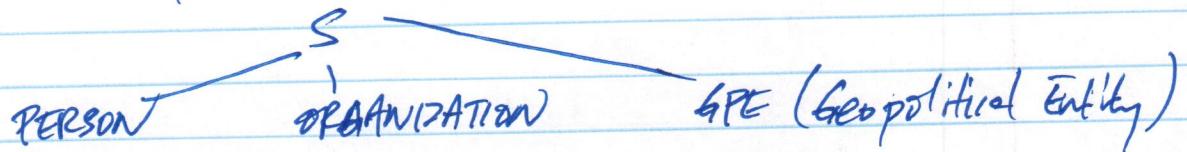


• Named Entity Recognition

from nltk import pos_tag, ne_chunk

from nltk.tokenize import word_tokenize

ne_chunk(pos_tag(word_tokenize(" ")))



• Stemming:

from nltk.stem import PorterStemmer

stemmed = [PorterStemmer().stem(w) for w in words]

remove - 'ing
- 'ed
- 'es

final and may be not a word

• Lemmatization:

is
being } is
was }

final word is a word based on
dictionary.

from nltk.stem.wordnet import WordNetLemmatizer

lemmed = [WordNetLemmatizer().lemmatize(w)
for w in words]

, pos='v'
→ verbs

2. Feature Extraction

• Bag of Words

unordered collection or set

	word 1	word 2	3	4	5	...	
document 1	1	1	0	0	1	...	Document -
2	n	m					Term Matrix
3							
4							

dot product - more similar 2 documents greater number
 but also true for non-matching terms

$$a \cdot b = \sum a_i b_i$$

$$\text{cosine similarity} - \cos(\theta) = \frac{a \cdot b}{\|a\| \|b\|}$$

$-1 < \cos \theta < 1$
 most 'dissimilar' most similar

• TF-IDF

- some words have more weights or importance.

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

$$\frac{\text{term frequency}}{\text{count}(t, d)} \quad \frac{\text{inverse document frequency}}{\log\left(\frac{|D|}{|\{d \in D : t \in d\}|}\right)}$$

- normalize or smooth values

Counter Vectorizer (Bag of Words)

from `sklearn.feature_extraction.text import CountVectorizer`

`vect = CountVectorizer(tokenizer=tokenize)`

`x = vect.fit_transform(corpus)`

`x.toarray()`

`vect.vocabulary_`

from `sklearn.feature_extraction.text import TfidfTransformer`

`transformer = TfidfTransformer(smooth_idf=False)`

`tfidf = transformer.fit_transform(x)`

`tfidf.toarray()`

$\text{TfidfVectorizer} = \text{CountVectorizer} + \text{TfidfTransformer}$
 from sklearn.feature_extraction.text import TfidfVectorizer
 vectorizer = TfidfVectorizer()
 $x = \text{vectorizer}.fit_transform(\text{corpus})$
 $x.\text{toarray}()$

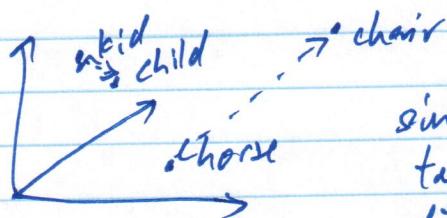
- One-Hot Encoding

- treat each word like a class
- create vector, can get large

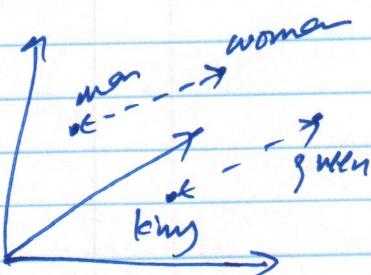
	word 1	2	3	4
doc 1	1	0	0	0
2	0	1	1	0
3				

- Word Embedding

- limit to fix size vector



similar words closer together and similar distance.

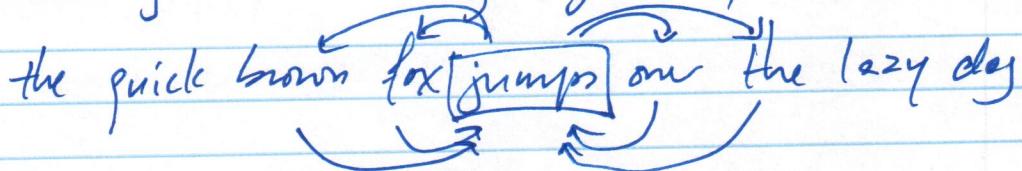


3. Modeling

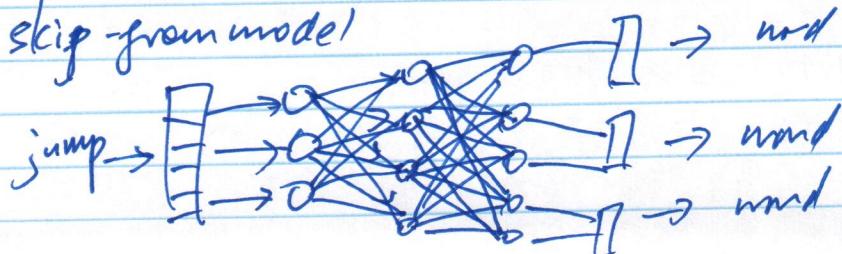
- fit parameters to training data
- use optimization procedure
- make predictions on unseen data

- Word2Vec

- predict neighbouring words given a word (continuous bag of words)
- predict given word using neighbouring words (skip-gram)



skip-gram model



- Glove - Global Vectors for Word Representation:

$$w_i \cdot w_j = P(j|i)$$

context | target
 word word

\ co-occurrence probability

- Embeddings for Deep Learning:

- words that occur in same context have similar meaning
 = Distributional Hypothesis

- t-SNE

- t-distributed stochastic neighbor embedding
- maintaining relative distance between objects
 (similar words have similar distances)
- explore grouping of similar words / images.

Lesson 4: Machine Learning Pipeline

e.g. def tokenize(text):

detected_uris = re.findall(url_regex, text)

for uri in detected_uris:

text = text.replace(uri, "placeholder")

tokens = word_tokenizer(text)

lemmatizer = WordNetLemmatizer()

clean_tokens = []

for tok in tokens:

clean_tok = lemmatizer.lemmatize(tok).lower()

clean_tokens.append(clean_tok).strip()

return clean_tokens

Pipelines:

```
pipeline = Pipeline([
    ('vect', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
    ('clf', Random Forest Classifier())
])
```

} transform
} classifier

```
pipeline.fit(X_train, y_train)
```

```
predicted = pipeline.predict(X_test)
```

} predict (and if last
is classifier).

Estimator - any object that learns from data.

- e.g. classification, regression, clustering, transforms.
- must have fit method

Transformer - transform method to apply transformation model
to new data.

- can include cleaning, reducing, expanding, or
generating features.

Predictor - predict method to predict test data.

Has fit method to train model on training data.

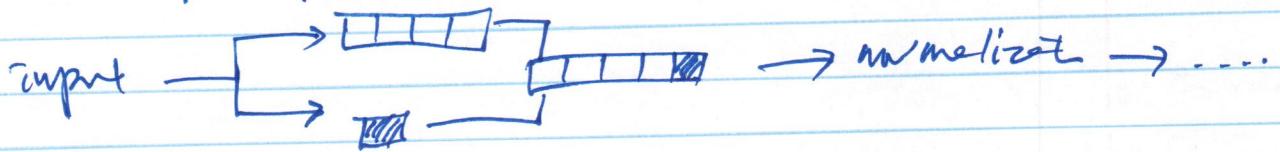
GridSearch - automates process of testing different hyper parameters
to optimize model.

- tune over entire pipeline

Prevent data leakage - contains all steps within each fold of
cross-validation!

Feature Union:

- perform steps in parallel and take union of results for next step.



```
pipeline = Pipeline([
    ('features', FeatureUnion([
        ('nlp-pipeline', Pipeline([
            ('vect', CountVectorizer()),
            ('tfidf', TfidfTransformer())
        ])),
        ('text-len', TextLengthExtractor())
    ])),
    ('clf', RandomForestClassifier())
])
```

```
pipeline.fit(Xtrain, y_train)
predicted = pipeline.predict(Xtest)
```

Custom Transformer:

```
from sklearn.base import BaseEstimator, TransformerMixin
class TenMultiplier(BaseEstimator, TransformerMixin):
    def __init__(self):
        pass
```

```
def fit(self, X, y=None):
    return self
```

```
def transform(self, X):
    return X * 10
```

- Fit method - $X = 2d$ array for feature data
 $y = 1d$ array for target
- transform method return self

- Transform method - does the transformation using X.

multiplier = Ten Multiplier()
 $X = np.array([1, 2, 3, 4])$
multiplier.transform(X)

class CoreNormalizer(BaseEstimator, TransformerMixin):

def fit(self, X, y=None):
return self

def transform(self, X):
return pd.Series(X).apply(lambda x: x.lower()).values

def __init__(self):
pass

core_normalizer = CoreNormalizer()
X = np.array([1, 2, 3, 4])
core_normalizer.transform(X)

Function Transformer:

from sklearn.preprocessing import FunctionTransformer
transformer = FunctionTransformer(np.log1p, validate=True)
X = np.array([1, 2, 3, 4])
transformer.transform(X)

GridSearch :

- using pipeline, run search on entire ~~the~~ workflow
- define grid of parameters, scores each combination with cross-validation, determine parameter that results in optimal model.
- also works on ≥ 2 classifier.

parameters = {
 'kernel': ['linear', 'rbf'],
 'C': [1, 10]
}

svc = SVC()

clf = GridSearchCV(svc, parameters)
clf.fit(X_train, y_train)

- avoid data leakage using pipeline,

pipeline = Pipeline([
 ('scaler', StandardScaler()),
 ('clf', SVC())
])

parameters = {
 'scaler-with-mean': [True, False],
 'clf-kernel': ['linear', 'rbf'],
 'clf-C': [1, 10]
}

Standardization happens when grid search starts. For each cross-validation, scaling is done only on data that is split and to be trained on, preventing leakage.

```
def main():
    X, y = load_data()
    X_train, X_test, y_train, y_test = train_test_split(X, y)

    model1 = Logistic_model()
    model1.fit(X_train, y_train)
    y_pred = model1.predict(X_test)

    display_results(model1, y_test, y_pred)
```

```
def Logistic_model():
    pipeline = Pipeline([
        ('features', FeatureUnion([
            ('', ''),
            ('', ''),
            ('', ''),
            ('', '')
        ]))
    ])
```

```
parameters = { }
```

```
cv = GridSearchCV(pipeline, param_grid=parameters)
```

```
return cv
```

```
def display_results(cv, y_test, y_pred):
    labels = np.unique(y_pred)
    confusion_matrix = confusion_matrix(y_test, y_pred, labels)
    accuracy = (y_pred == y_test).mean()
```

5. Experiment Design + Recommendation

Lesson 2: Experimental Design Concepts

Correlation does not imply Causation

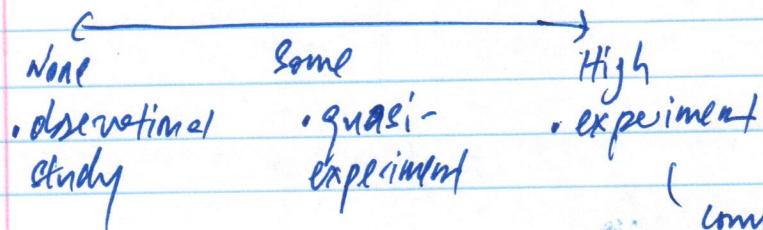
1. what is experiment?
2. what are types of experiments?
3. how to measure outcome?
4. common pitfalls of experiment design?

Key feature of experiment

1. Comparison between groups
2. Control for other variables.
 ↳ randomization

1. Types of Experiment.

Control over variables



compare outcomes between ≥ 2 groups

lack control

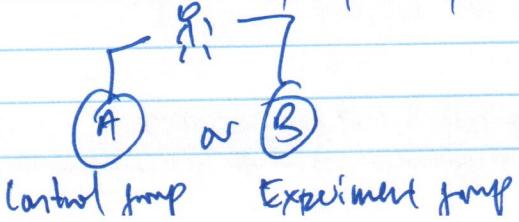
due to ethical concerns

directionality of cause might be uncertain

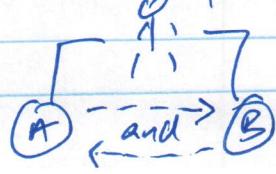
A/B test - compare performance between group A and B.

Between-subjects design:

1 condition per participant.



Within-subjects design:
many conditions per participant



Factorial designs.

- control, X, Y, X and Y

2. Types of Sampling

{ - Simple random sampling

all population have equal chance of selection

{ - Stratified random sampling

divide population into distinct subgroups before sampling.

proportional sampling - maintain distribution

nonproportional sampling - fix number sample from each group.

probabilistic sampling - sample representation of whole population.

Convenience Sample - limited to group of people

e.g. college students, elderly, military.

3. Measuring outcomes

- objective features to evaluate performance - evaluate metrics
- goal of study must be separated from enablers metrics
 - metric is not main point of study
 - metric implies of metric relative to goal
 - why want to run experiment.

- metrics - key numbers (KR)

key performance indicators (KPI)

operational definition

definition - constraint

User funnel - path / steps taken by user from start of experience to desired destination state.

- drop off in user as move from step to step. Mitigate experiments to improve drafts.

Unit of division - observation unit by which users are split into experimental groups.

1. Event-based division

- \ experiment condition randomly rolled
- \ ignores previous visits

2. Cookie-based division

- \ stored on user device
- \ determines experimental condition
- \ unobtrusive by deleting cookie / anonymous browsing.

3. Account-based division

- \ randomly divides user IDs into condition
- \ requires accounts + logged in.
- \ limit pool of data, risks personal identifiable information.

Consider consistency of experience when selecting dimensions.

Evaluation Metrics

- Features that provide an objective measure of outcome of an experimental manipulation.
- e.g. proportion of visitors made purchases, average amount purchased.

Invariant metrics

- Features that objectively check for equivalence between experimental groups, makes sure experiment is running as expected, no biases!
- e.g. number of cookies, logged in users.

4. Controlling Variables

- Need to have 1. one variable feature of interest.
- 2. none other features change.

- all
- if unable to control other features, then susceptible to confounding variables.
 - outcome observed may be due to 3rd variable, neither than feature of interest.
 - \ causal relationship between 3rd variable/feature of interest,
indirect relationship mediated by 3rd, intermediate variable
 - diversion of effect - A to B
B to A
A, B to C

5. Experiment validity - degree which experiment actually accomplishes conclusion of study.

1. Construct Validity

- objective metric worth support goals of study.

2. Internal Validity

- claims of causality can be supported.

3. External Validity

- results can be generalized to other article experiment.

6. Bias

- systematic error that affect interpretation of results.

• Sampling bias - not representative of population.
- e.g. self-selection bias, survivor bias.

• Novelty bias - change in participant behavior solely due to novel stimuli.

- order bias - condition order, whether occur first (primacy bias) or last (recency bias).
 - avoid by randomizing order of conditions.
- experimenter bias - presence of experimenter introduces bias from subject, or from experimenter him/herself.
 - blinding participants and experimenters from condition to help avoid subconscious bias.

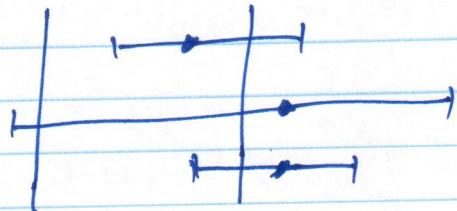
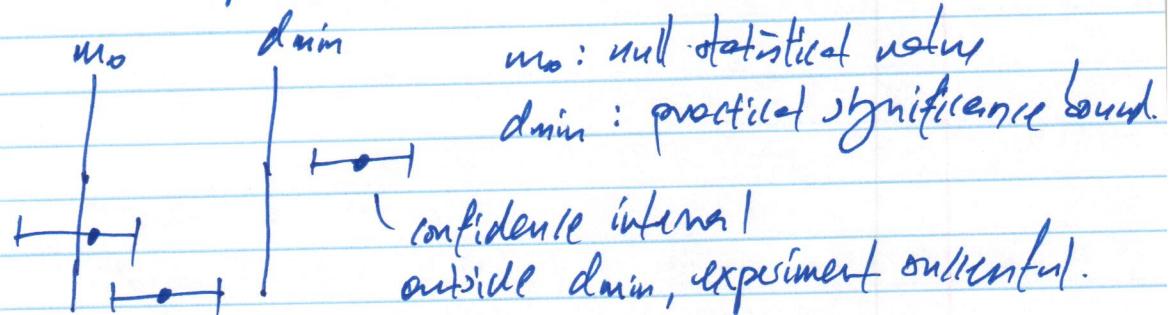
7. Ethics

- minimize participant risk. - physical, emotional.
- clear benefits for risks taken.
- provide informed consent
 - ↳ if deception is needed, must debrief at end of experiment.
 - ↳ included in terms of agreement.
- handle sensitive data appropriately.
 - ↳ personal identifiers must be anonymized.
 - ↳ survey + census results aggregated to avoid tracing outcomes back to individual person.
- Institutional Review Board (IRB) to approve proposals.

S - specific - goals are specific.
 M - measurable - outcomes measurable
 A - achievable - steps + goals realistic
 R - relevant - purpose
 T - timely - ~~reasonable~~ reasonable time frame

Lesson 3: Statistical Considerations in Testing

Practical significance - level of effect used to observe for experiment to be a true outlier.



- uncertainty or confidence interval straddles d_{\min} bounds.
- option to collect more data

statistical power - given true mean, probability of rejecting null hypothesis.

Hypothesis Testing

null hypothesis : $H_0: \mu = 191$
 research hypothesis : $H_1: \mu > 191$

$n = 100$
 $\bar{X} = 197.1$
 $\text{std dev} = 25.6$

} does new sample data support null or research hypothesis?

probability using Central Limit Theorem

$$P(\bar{X} > 197.1) = P\left(z > \frac{\bar{X} - \mu_0}{\text{std dev}}\right) = P(z > 2.38)$$

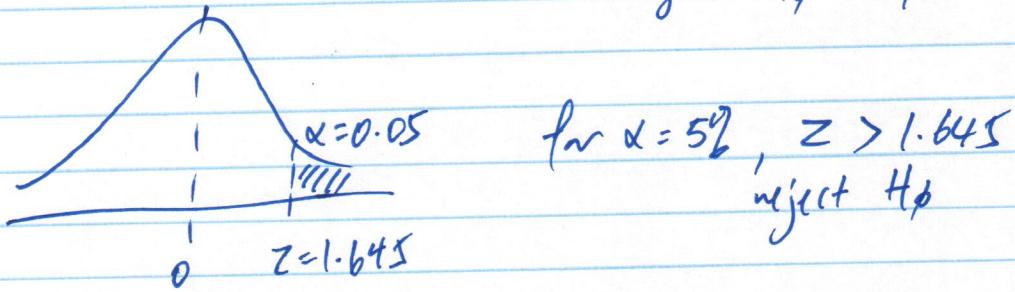
$$= 1 - 0.9913$$

$$= 0.0087$$

$$\left. \begin{array}{l} n = 100 \\ \bar{x} = 192.1 \\ s_d = 25.6 \end{array} \right\} P(X > 192.1) = P\left(z > \frac{\bar{x} - \mu_0}{\frac{s_d}{\sqrt{n}}}\right) = P(z > 0.43) = 1 - 0.6614 = 0.3336$$

33.36% probability
sample means > 192.1 .

$$\alpha - \text{level of significance} = P(\text{Type I error}) = P(\text{Reject } H_0 \mid H_0 \text{ is true})$$



Step 1: Test hypothesis.

1. Upper-tailed test

$$- H_1 : \mu > \mu_0$$

2. Lower-tailed test

$$- H_1 : \mu < \mu_0$$

3. 2-tailed test

$$- H_1 : \mu \neq \mu_0$$

Step 2: Test statistic

$$z = \frac{\bar{x} - \mu_0}{\frac{s_d}{\sqrt{n}}}$$

Step 3: Set up Decision Rule.

1. type of test

upper - reject H_0 if $z \geq z_\alpha$

lower - " " " " $z \leq z_\alpha$

2-tailed - " " " " $z_\alpha \leq z \leq z_\alpha$

2. type of test statistic
- z standard norm distribution
 - t distribute.

3. α - level of significance.

Step 4: Compute test statistic

$$H_0: p_1 - p_2 = 0$$

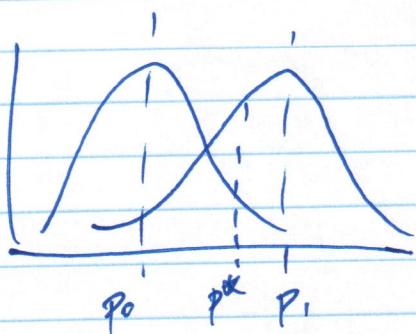
$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

where $\hat{p} = \frac{Y_1 + Y_2}{n_1 + n_2}$

sample proportion

standard error $= \sqrt{\frac{p_0(1-p_0)}{n}}$ $\approx \hat{p} = \frac{Y}{n}$

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{\hat{p} - p_0}{SE_0}$$



Number of samples

$$p^* - p_p = Z_{1-\alpha} SE_0$$

$$p_1 - p^* = -Z_\beta SE_1$$

$$p_1 - p_p = Z_{1-\alpha} SE_0 - Z_\beta SE_1$$

$$p_1 - p_p = Z_{1-\alpha} \sqrt{p_0(1-p_0)} - Z_\beta \sqrt{p_1(1-p_1)} \quad \text{--- (1)}$$

$$n = \left(\frac{Z_{1-\alpha} \sqrt{p_0(1-p_0)} - Z_\beta \sqrt{p_1(1-p_1)}}{p_1 - p_p} \right)^2$$

Dummy test - experiment on against equivalent groups

- A-A test

- provide distributional data if historical data not available.

- test for any errors

replicum time + resources.

• quantile_ci

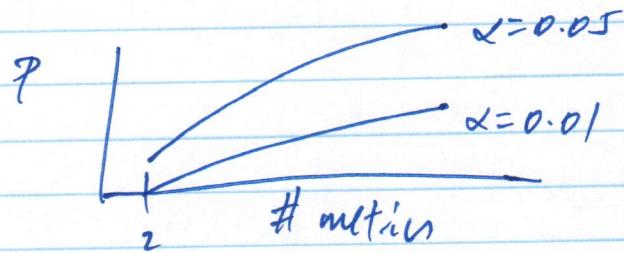
• rank-sum test (mann-whitney)

 `scipy.stats.mannwhitneyu(x, y, alternative='greater')

• sign test

- Probability of making at least one Type I error

$$\Rightarrow \alpha_{\text{over}} = 1 - (1 - \alpha_{\text{ind}})^n$$



introducing comlti factor, $\alpha_{\text{ind}} = \frac{\alpha_{\text{over}}}{n}$

Bonferroni correction

Sidak correction, $\alpha_{\text{ind}} = t(1 - \alpha_{\text{over}})^{1/n}$

$$p = 0.5$$

$$\text{std dev} = \sqrt{p * (1-p) * \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

$$z = \frac{(n_1 + 0.5) - (p * n)}{\text{std dev}} \quad \text{or} \quad \frac{p_1 - p_2}{SE}$$

$$p = 1 - \text{stats.norm.cdf}(z)$$

Lesson 4: A/B Testing Case Study

Goal: Revise structure of homepage will increase number of people download software, eventually number of people that purchase a license.

Funnel:

1. homepage
2. click with addit. information, gallery
3. page for user to download software
4. " " " purchase license
5. support page with documentation + FAQ for software.

- user account required.
- program free to use for 7 days trial.
- end of 7 days, direct user to a dialog box to license pg.
- purchase code, receive unique code to register account with program.

A/B testing - control = group A
- experiment = group B

Unit of division

- event-based - good for when changes are not user visible.
- account-based - requires user accounts created.
- cookie-based - assign on first page view.
 - track whether download page is reached or not.
 - account is reached or not.

Metric - number of cookies on homepage

- download page

- account creation page

number of purchases - user account database

invariant metric - number visitors @ home page
evaluative metric - ratio # downloads / # visitors
- ratio # licenses / # visitors

Lesson 6: Introduction to Recommendation Engines

Recommendation types:

1. Knowledge-based
2. Collaborative filtering-based — model-based collaborative filtering
3. Content-based
 ↳ neighborhood based
 /

Similarity Metrics:

1. Pearson's correlation coefficient
2. Spearman's " "
3. Kendall's Tau
4. Euclidean Distance
5. Manhattan Distance

Business cases for recommendations:

1. Relevance
2. Novelty
3. Serendipity
4. Influenced diversity

Rank-based recommendation

- orders based on metric - cost, highest ratings, most popular

• Knowledge-based

- user provides preferences or filter to make recommendation.

• Collaborative filtering - recommendations based on interaction between user + item.

• Content-based - use information about users or items to assist in recommendations.

Collaborative Filtering:

1. Model-based

 ↳ ML modeling

2. Neighborhood-based

 ↳ Similarity-Based methods

- Pearson's correlation
- Spearman's correlation
- Kendall's tau

most dissimilar

range $-1 \rightarrow +1$

\downarrow
most similar

 ↳ Distance Based methods

- Euclidean distance
- Manhattan distance

most dissimilar

range $0 \rightarrow \infty$

\downarrow
most similar

$$\text{Pearson's, } \text{CORR}(x, y) = \frac{\sum_i^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i^N (x_i - \bar{x})^2 (y_i - \bar{y})^2}}$$

$$\text{Spearman's, } \text{SCORR}(x, y) = \frac{\sum_i^N (x_i^r - \bar{x}^r)(y_i^r - \bar{y}^r)}{\sqrt{\sum_i^N (x_i^r - \bar{x}^r)^2 (y_i^r - \bar{y}^r)^2}}$$

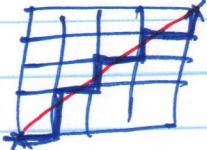
$x \rightarrow x'$, ranked
 $y \rightarrow y'$, ranked

$$\text{Kendall's Tau, } \text{tau} = \frac{2}{n(n-1)} \sum_{i < j} \text{sign}(x_i^r - x_j^r) \text{ sign}(y_i^r - y_j^r)$$

$$\text{sign}(x_i^r - x_j^r) = \begin{cases} -1, & x_i^r < x_j^r \\ 0, & x_i^r = x_j^r \\ 1, & x_i^r > x_j^r \end{cases}$$

$$\text{euclidean distance, } \text{EUC}(x, y) = \sqrt{\sum_{i=1}^N (x_i - y_i)^2}$$

$$\text{manhattan distance, manhat}(x, y) = \sum_{i=1}^N |x_i - y_i|$$



euclidean

Movie recommendation example :

1. Data cleaning
 2. Rank by mean ratings
 - count ratings
 - data ratings
 3. Filter by year
 - genre
- } knowledge-based

1. User-item matrix

	movieid1	2	3
user1	rating	" "	" "
2	" "	" "	:
3	" "	" "	" "

} user-by-movie

} movies watched (user-id)

• movies seen } dictionary
of each user (key),
list of movies (value).

• movie-to-analyze } dictionary of users
to >2 movies.

• Compute_euclidean distance between user 1, 2

movie 1 = movie_to_analyze [user1]

movie 2 = " " [user2]

same_movie = np.intersect1d (movie1, movie2, assume_unique=True)

$df = \text{user_by_movie.loc}[(\text{user1}, \text{user2}), \text{same_movies}]$
 $dnt = \text{np.linalg.norm}(df.\text{loc}[\text{user1}] - df.\text{loc}[\text{user2}])$

- $\text{find_closest_neighbors(user)}$
ranked by euclidean distance,
- $\text{movies_liked}(\text{user_id}, \text{min_rating} = 7)$
 $\text{user_items}[(\text{user_items['user_id']} == \text{user_id}) \& (\text{user_items['rating']} >= \text{min_rating})][\text{'movie_id'}]$
- $\text{movies_names(movie_ids)}$

$\text{make_recommendation(user, num_rec = 10)}$

$\text{movies_seen} = \text{movies_watched(user)}$

$\text{closest_neighbors} = \text{find_closest_neighbors(user)}$

$\text{neocom} = \text{np.array}([])$

for neighbor in closest_neighbors:

$\text{neighbor_liked} = \text{movies_liked(neighbor)}$

$\text{new_neocom} = \text{np.setdiff1d}(\text{neighbor_liked},$

movies_seen

$\text{assume_unique} = \text{True}$)

$\text{neocomm} = \text{np.unique}(\text{np.concatenate}([$

$\text{neocom}, \text{new_neocom}],$

$\text{axis} = 0))$

if ($\text{len(neocomm)} > \text{num_rec} - 1$):
break

$\text{movie_names(neocomm)}$

Content-based recommendation:

matrix of similar movies by dot product

not similar

↑
range from 0 to 1

↓ more similar

movies [electrino, ... , gmae, ...]
dot_product = movies . dot (np.transpose (movie))

	movie_1	2	3	...
movie_1	0	1	1	
2	1	2	0	
3	1	7	- - -	

for user in users:

{ reviews_temp = ranked_reviews [ranked_reviews ['user_id'] == user]
movie_id_temp = np.array (reviews_temp ['movie_id'])
movie_temp = np.array (get_movie_names (movie_id_temp))
movie_list = list (movies [movies ['movie_id'].isin (movie_id)] ['movie'])
movie_idx = np.where (movie_idx ['movie_id'] == movie_id) [0] [0]

recs = defaultdict (set)

for movie_id in movie_id_temp:

rec_movies = find_similar_movie (movie_id)

temp_recs = np.setdiff1d (rec_movies, movie_temp)
users [user].update (temp_recs)

movie_idx = np.where (movie_idx ['movie_id'] == movie_id) [0] [0]

similar_idx = np.where (dot_product [movie_idx]

= np.max (dot_product [movie_idx])) [0]

similar_movie = np.array (movies . iloc [similar_idx,] ['movie'])

Ratings

- interactions - what question to ask?

- how many to ask?

- what type of scale to use?

 \ 0 to 1 } binary, no mid-way/neutral

 \ 1 to 4 }

 \ 1 to 7 } neutral ratings.

 \ 1 to 3 }

Goals of recommendation system

1. Relevance - relevant to users

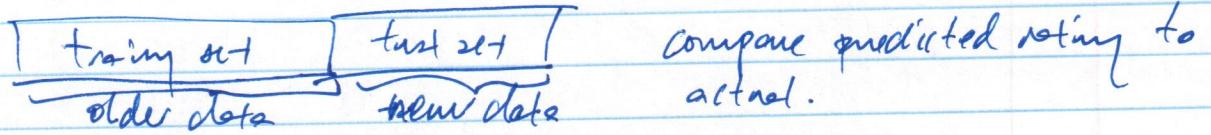
2. Novelty - something user is not already exposed.

3. Serendipity - something user unknown to, different, but likes

4. Diversity - diverse things to recommend, rather than a bunch of similar ones.

Lesson 7: Matrix Factorization for Recommendation:

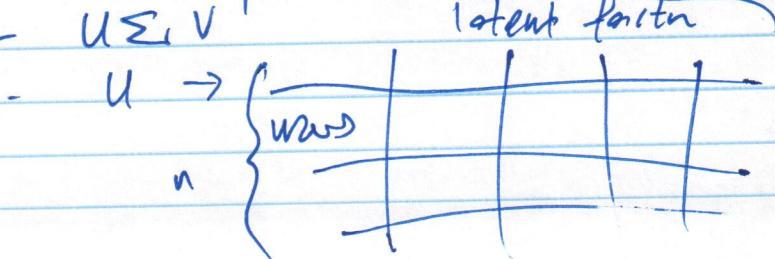
Test recommendation engine before going live. - offline testing.



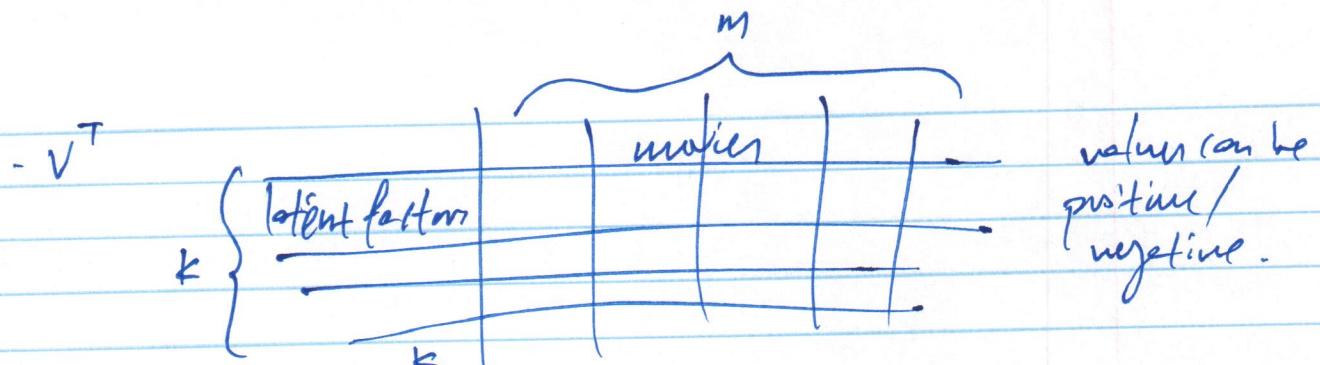
Latent factors - not observed in data, but inferred based on relationships that occur.

a. use singular value decomposition (SVD) to find latent factors common to movies + customers.

$$\text{SVD} - U \Sigma V^T$$



was - latent factors matrix



$- \Sigma$ $\left\{ \begin{array}{c} \text{diagonal matrix} \\ k \end{array} \right.$

- always positive, sorted largest to smallest.
- diagonal matrix

$- U \Sigma V^T = \text{user-item matrix}$

$U, S, V^T = \text{np.linalg.svd}(\text{user-movie})$

(will not work if any NaN present.)

total variability of data = $\sum s_i^2$

variability of 1st n components = $\sum_i^n s_i^2$

% variability of 1st n = $\frac{\sum_i^n s_i^2}{\sum_i^m s_i^2} \times 100\%$

→ reduce # of latent factors by their % variability

$$\text{predicted} = U \cdot S \cdot V^T$$

$$\text{_squared error} = \sum_i \sum_j (\text{predicted} - \text{actual})^2$$

latent factors - are not labeled

- ~~only~~ know their weights

- solve SVD analytically treating it as a closed loop solution.
- But only works without missing data.

44 2 Solve SVD using gradient descent method to handle missing data

↳ FunkSVD

$$UV^T = A$$

- Gradient Descent

$$\text{User matrix : } \frac{\partial}{\partial u_i} (y - uv)^2 = 2(y - uv)v_i$$

$$\text{Movie matrix : } \frac{\partial}{\partial v_i} (y - uv)^2 = 2(y - uv)u_i$$

$$\text{update : } u_i^{\text{new}} = u_i^{\text{current}} + 2\alpha (y - uv) v_i^{\text{current}}$$

- Only apply to non-NaN values in U and V

def FunkSVD(ratings_mat, latent_features=4, learning_rate=0.0001, iters=100):

$$n_{users} = \text{ratings_mat.shape}[0]$$

$$n_{movies} = \text{...} \quad [1]$$

$$\text{num_ratings} = \text{np.count_nonzero}(\text{np.isnan}(\text{ratings_mat}))$$

$$\text{user_mat} = \text{np.random.rand}(n_{users}, \text{latent_features})$$

$$\text{movie_mat} = \text{...} \quad (\text{latent_features}, n_{movies})$$

$$\text{sse_accum} = 0$$

for iter in iters:

 for i in range(n_users):

 for j in range(n_movies):

 if ratings_mat[i, j] > 0:

$$\text{diff} = \text{ratings_mat}[i, j] -$$

$$\text{np.dot}(\text{user_mat}[i, :], \text{movie_mat}[:, j])$$

$$\text{sse_accum} += \text{diff} ** 2$$

 for k in range(latent_features):

$$\text{user_mat}[i, k] += \text{learning_rate} * 2 * \text{diff}$$

$$* \text{movie_mat}[k, j]$$

$$\text{movie_mat}[k, j] += \text{learning_rate} * 2 * \text{diff}$$

$$* \text{user_mat}[i, k]$$

```
def main():
    import recommender
    import recommender as r
    rec = r.Recommender()
    rec.fit(reviews_path = "train-data.csv",
            movies_path = "movies-clean.csv",
            learning_rate = 0.01,
            iters = 1)
```

```
rec.get_id_rating(user_id=8,
                  movie_id=2844)
```

```
rec.make_recommendation(8, 'user')
" " " " (1, 'user')
" " " " (1853728)
" " " " (1)
```

```
class Recommender():
```

```
    def __init__(self):
```

```
        def fit(self, reviews_path, movies_path, learning_rate,
                latent_features=12, iters=100):
```

```
            self.reviews = pd.read_csv(reviews_path)
            self.movies = pd.read_csv(movies_path)
```

```
            user_item = self.reviews[['user_id', 'movie_id',
                                      'rating', 'timestamp']]
```

```
            self.user_item_df = user_item.groupby(['user_id',
                                                   'movie_id'])['rating']
                                         .max().unstack()
```

```
            self.user_item_mat = np.array(self.user_item_df)
```

```
            self.latent_features = latent_features
```

```
            self.learning_rate = learning_rate
```

```
            self.iters = iters
```

`self.n_users = self.user_item_mat.shape[0]`
`self.n_movies = " " " [1]`

`self.num_ratings = np.count_nonzero(~np.isnan(self.user_item_mat))`
`self.user_ids_series = np.array(self.user_item_df.index)`
`self.movie_ids_series = np.array(self.user_item_df.columns)`

`user_mat = np.random.rand(self.n_users, self.latent_features)`
`movie_mat = np.random.rand(self.latent_features, self.n_movies)`

`sse_accum = 0`

`for iteration in range(self.epochs):`

`sse_accum = 0`

`for i in range(self.n_users):`

`for j in range(self.n_movies):`

`if self.user_item_mat[i, j] > 0:`

`diff = self.user_item_mat[i, j] -`
`np.dot(user_mat[:, i], movie_mat[:, j])`

`sse_accum += diff ** 2`

`for k in range(self.latent_features):`

`user_mat[i, k] += self.learning_rate * 2 * diff *`
`movie_mat[k, j]`

`movie_mat[k, j] += self.learning_rate * 2 * diff *`
`user_mat[i, k]`

`print("%d %f" % (iteration + 1), sse_accum / self.num_ratings)`

```
self.user_mat = user_mat  
self.movie_mat = movie_mat
```

```
self.ranked_movies = rf.ordered_ranked_df(self.movies,  
                                         self.reviews')
```

```
def predict_rating(self, user_id, movie_id):
```

```
try:
```

```
    user_row = np.where(self.user_ids_series == user_id)[0][0]  
    movie_col = " ".join(["movie_ids_series == movie_id"])
```

```
    pred = np.dot(self.user_mat[user_row, :],  
                  self.movie_mat[:, movie_col])
```

```
    movie_name = self.movies[self.movies['movie_id'] == movie_id][  
                           'movie'].to_list()
```

```
    print("User {}, predict rating of {}, for movie {}.".format(  
          user_id, str(pred), str(movie_name[0])))
```

```
return pred
```

```
except:
```

```
    print("User-movie pair not found")  
    return None
```

```
def make_recommendations(self, -id, -id-type = 'movie',  
                        rec_num = 5):
```

rec_ids = None
rec_names = None

if -id-type = 'user':

if -id in self.user_ids_series:

idx = np.where(self.user_ids_series == -id)[0][0]

preds = np.dot(self.user_mat[idx, :], self.movie_mat)

indices = preds.argsort()[-rec_num][:-1]

rec_ids = self.movie_ids_series[indices]

rec_names = rf.get_movie_name(rec_ids, self.movies)

flip order

else:

rec_names = rf.popular_recommendations(-id, rec_num,
 self.ranked_movies)

print("Offering top recommendation, since user is not in DL")

elif -id-type = 'movie':

if -id in self.movie_ids_series:

rec_names = list(rf.find_similar_movies(-id,
 self.movies))[:rec_num]

else:

print("Movie does not exist in DL")

return rec_ids, rec_names

recommender_function.py

```
def get_movie_names(movie_ids, movies_df):
```

```
    movie_list = list(movies_df[movies_df['movie_id'].isin(movie_ids)][['movie']])
```

```
    return movie_list
```

```
def create_ranked_df(movie, reviews):
```

```
    movie_ratings = reviews.groupby('movie_id')[['rating']]
```

```
    avg_ratings = movie_ratings.mean()
```

```
    num_ratings = movie_ratings.count()
```

```
    last_ratings = reviews.groupby('movie_id').max()[['date']]
```

~~```
 avg_ratings_df = pd.DataFrame({avg_ratings}); avg_ratings
```~~~~```
    num_ratings_df = " " (num_ratings); num_ratings
```~~~~```
 last_ratings_df = " " (last_ratings); last_ratings
```~~

```
 ratings_df = pd.merge(avg_ratings_df, num_ratings_df,
```

```
 on='movie_id', how='outer')
```

```
 ratings_df = " " (ratings_df, last_ratings_df,
```

```
 on='movie_id', how='outer')
```

```
 ranked_movies = pd.merge(movie, ratings_df,
```

```
 on='movie_id', how='right')
```

```
 ranked_movies['avg_ratings'], 'num_ratings',
```

```
 'last_ratings', ascending=False, replace=True)
```

```
 ranked_movies = ranked_movies[ranked_movies['num_ratings'] > 4]
```

```
 return ranked_movies
```

↗ Find similarity based on dot product of <sup>ratings</sup>  
content. <sup>highest</sup>

def find\_similar\_movies(movie\_id, movie\_df):

movie\_content = np.array(movie\_df.iloc[:, 4:])

dot\_product = movie\_content.dot(np.transpose(movie\_content))

movie\_idx = np.where(movie\_df['movie\_id'] == movie\_id)[0][0]

similar\_idxs = np.where(dot\_product[movie\_idx] ==

= np.max(dot\_product[movie\_idx]))[0]

similar\_movies = np.array(movie\_df.loc[similar\_idxs, 'movie'])

return similar\_movies

def popular\_recommendations(movie\_id, n\_top, ranked\_movies):

top\_movies = list(ranked\_movies['movie'][n\_top:])

## Recommendation Systems

### 1. Knowledge-based

↳ filter by known parameters - year, genre ...

### 2. Collaborative-filtering

↳ model based

↳ neighbor based - similarity } based on user ratings /  
- distance content values

### 3. Content-based

↳ similar content - dot product of content values / ratings.

↳ latent factors - SVD  $A = U \Sigma V^T$

↳ Frank SVD  $A = U V^T$

↳ recommender validation - split training, validate

↳ cold start - new users do not have ratings, missing information