

Avoid overfitting

- Increase size, I quality of your data
- Data augmentation
- Early stopping
- Regularization
- Dropout
- Decrease model complexity
- Feature selection
- Ensemble
- Cross Validation
- Transfer learning

Gradient accumulation

- with gradient grad-zeno

Avoid underfitting

- Increase model complexity
- increase amount of training data
- Reduce ensembling
- Increase # of epochs
- Use regular (reduce regular strength)
- Cross validation
- Ensemble learning

(overfit)

high variance - fits training data very closely, but it may not generalise well to unseen data

→ L₁ and L₂ regularization

L₁ regularization penalizes sum of absolute values of parameters. Model encourages to have a smaller number of large weights, rather than ~~large # of~~ small # of smaller weights. This leads to more sparse model where many weights are zero.

L₁ - large # of features (other's weight & regularization selected)
(L₂ - small #) $\{$ to increase model stability

L₂ regularization

penalizes the sum of squares of model's weights. Model is encouraged to have a large # of smaller weights, not dense model with all non-zero weights.

$$\text{Err}(x) = \text{Bias}^2 + \text{Variance} + \underbrace{\text{Irreducible error}}$$

↳ measure of ~~size~~ amount
of noise in data

Explain. Meth.

Accuracy - measures how often classifier correctly predicts

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

→ predicted actual class

std dev values, median value across all cells
const - immutable

Uniq - keys can include null

Precision - how much of positive predicted are actually positive

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

Recall - how much of actual the cases we were able to predict correctly

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1 Score

$$F1 = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

ROC - the rate of false positive

2) Stemming Stemming - multiple words same root
brings prefixes & suffixes from word.
to root word doesn't have to be word

Lemmatization

Is aim to identify actual word from
dictionary



POS (part of speech) tag

bots and confusion b/w two
same words that have diff. meaning



TF * IDF

Doesn't consider context -
→ give more weightage to
rare words



how often a
word appear in
specific doc.

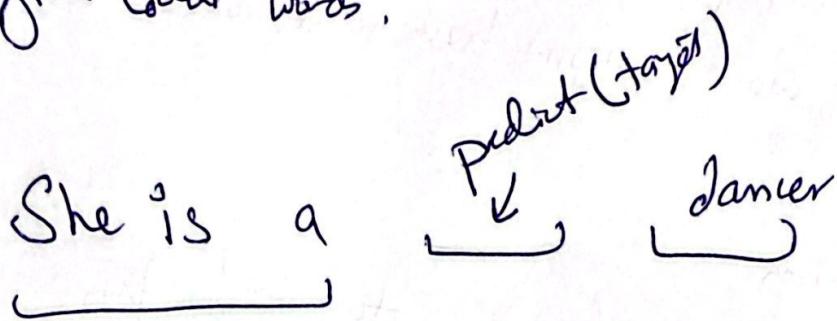
how come
word is in
entire doc
collection

Out-of-Vocab (OOV) - word not in vocab

CBOW (NN)

NLTK / SPACY

↳ NN based algo that predicts a target word given context words.



↳ Focuses on many of words for its similarity.

Skip-gram (NN)

↳ predicts surrounding words given a single target word. It focuses on underlying meaning words related to the word. It predicts words using distributed representation of input word.

Glove Embeds (not NN)

↳ Focuses on co-occurrence

↳ co-occur means syntactically similar

Chunking / Shallow Parsing

↳ split sentence into chunks

In Bow to score word

vector = root length

Bow score word

- count
- frequency (count/total word in doc)

N-GRAM

↳ contains square of n-items

↳ n-gram captures the context information & relationships

↳ sliding = widow of n words across a text corpus

↳ Analyze frequency of certain word sequences.

TOKENIZATION ~~SLIDES~~ HOPPY

→ There can be overlapping in N-gram

→ WINDOW will be SLIDE

ONE word at a time.

→ One-hot-encoding = is applied on category

→ BoW

one-hot encoding = category dict

→ BoW (we look at histogram of total words within the text. Considering word count as feature)

Similar docs have similar counts

→ Ignores grammar & word order but keeps track of frequency of words.

→ At vocab of known words and a frequency of known words

Feature extraction from = convert encoding words to numbers

Corpus is full Text including repetition

In BoW , next step is to score the words in each doc.

As vocab has 10 words

→ [It was the best of times]
vector would be
[1 1 2 1 1 1 0 0 0 0]
←
words: It was the best of times was
→ If word is encountered unknown,
it would be ignored

Space - Not a word

SQL 2

- `Select DISTINCT (Name) From table
City IN ('Paris', 'London')
ORDER BY Price DESC`
 - wher
 - ↳ between
 - like (%)
 - Fin

'%' - n characters after S.
- `Select * From Customers Where
ID = 1 And (name like 'G%' OR NOT ID = 'F')`
 - with bracket AND will get precedence
- `NOT LIKE, NOT BETWEEN, NOT IN`
 - : size character
 - ': n - characters
- `INSERT INTO table (col1, col2),
VALUES ('x', 'y')
VALUES ('a', 'b')`
- `SELECT * where col-name From table Where
col-name is NULL`
- `UPDATE table
SET col1 = 'a', col2 = 'c'
Where ID = 1.`
 - without where, it will be applied to whole database
- `Delete From table
where ID = 2`
- `Drop Table`

Aggregate fn - Perform operation on set of values
→ SUM, MAX, AVG, COUNT, SUM

→ Select MIN(PRICE) as smallest Price from table

→ Select MIN(PRICE) as smallest Price, Category ID
FROM table

Group by Category ID

IN (SELECT)
 └ sub query

Select * FROM Customer

where Customer_ID IN (SELECT customer_ID FROM orders)

OR

Select Cx FROM Customer C

INNER JOIN ORDERS O ON C.customer_ID = O.customer_ID

UNION to join results of two sets SELECT

Select * FROM t1
UNION

SELECT * FROM t2

Q2)

Where can not be used with aggregate fn

So, HAVING

is
before
aggregate
appears

↳ where
comes
in rows

SELECT *
FROM table
where cond
Group by column
Having count
Order by name

Select count (1) from Customers
Group by city where last_name = "D"
Group by last_name

HAVING COUNT(order_id) > 2

EXISTS

SELECT * from Table
where exists (select
ANY exists (

→ SELECT P from T where P = ANY (select)

STORED Procedure.

CREATE Problem

NAME OR #

Qcby Vard(20)

② Pot $v_{\text{arch}}(r_0)$

As

SELECT . . .

403

→ EXEC NAMENRDN @city = "C", @pt = "P"

Alter take x,

Add column one datatype
DROP COLUMN C1

CONSTRAINT, To limit type of job that goes inside table

SQL 3)

UNID - allow NULL, not NULL

Foreign key - link b/w tabs

Indexes help in retrieving data faster

{ Create view as

CHAR - fixed length

VARCHAR

CHAR - 1 byte

INT - 4 bytes

float - "

Double - 8 bytes

bool - 1 byte

Python

- interpretation lang

↳ doesn't have to compile before run

→ FUNCTIONS/^{class} are first class. They can be assigned to variables, return from other fns and passed into fns.

3. ~~INTER~~ COMPILE means to convert into machine code. Python is not in machine lang code. If you run the code

→ Name space → each obj is created in name's address of object is created. object is name which is assigned to each class in Python.

→ object are func/class

→ Python Decorator

↳ are used to add some designed path to a function without changing its structure

MVC

- model represents data - data not UI or UI
- view handles display
- Recv input, interact with model, update view

MVT

- Controller part is big job care of by framework
- Templates act as a present layer and are easy HTML code that render data. Templates are try HTML files that contain placeholders for dynamic data

break, exit, loop
continue, skip, return
sentinels

pickling: convert to string

Function that return iterable of 1 item are generator

Tremor spent

* if — else —

Package names ~~not~~ ^{can} multiple
modules



Shallow GPS

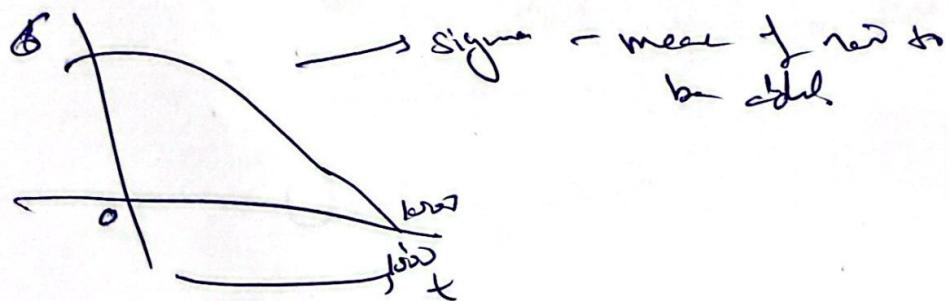
↳ GPS of few points
↳ changes will be small

sol 6

Sol 6)

we added diff lvl of noise at diff steps

↳ create a noisy schedule with
- mostly decaying fm



random pick $\pm \frac{1}{2}$ times of the
odd noise accepted

Δt start $t=0$, pure noise

Momenta:

its change same parallel multiple
times in multiple steps. You should
just increase amount you change the

So we assume we have $\boxed{3}$, we add noise to it and then try to predict that noise so we can get back $\boxed{3}$

$$\boxed{\tilde{f}_3} = \boxed{f_3} + \boxed{\epsilon_{3,y}}$$

$$\text{MSE} = \frac{(y - \tilde{y})^2}{n}$$

more dots

Net predict.

it would be much easier if we can already input a test '3' to get density

Guidance:-

If we feed in '3' {noise} but with syn noise is every except 3.

the i guess

3

we will gradient of
 x_3 is a handwritten digit with
 $\frac{\nabla p(x_3)}{\nabla x_3}$ to pixels of x_3

charge = pixels of x_3 pixels
of x_3 until charge = pixels of x_2

For - can't do it for every pixel, so will charge
pixel values accordingly to gradient

So instead of charging input to weight
of model we are charging input to
model (imagine pixel)

So, we'll take each pixel and subtract
from it a little bit of gradient

$$x_3 - c \times \frac{\nabla p(x_3)}{\nabla x_3}$$

Loss Function

Classification

→ BCE

→ Categorical cross entropy

as sigmoid function
inside of largest element is a vector

→ Softmax!

↳ Normalise while vector

Regression

→ MSE

→ MAE

overfitting (regularization)

→ L1 loss, L2 ridge, Dropout, Early stopping, batch Normalization, gradient penalty, cross validation, Data Augmentation, Ensemble

validation set for hyperparameter tuning.
Test set: model doesn't see directly

Loss shirts less import factor

sparse mesh

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum |\partial|$$

large λ near underfit
low, mean ~~now~~ ^{say} ~~old~~ ^{new} ~~last~~ ^{last}
fig

Ridge

$$+ \lambda \sum |\partial|^2 \rightarrow \text{dense mesh}$$

λ
encourages smaller

weights

SGD/WD \rightarrow to compute gradient of loss function

Gradient accumulates

- ⇒ An epoch means we have passed each sample of training set one time through the network to update the parameters.
 - ⇒ No. of epochs is hyper parameter that defines the times that GD will pass the entire dataset
- In batch gradient, one epoch corresponds to a single iteration
- In stochastic GD, one epoch corresponds to n iterations where n is # of training samples
- In minibatch one epoch corresponds to $\frac{n}{b}$ iterations where b is size of minibatch

Gradient Accumulation

- ↳ does only calculate but doesn't update model parameters in fine-grained accounts
- ↳ Finally use accountable gradient to update parameters.

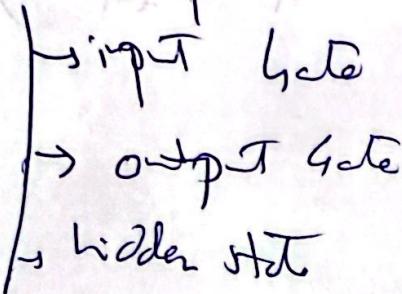
Split batch size into mini batches, so all mini batches use same model rapidly to calculate the gradient, on which is equivalent to using original batch without splitting.

- » zero-grad to reset
- » grad : to compute grad

3)

RNN

RNN component



RNN - to process sequence of data

$$x(t) = x(1), \dots, x(t) \text{ with } t \in \mathbb{N}$$

Time step t ranging from $(1 \rightarrow T)$

→ RNN perform same task for every seq element of sequence, with output being dependent on previous seq computed

weight matrix

→ input matrix

→ hidden weight matrix → captures temporal dependence.

It multiplies the previous hidden state (h_{t-1}) at each time step.

ReLU's bent

- not convex
so derivative is not possible b/c it is not defined there

∴ Gradient Descent won't work

Gradient Descent is optimization technique

Back propagation to calculate gradient for GD

$$\text{Hidden state } (h_t)_2 = f(W_x \cdot x_t + W_h \cdot h_{t-1} + b)$$

$$\text{Hidden } h_t = f(W_x^T x_t + W_h h_{t-1} + b)$$

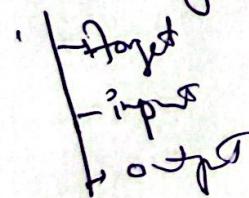
state
↓
activtis f_t

W_x - input weight mat - to ~~convert~~

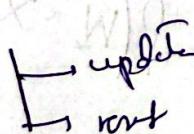
it converts x_t (input mat) into format suitable
for combining with previous previous hidden state

~~$$h_t = f(h_{t-1}, x_t)$$~~

LSTM → three gates structure



h_{t-1}



R-CNN

When can not use standard CNN for detect

b) Fixed Fully convl by fixed

- No info about location,
- multiple obj

image \rightarrow all region proposals \rightarrow Compt CNN \rightarrow classif. fct. \rightarrow class

- \rightarrow classifier loss - & Soft
- \rightarrow obj regre loss

Fast R-CNN

Img \rightarrow CNN map \rightarrow region proposal

Faster R-CNN

instead of select search, CNN gives region proposal

YOLO

- ↳ takes CNN predict box, 9 class
- ↳ image into bounding box
- ↳ width in bounding box



SVM

→ find hyperplane in n-dimension space



in AP - does detection

5)

Severe to Severe

↳ Median to Trans

ONET

↳ cannot create 1 decoder

CLIP - Multimodal

for no grad - disallow competing gradients

Model. eval dead

↳ dropout layers are bypassed

↳ batch normalizations are computed mean

↳ var.

Unit consists of

- contracting path (encoder)
- Expansive path (decoder)

Skip connections
↳ are concatenations

↳ use skip connection to
generate segmentation

Depth allows network to capture high-level features

Each pixel in output image represents a
label that corresponds to a class in output input

loss:

- pixel wise cross entropy
- SLD.

Max pool

- dimensionality reduction
- noise suppression
- translation invariance
- robustness to small transformations

Unit consists of

- contracting path (encoder)
- expansive path (decoder)

skip connections
↳ are concatenations

→ use skip connection to
generate segmentation

Depth allows network to capture high-level features

Each pixel in output image represents a
label that corresponds to a class in ~~input~~ input

loss:

- pixel wise cross entropy
- SLD.

→
→

Max pool

- dimensionality reduction
- noise suppression
- translational invariance
- robustness to small transforms

1) FCNN

↳ deep \Leftrightarrow features can be obtained from going deeper but speed goes down

mean

Shallow layer have

↳ high speed if
↳ low level features

So, we need to we connect them

→ output of conv layer have same depth as filter

→ Filter has same depth as input

Dilated conv has gaps when applied

to input. resulting in layers receptive fields

↳ Gap is introduced in kernel.

1x1 Conv

↳ to reduce # of channels while introduces non-linearity without loss of info

↳ It means a single number (1x1) as opposed to matrix.

This 1x1 pixel will convolve one entire input image pixel-by-pixel.

→ 1x1 filters will only have a single parameter for each channel in input

Depth wise Conv

↳ single channel of the unlike CNN

Point conv

↳ 1x1

10) - loss in ^{CS} how well img is segmented

↳ area of overlap
" " union

DICE as loss for
overlap

DICE - loss for good segments

↳ for undetected obj.

↳ measure of dissimilarity of predicted & the
actual segmt

$$DL(y, \hat{y}) = 1 - \frac{2y\hat{y} + 1}{y + \hat{y} + 1}$$

$$R^2 = \frac{1 - \text{sum of square residuals}}{\text{Total sum of squares}}$$

R^2 - end of regression model
 $R^2 = 1 - \text{Partial } R^2$ - mean
 $R^2 = 0$ - no effect = target value
 total var explained by model
 $\frac{\text{total var explained by model}}{\text{total actual var}}$

actual value vs. target value
 Mean
 bias & variance

Q) Two variables are correlated, no variance of all.

before mba

Actual

	0	1
0	TP	FN
1	FP	TN

Pred

	0	1
0	TP	FP
1	FN	TN

Model predicted +ve, ~~actl~~ is -ve

TP,

FN: model

TN, "

FP,

" -ve, actl is +ve

" -ve, " " -ve

+ve, " -ve

Precision

$$= \frac{\cancel{TP+FN}}{TP}$$

→ same point (to 5)
→ right pred

$$\frac{TP}{TP + FN}$$

→ kth cell to the main diag.
→ how far you fill all the

Recall

$$\frac{TP}{TP + FN}$$

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

$$F_1 = \frac{2 \cdot \text{Prec.} \cdot \text{Recall}}{\text{prec} + \text{recall}}$$

P - are ~~for~~ ^{under} pre-reals

$$\text{Precision} = \frac{2}{5} = 0.6$$
$$\text{Recall} = \frac{2}{5} = 0.4$$

$$\begin{array}{ll} T & 1.0 \\ T & 1.0 \\ R & 0.6 \end{array} \quad \begin{array}{ll} , 0.2 \\ , 0.4 \\ , 0.4 \end{array}$$

sent
Swig
+ tank + signal

BLEU score \rightarrow generic
 \rightarrow evaluate test \rightarrow similarity \rightarrow context test

\hookrightarrow scaj LLM

BLEU - biling

\hookrightarrow make translate test by copy it with
human translate

GLUE - how well been made understand human

day.

People

BLEU \rightarrow monotic
 $\hookrightarrow \phi(x)$
 \downarrow context test for

ARC . A12 Ready chky. - ready rd

Tatty lds - coherent & well studied resps
is answer right

sol

Human & fundilly cont. able

MTC bench - multi turn digits

char of thg, you g. ready step by step

Tree of thg.

React based

char of thg

DEPT - ideally mat

sight paths & file in

LORD right handle
for r/t decomp mtr

R-CNN

• You can not use standard CNN for detect

↳ Fixed Fully conv by fixed

- No info about location,
- + multiple object

image \rightarrow all region
proposals \rightarrow Compute CNN
feats \rightarrow classify

- > classification loss —> Soft
- > bbox regression loss

\longleftrightarrow

Fast R-CNN

Img \rightarrow CNN map \rightarrow region proposals

\hookrightarrow

Faster R-CNN

instead of select search, CNN gives
reject proposals

Yolo

↳ solve CNN predict box, 4 dim

↳ image in SVD grid

↳ width in boundary box

→

SVM → find hyperplane in an n-dimes
space

←

in AP - do detect

Common LMs (1)

completion - output of model

Zero shot - ^{inference} No examples in prompt

One shot - One ↴ ↴ ↴

Few shot - ^{inference} Multiple examples

↳ Small models might need one/few shot inferences

→ Gen config - inference Parameters

↳ max_tokens → max tokens to choose from

↳ top_k → select an output from top k results after applying random weighted strategy using pns

↳ top_p → pns to choose from (confidence pns)

↳ Temperature.