

# PERSONAL NOTES

## PRINCIPLES OF LARGE SCALE ML

Hanfeng Zhai

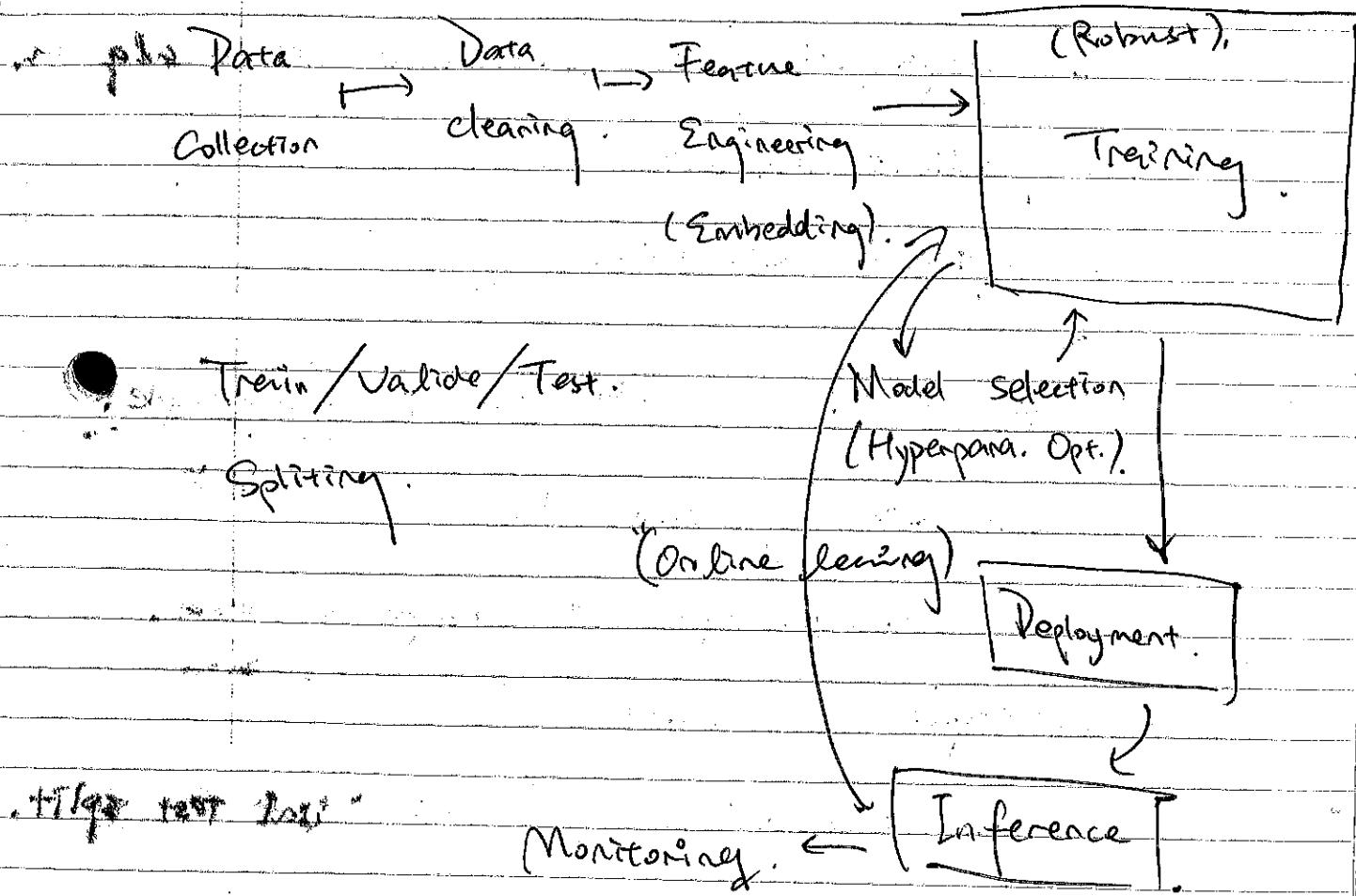
**Disclaimer:** These notes are intended solely for personal reference and study purposes. They represent my own understanding of the course material and may contain errors or inaccuracies. The content presented here should not be considered as an authoritative source, and reliance solely on these materials is not recommended. If you notice any materials that potentially infringe upon the copyright of others, please contact me at [hz253@cornell.edu](mailto:hz253@cornell.edu) so that appropriate action can be taken. Your feedback is greatly appreciated.

2023

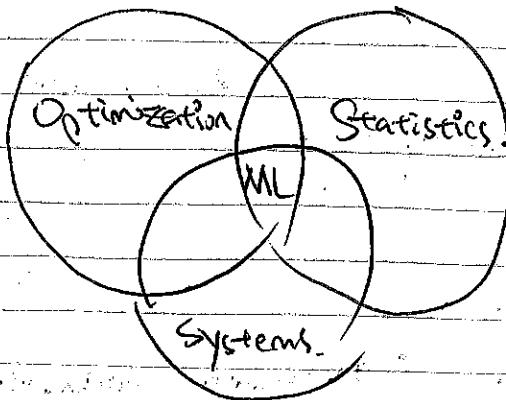
Scaling → ML

↑  
Principles.

How Scale impact performance of ML System.



Learn



### Principle #1: Optimization. Tools

Solve it using gradient-base algorithms

there are fast & "canned"

backpropagation + numerical linear algebra,

- gradient descent.

- empirical risk min.

### Principle #2: ~~X~~ whole dataset, $\rightarrow$ Subsample

- Stochastic gradient d.

$\hookrightarrow$  subsample the component losses.

$\hookrightarrow$  compute faster.

- Cross-validation.

train / val test split.

- bagging.

$\hookrightarrow$  subsample with ~~random~~

- kernel Sampling.

in CNN

- dropout.

- data augmentation.

VAE. (dimension reduction)

Principle #3 Use algorithms there are compatible w/ hardware or vice versa.

- batch size selection.

- numerical precision.

- GPU  $\hookrightarrow$  NN.

- Caching

- data loader / streaming

- Using fast numerical GPU.

- linear algebra kernels.

parallel / distributed computing.

~~https://~~ a neuromorphic computing. (chip  $\rightarrow$  NN).

CS. cornell.edu/courses/6.S087.

- Canvas.

- Gradescope  $\rightarrow$  problem.

- ed discussion.

- CMS - programming a. (Python, numpy, torch)

Week 1 - 2

Aug. 24.

Grading:

Psots

PS.

Prelims

Final

Paper reading.

Review sets - linear algebra

▷ Vector calculus

▷ logistics

▷ computing w/ python

- office hours

Wednesday : 2-3 pm

Gates 426.

# Principle 1 : ML or optimization.

Continuous optimization.

optimizing real numbers.

Data.

→ Embedding steps →  $\mathbb{R}^d$

vectors

e.g. pred. ppl.

33

g. wog.

2013 . current job / current dur.

58,000 / yr.

$$U \rightarrow \mathbb{R}^4$$

Vector = element in vector space

$$x, y \in V \quad xy \in V$$

+ associated & commutative.

multiplication

standard property of a vector space.

Typical example:  $\mathbb{R}^n$ :  $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$

e.g. Matrix  $\mathbb{R}^{n \times n}$ .

$\mathbb{R}^{n \times n \times p}$  can all be interpreted as

vectors.

Core idea under: basis in linear algebra.

how many numbers does it

take to point that "vector space".

A basis for  $V$  is the subset for

$B = \{x_1, x_2, x_3\}$ . s.t. for every  $v \in V - v$  can be written as a linear combination of vectors in  $B$ . exists some  $a_1, a_2, \dots, a_n \in \mathbb{R}$ , and  $a_1 x_1 + a_2 x_2 + \dots + a_n x_n = v$ .

$$\text{s.t. } v = \sum_{i=1}^n a_i x_i \quad V = \text{span}(B)$$

and  $\forall x \in B$ ,

$|B| = \text{"Dimension"}$

terminology confusion

Numerical linear algebra

$$V = \mathbb{R}^n$$

$$B = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$x \notin \text{span}(B \setminus \{x\})$

numpy, sympy, ...

Week 2 (1) Aug. 29.

Automatic diff.

Numpy faster python.

- multi-processing

- C++

- Cache

\* & (2) early nested

$\downarrow$

multiply tensor / matrix multiply

per element

broadcast  $\rightarrow$  faster than manually write.

(1, 3, 7)

(1, 3, 7)

i.e. loop

(3, 1, 7)

(3, 1, 7)

repeat

(3, 1, 5)

(3, 1, 6)

X

(3, 3, 7)

→ numpy broadcasting.

Definition -  $\nabla f(w)^T$  is a linear map.

such that:  $\Delta \in \mathbb{R}^d$

$$\nabla f(w)^T \Delta = \lim_{a \rightarrow 0} \frac{f(w + a\Delta) - f(w)}{a}$$

$$f(x) = x^T A x.$$

$$\lim_{a \rightarrow 0} \frac{(x + a\Delta)^T A (x + a\Delta) - x^T A x}{a}$$

$$\lim_{a \rightarrow 0} \frac{x^T A x + a\Delta^T A x + a x^T \Delta x + a^2 \Delta^T A \Delta - x^T A x}{a}$$

$$\lim_{a \rightarrow 0} \Delta^T A x + x^T A \Delta + a \Delta^T A \Delta.$$

$$\Delta^T A x + x^T A \Delta$$

$$\hookrightarrow \Delta^T (A x + A^T x).$$

$$\nabla f(x) = A x + A^T x$$

all elements

$$f(x) = \|x\|_2^2 = \sum_{i=1}^d x_i^2$$

$$\frac{\partial f}{\partial x_i} = \frac{\partial}{\partial x_i} x_i^2 = 2x_i$$

$$\nabla f(x) = 2x,$$

$$f(x + \alpha \Delta) = \frac{(x + \alpha \Delta)^T}{(x + \alpha \Delta)}$$

Euclidean norm

$L_1$  Norm.

$$f(x) = \|x\|_1 = \sum_{i=1}^d |x_i|$$

$$\frac{\partial f}{\partial x_i} = \frac{\partial}{\partial x_i} |x_i| = \text{sign}(x_i) = \frac{x_i}{|x_i|}$$

$$\nabla f(x) = \text{sign}(x).$$

→ Symbolic differentiation:

write function in mathematical expression.

apply chain rule, ...  
↓  
node

→ Problems.

had to do it manually.

- converting code > no methe. Not efficient
- no guarantee → comp. Structure

Why important?

differentiating  $f$  fraction.

→ Numerical differentiation.

- high-order different.

◦

- error accumulating

## Problems:

- numerical imprecision.
  - number of operations increases.

- if not smooth  $\rightarrow$  wrong results.

- unclear how to setup  $\varepsilon$

- for a vector  $\rightarrow$  scalar function,

You have to compute each partial

individually, meaning  $O(d)$  blowup in cost.

- if  $\varepsilon$  too small, may get zeros

or NaN

## Automatic differentiation:

- forward mode

- reverse mode

Replace  $y$  with a tuple & a derivative

\* operator overloading

## Problems w/ Forward Mode AD:

- Differentiate with respect to

one input

dimension

- Blow-up proportional to  $d$ .



if we are computing a gradient of  $f$ :

$$\mathbb{R}^d \rightarrow \mathbb{R}$$

Week 2 - 2

Aug 31.

## Back propagation.

Review for Forward Mode AD.

$$x \in \mathbb{R}$$



$$\frac{\partial y}{\partial x} \rightarrow \text{store the tuple: } (y, \frac{\partial y}{\partial x}).$$

\* Python supports operator overloading

## Reverse - Mode AD

→ fix one output  $l$  over  $\mathbb{R}$

→ compute partial derivative  $\frac{\partial l}{\partial y}$

$$(y, \nabla_l).$$

Same shape

- Derive backprop. thru chain rule.

$$l = f(u), u = g(y)$$

$$h = P(u_1, u_2, \dots, u_k), \quad u_i = g_i(y)$$

$$\nabla_y h = \sum_{i=1}^k Dg_i(y)^T \nabla_{u_i} h$$

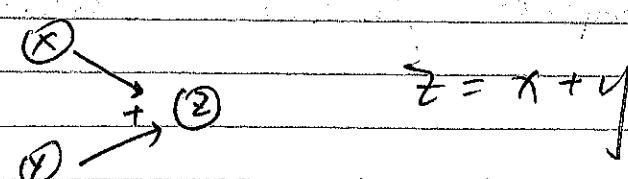
derivative

is a Jacobian

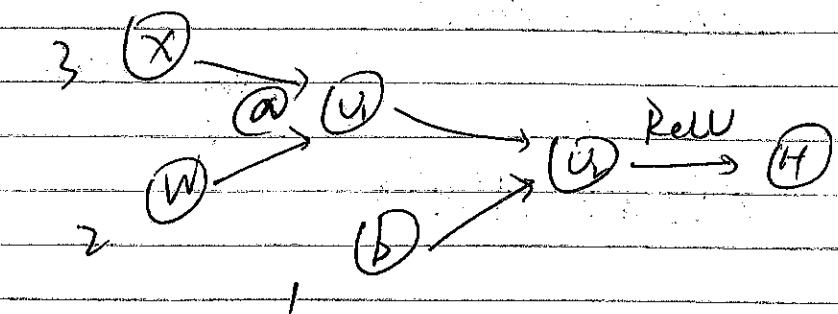
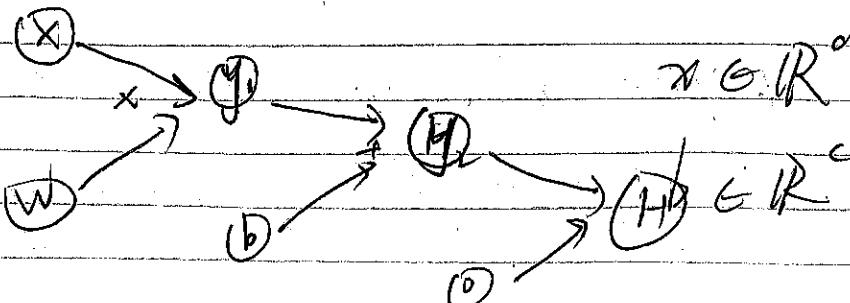
compute gradient w.r.t.  $y$ :

$$\nabla_y h = \sum_{i=1}^k Dg_i(y)^T \nabla_{u_i} h$$

process  $y$                       process  $u_i$



$$H = \max(0, Wx + b). \quad W \in \mathbb{R}^{d \times d}$$



$$\nabla_x H = \frac{\partial H}{\partial U_1} \cdot \frac{\partial U_2}{\partial g} \cdot \frac{\partial g}{\partial V_1} \cdot \frac{\partial V_1}{\partial x}$$

$$H = f(U_1), \quad U_1 = g(U_2), \quad U_2 = h(x).$$

$$\nabla_x H = \text{ReLU}'.$$

interpretation:

$$\nabla_{U_1} H = \nabla_{U_2} H = \nabla_b H = \nabla_h H = 0$$

$$\nabla_{U_2} H = \text{ReLU}'(U_2) \cdot \nabla_h H = 1 \cdot 1 = 1.$$

## Problems

1.2. Derivative of  $x^T A x$ .

b.  $\Delta^T A x + x \Delta A$

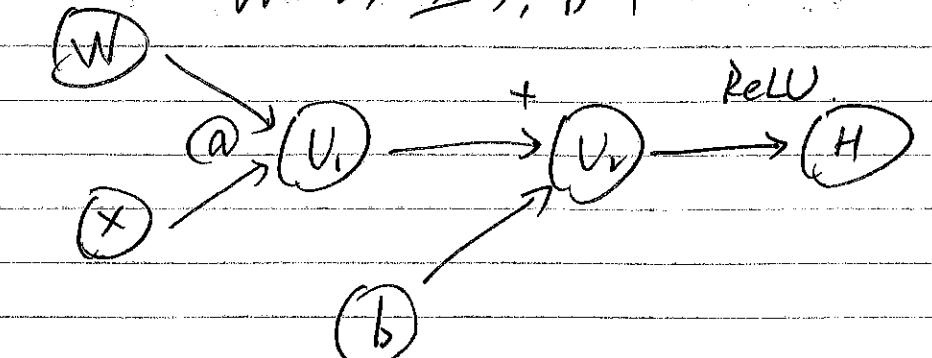
$\downarrow$  why  $x^T A \Delta = (-x^T A \Delta)^T$ ?  
 $\Delta^T A x + \Delta^T A^T x$

d.  $\rightarrow$  a Jacobian.

2.5. Why  $\text{sum}(\log)$  in np.

Review: Backprop.

$W=2, x=3, b=1$



•  $U_1 = W * X = 2 * 3 = 6$ .

Forward •  $U_2 = U_1 + b = 7$ .

•  $H = \text{ReLU}(U_2) = \text{ReLU}(Wx + b) = \text{ReLU}(7)$

•  $\nabla_H H = 1$ .

• Initialize:  $\nabla_{U_2} H = \nabla_{U_1} H = \nabla_{Wx} H = \nabla_x H = \nabla_b H = 0$

•  $\nabla_{U_2} H += \text{ReLU}'(U_2) \nabla_H H = 1 * 1 = 1$ .

•  $\nabla_{U_1} H += \nabla_{U_2} H = 1$ .

•  $\nabla_b H += \nabla_{U_2} H = 1$

•  $\nabla_W H += x \nabla_{U_1} H = 3 * 1 = 3$

labor day  
Week 3 - 1/2.

$$\nabla_x H + = W \nabla_{\theta_i} H = 2 * 1 = 2$$

Machine learning systems

$$f: U \rightarrow V$$

$$DF(x) \in \mathcal{L}(U, V)$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$DF(x) \in \mathcal{L}(\mathbb{R}^m, \mathbb{R}^n).$$

$$\mathbb{R}^{m \times n}$$

$$\text{if } F(x) = x^2$$

$$\lim_{\alpha \rightarrow 0} \frac{F(x+\alpha \delta) - F(x)}{\alpha}$$

$$A \in \mathbb{R}^{n \times n}$$

$$\lim_{\alpha \rightarrow 0} \frac{x^2 + \alpha \delta x + \alpha x \delta + \alpha^2 \delta^2 - x^2}{\alpha}$$

$$= 2x + \alpha \delta$$

$$DF(x) = (\delta \mapsto \delta x + \alpha)$$

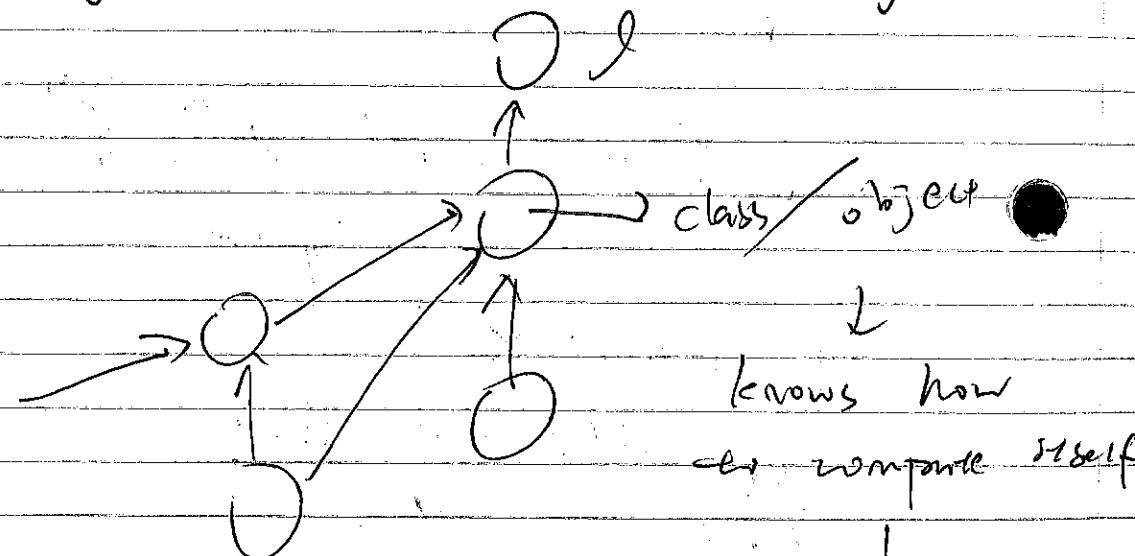
$$\text{or } (DF(x))(\delta) = \delta x + \alpha$$

Why Backprop.

Gradient has + output  $\delta$ .

Backprop manifest gradient

By  $\delta$  for intermediates  $y_j$ .



def \_\_add\_\_(self, other):

- two ways:

(lazy)

SAGG

v.s. GRAPH (only later)  
compute immediately

don't compute, but  
know the dependencies

Eager

I cutoff

I error

, force storage

sidebussing &

development (x),

I If / branch / control flow  
(v)

↓  
convert it  
into some form  
that "in python"

- optimization

immutabilized  
Graph.

o first computation.

o simplification  
(transform. ✓)

o reuse the graph.  
(✓)

o good for caching  
(?).

o memory locality  
(✓).

o less python  
overhead (v)

↓ Deployment.

TI  $\xrightarrow{\text{Default}}$  Graph  
PI  $\longrightarrow$  Eager

## ML framework.

→ numerical linear algebra library

→ hardware accelerator support.

o backprop engine. (gradient. auto.).

o library for writing Deep Notes.

PyTorch & TensorFlow

MxNet  
(C++)

Jax.

Flax (Julia)

Caffe (ancient). → 1st GPU use.

```
model = torch.nn.Sequential(  
    nn.Linear(...),  
    nn.Linear(...),  
)
```

hardware - accelerator.

X. to ("cuda")

X. to ("cpu")

Y. to ("cpu").

Z = X + Y

throw an exception when two diff

vars. on diff hardware.

# Programming Assignment #1.

def get-next-order

return to a "+1" order

def to-backprop

convert array to backprop

\* class Backprop Array.

def \_\_init\_\_ (np.array).

def \_\_repr\_\_ ( )

→ string.

use graph

Search alg.

def all\_dependencies ( )

find all the dependencies

1.2  
the dependency  
reverse mode

def backward ( )

backward → find gradients

def grad\_fn ( )

"Define math operators in backprop"

def add

→ BA-Add( )

def sub

→ BA-Sub( )

def mul

→ BA-Mul( )

def true\_div

→ BA-Div( )

def sum ( )

→ BA-Sum( ).

def reshape ( )

→ BA-Reshape( ).

def transpose ( )

→ BA-Transpose( ).

Tensor operation

2.1

2.2

2.3

compute "grad"

implement for  
Different classes, scalar - array.

1.1.3

compare ad., nds,  
 (1.5) sd  
 - def numerical\_diff()  
 global  
 fns.  
 - def numerical\_grad()  
 - def backprop\_diff()  
 (2.3) (2.4) (2.5) (2.6)

Store took functions in a class.

class TestFxs():

def {f<sub>i</sub>, df<sub>i</sub>/dx, f<sub>2</sub>, ...}:

Wish'd test implementation: Test Fxs. df<sub>3</sub>/dx (1.4)

input if --name == "Main":

J.

Scalar write the script to execute

output

Concrete & tf. main functions.

compare ad., nds,  
 (1.5) sd

for #3:

$$\frac{\log(x + \alpha\Delta) - \log(x)}{\alpha} = D(f(x)) \quad \text{for } \alpha \gg 0$$

$$\frac{\exp(x + \alpha\Delta) - \exp(x)}{\alpha}$$

$$\frac{\exp(x) \cdot \exp(\alpha\Delta) - \exp(x)}{\alpha}$$

$$\frac{\exp(x) \cdot (\exp(\alpha\Delta) - 1)}{\alpha \cdot \exp(0)}$$

## Week 4 - Lecture 7.

G1

We compute gradient in large-scale

opt. problem  $\rightarrow$  size learning task

examples will most be on supervised

learning, but not limited to ...

$$\text{In SL: } f(w) = \frac{1}{n} \sum_{x,y \in D} L(h_w(x), y)$$

size of dataset.

$$= \frac{1}{n} \sum_{x \in D} f(w; x)$$

In framework, we strive to:

minimize  $f(w)$  over  $w \in \mathbb{R}^d$

computing  $f$  takes  $O(n)$  time.

Init  $w_0 = 0$   $\Rightarrow$  1st nec.  $\nabla f(w_0)$   
gradient.

$$w_{t+1} = w_t - \alpha \nabla f(w_t)$$

Gradient descent

How much time does it cost?

$\Rightarrow O(ndk)$ .

How much memory required?

$\Rightarrow O(nd)$ .

or  $O(d)$

parallel way

running background thru whole thy.

In a NAIVE way  $\rightarrow O(nd)$ .

Newton's Method

$$w_{t+1} = w_t - (\nabla^2 f(w_t))^{-1} \nabla f(w_t)$$

$\hookrightarrow$  converges faster.

"BUT" more expensive than G1

store Hessian  
matrix.

Q: How much time? & How m. Mem?

$\Rightarrow O(nd^2k)$ .

$\Rightarrow O(d^2)$

GD scared away from Newton's method due to the large memory required, i.e. quadratic GD.

\* why are we confident that GD will converge?

the value of gradient is computed based on particular weight,

if move a lil' bit, it will drastically change grad

we want to make it robust.

How close GD ~?

$\rightarrow$  bound. Derivative

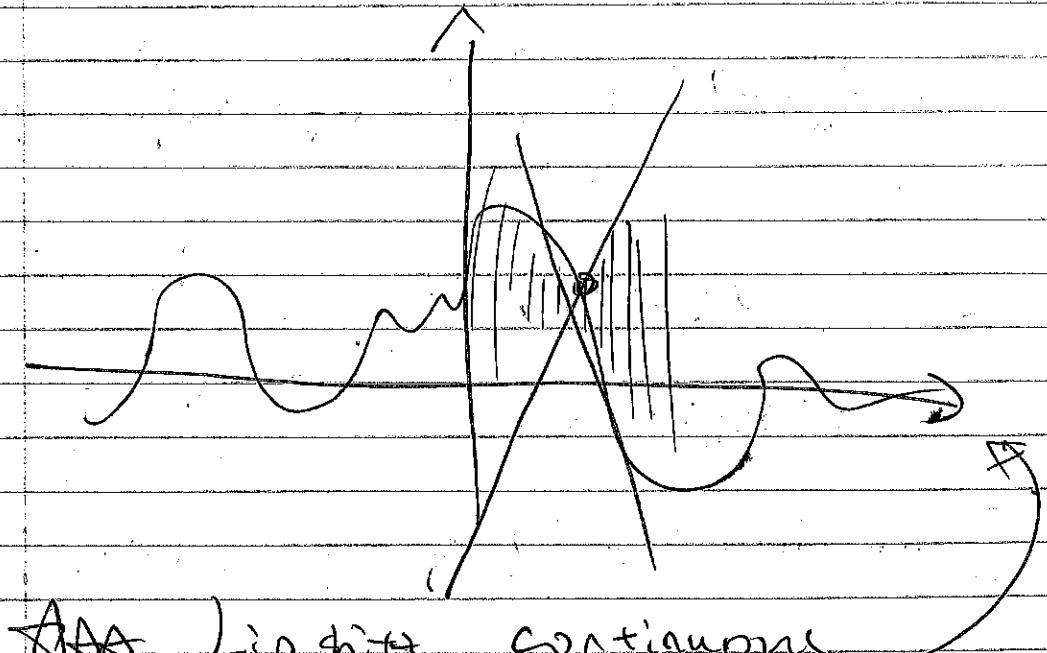
$\hookrightarrow$  Second-order D. is also  $\exists$ .

Gradient  $\nabla f$  is L-Lipschitz continuous

$$\|\nabla f(u) - \nabla f(v)\|_2 \leq L \|u - v\|_2$$

(Assume  $\exists f$  s.t.  $f(w) \geq f^*$ .  
 $\Rightarrow Hw, Hv, Hwv = 1$ .

$$\left\| \frac{\partial^2}{\partial w^2} f(w + \alpha u) \right\| \leq L$$



$\star$  Lipschitz continuous

$\hookrightarrow$  L-smooth

how to prove GD converges ???

$$f(w_{t+1}) = f(w_t - \alpha \nabla f(w_t))$$

fundamental thm. of calc.

$$= f(w_t) + \int_0^{\alpha} \frac{\partial}{\partial \eta} f(w_t - \eta \nabla f(w_t)) d\eta$$



$$\{ f(b) - f(a) = \int_a^b \frac{\partial}{\partial x} f(x) dx.$$

$$= f(w_t) + \int_0^{\alpha} (-\nabla f(w_t))^T \nabla f(w_t -$$

$\eta \nabla f(w_t)) d\eta$ . how?  
Imagine = 0. close to grad.

$$\text{Let } h(\eta) = f(w_t - \eta \nabla f(w_t))$$

$$\{ h(\alpha) = h(0) + \int_0^{\alpha} h'(\eta) d\eta$$

Expect  $\eta$  to be small ← expect

$\Delta x$  small

$$\frac{\partial}{\partial \eta} f(w_t - \eta \nabla f(w_t))$$

$$= \lim_{\delta \rightarrow 0} \frac{f(w_t - (\eta + \delta) \nabla f(w_t)) - f(w_t - \eta \nabla f(w_t))}{\delta}$$

$$= \lim_{\delta \rightarrow 0} \frac{f(w_t - \eta \nabla f(w_t) - \delta \nabla f(w_t)) - f(w_t - \eta \nabla f(w_t))}{\delta}$$

$$= \lim_{\delta \rightarrow 0} \frac{f(\hat{w} + \delta (-\nabla f(w_t))) - f(\hat{w})}{\delta}$$

$$= (-\nabla f(w_t))^T \nabla f(\hat{w}).$$

→ close to norm. → always

Positive ⇒ sd works → min

coat.

$$\frac{\alpha^2}{2} = \int_0^\alpha \eta d\eta$$

$$= f(w_t) + \int_0^\alpha (-\nabla f(w)^T \nabla f(w_\tau)) d\eta$$

$$+ \int_0^\alpha (-\nabla f(w_\tau))^T (\nabla f(w_\tau - \eta \nabla f(w_\tau)) \\ - \nabla f(w_\tau)) d\eta.$$

Apply Cauchy theorem:

$$A \cdot B \leq \|A\| \|B\|$$

$$\{ \|AB\|\}$$

$$\leq f(w_t) + \int_0^\alpha (-\nabla f(w)^T \nabla f(w_\tau)) d\eta$$

$$+ \int_0^\alpha \|(-\nabla f(w_\tau))\| \|\nabla f(w_\tau - \eta \nabla f(w_\tau)) \\ - \nabla f(w_\tau)\| d\eta.$$

$$\leq f(w_t) + \int_0^\alpha \|(\nabla f(w_\tau))\| d\eta$$

$$+ \int_0^\alpha \|(-\nabla f(w_\tau))\| \cdot L \cdot \|\eta \nabla f(w_\tau)\| d\eta.$$

$$f(w_\tau) \leq f(w_t) - \alpha \|\nabla f(w_t)\|^2.$$

$$+ \frac{\alpha L}{2} \|\nabla f(w_t)\|^2$$

$$\leq f(w_t) - \alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(w_t)\|^2$$

$\Rightarrow$  How does constant  $L$  affect how we learn tasks in scales

that's why

assume:  $\alpha L < 1$

$$\leq f(w_t) - \frac{\alpha}{2} \|\nabla f(w_t)\|^2$$

$$\frac{\alpha}{2} \|\nabla f(w_t)\|^2 \leq f(w_t) - f(w_{t+1})$$

Single iteration of GD.

$$\frac{\alpha}{2} \sum_{t=0}^k \|\nabla f(w_t)\|^2 \leq \sum_{t=0}^{k-1} (f(w_t) - f(w_t))$$

$k$  iterations

$$< f(w_0) - f(w_k).$$

$$\frac{\alpha}{2} \sum_{t=0}^{k-1} \|\nabla f(w_t)\|^2 \leq f(w_0) - f(w_k)$$

assume  
 $f$  is bounded from below: ~~\*\*\*~~

$$= \leq f(w_0) - f^*.$$

$$\frac{\alpha}{2} k \min_{t \in \{0, 1, \dots\}} \|\nabla f(w_t)\|^2 \leq f(w_0) - f(w_k)$$

$$\min_{t \in \{0, 1, \dots, k\}} \|\nabla f(w_t)\|^2 \leq \frac{2(f(w_0) - f^*)}{\alpha k}$$

$\Rightarrow$  require  $\alpha k \leq 1$ .

LHS: sum to iteration of (i)  
 return the optimal weights.

RHS: upper bound of func.

Step size  $\alpha \rightarrow \infty$  to large as possible  $\Rightarrow \frac{1}{\alpha} \rightarrow \text{small}$

the largest step

$$\hookrightarrow \leq \frac{2(f(w_0) - f^*)}{k}$$

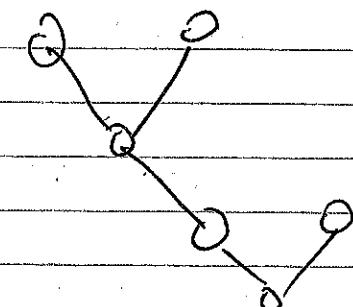
GD only find local opt.

OH (Sep. 13)

(1.1)

o Breadth-first Search ?

what are the dependencies



(1.2)

what is the utility of the grad-fn

why zero out "grad" ?

(1.3)

1.3.4. - why ?

(2)

Generally, how write the help

function to manipulate the

dimensions ?

(2.5) numerical grad - hints ?

(2.7).

$$\mathbb{R}^d \rightarrow \mathbb{R}$$

confused on the implementation.

② writing a for loop

white the for

① add empty row / cols. for

extra dim.  $\rightarrow$  new axis.

check the dims are targeting the same

$\Rightarrow$  broadcasting.

## Vecture 8 . . (Wk. 4).

HW Review:

$\nabla \mathbf{f}$  has same shape as  $\mathbf{X}$

$\mathbf{X}.\text{grad}.\text{shape} = \mathbf{X}.\text{data}.\text{shape}$

\* " += " in numpy mutates the vector

$\text{self.X.grad} += \text{stuff}$

$\text{self.X.grad} = \text{self.X.grad} + \text{stuff}$

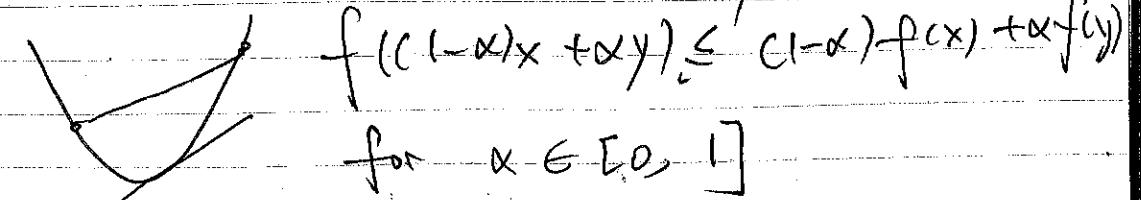
\* init both as arrays as all zeros

GD continued.

"last time"

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla f(\mathbf{w}_t)$$

$\Rightarrow$  most simple condition - convexity  
 ${}^{(n)}\text{-order convexity}$



1st order def. considering whether  
func. cont.

$$f(x) + (y-x)^T \nabla f(x) \leq f(y)$$

2nd order def.  $\frac{\partial^2}{\partial x^2} f(x+\alpha u) \geq 0$  for any  $\alpha, u \in \mathbb{R}^d$

"parabola has a const. sec. deriv."

- convex has a unique global optimum.

"that's why we call 'convex'"

Strongly convex function

"a function that can be bounded  
from below"

$\mu$ -strongly convex function

scalar  
positive num.  $\frac{\partial^2}{\partial x^2} f(x+\alpha u) \geq \mu$  (for  $\|u\|=1$ )

e.g. "Strongly convex".

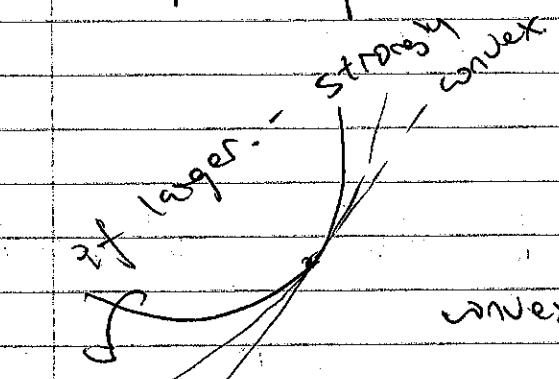
$$f(x) = \max(x, 0)$$

Convexity: vector space  $\rightarrow$  real num.  
i.e. scalar.

Strictly convex vs. Strongly convex

e.g.  $f(x) = e^x$  about not larger  
than a positive val.

$$f(y) \geq f(x) + (y-x)^T \nabla f(x) + \frac{\mu}{2} \|y-x\|^2$$

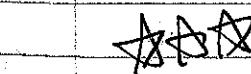


$f$  is  $\mu$ -strongly  
convex, iff  $\exists g$  s.t.

$$f(w) = g(w) + \frac{\mu}{2} \|w\|^2$$

parabola

$\& g$  is convex.



Use in a lot of analysis.

Polyak - Lojasiewicz (PL)

TBC

$$\|\nabla f(w)\|^2 \geq 2\mu (f(w) - f^*)$$

min val.  
↑

$\nearrow$   
but does not  
apply strongly convex  
strongly convex.

$$\text{e.g. } f(w) = \frac{\mu}{2} \|w\|^2$$

$$\nabla f(w) = \mu w$$

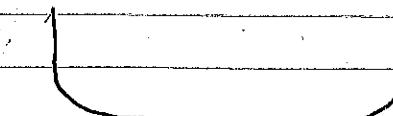
$$\|\nabla f(w)\|^2 = \|\mu w\|^2 = \mu^2 \|w\|^2$$

$$= 2\mu \left( \frac{\mu}{2} \|w\|^2 \right) = 2\mu (f(w) - f^*)$$

★ the gradient is only close to zero

when it's only ~~at~~ around global optimum

satisfies PL yet not strongly CV.



$$f(w_{t+1}) = f(w_t) - \alpha \nabla f(w_t)$$

$$\leq f(w_t) - \alpha \left( 1 - \frac{\alpha L}{2} \right) \|\nabla f(w_t)\|^2$$

non-convex cover of gd.

$$\leq f(w_t) - \alpha \left( 1 - \frac{\alpha L}{2} \right) (f(w_t) - f^*)$$

$$f(w_{t+1}) - f^* \leq (f(w_t) - f^*) - 2\mu \alpha$$

$$\left( 1 - \frac{\alpha L}{2} \right) (f(w_t) - f^*)$$

$$\leq \underbrace{\left( \frac{1}{2} \mu \alpha \left( 1 - \frac{\alpha L}{2} \right) \right)}_{\text{call it } \delta} (f(w_t) - f^*)$$

call it  $\delta$

$$\leq \delta (f(w_t) - f^*)$$

$$f(w_k) - f^* \leq \delta^k (f(w_0) - f^*)$$

if  $\delta > 0$ .

★ How should we choose  $\alpha$ ?

→ show  
 $1 - \frac{\alpha L}{2}$  is non-in

$$2\mu\alpha - \frac{2\mu\alpha^2 L}{2} \rightarrow 0$$

$$1 - \alpha L = 0$$

$$\alpha L = 1$$

$$\nu = \frac{1}{L}$$

$$\text{we get } \gamma = 1 - \ln \frac{1}{L} \left(1 - \frac{1}{2}\right) = 1 - \frac{\mu}{2}$$

$$f(w_k) - f^* \leq \left(1 - \frac{\mu}{2}\right)^k (f(w_0) - f^*)$$

$$\leq \exp\left(-\frac{\mu k}{2}\right) \cdot (f(w_0) - f^*)$$

convergence @ a linear rate

even the is decreasing exponentially

converges with errors  $\propto$  num. of digits.

going to be linear wrt.

"roughly linear with digits int. scale"

Goal: output  $w$  with

$$f(w) - f^* \leq \varepsilon \rightarrow \text{error tol.}$$

$$(\varepsilon \in \mathbb{R}, \varepsilon > 0).$$

It suffices to run  $T$  iterations of  
(num. step)

(3), s.t.

$$\exp\left(-\frac{\mu T}{L}\right) (f(w_0) - f^*) \leq \varepsilon.$$

we can solve this for  $T$ :

$$\exp\left(\frac{\mu}{L}T\right) \leq \frac{f(w_0) - f^*}{\varepsilon}$$

$$\frac{\mu}{L}T \geq \log\left(\frac{f(w_0) - f^*}{\varepsilon}\right)$$

$$T \geq \left(\frac{\mu L}{\mu}\right) \log\left(\frac{f(w_0) - f^*}{\varepsilon}\right)$$

ratio of  
largest open  
val of  $L$   
the smallest  
val.

$$L = K$$

hessen  
for sc.

pos. def.

symm.

"condition  
number".

$$\|\nabla f(x) - \nabla f(y)\| \leq \|x - y\|$$

$$\left\| \frac{\partial^2}{\partial x^2} f(x + \alpha u) \right\| \leq L$$

(for  $\|u\| = 1$ )

$$\text{eqly. } \left\| \frac{\partial^2}{\partial x^2} f(x) \right\| \leq L$$

$\nabla^2 f$

$k'$  range  $\rightarrow$  lot of steps

$k$  small  $\rightarrow$  a few steps.

$k \geq 1 \rightarrow k=1$  isotropic quadratic

Condition number  $\propto$  num. of steps

as  $\alpha$

to solve the

opt. problem,

the other thing impact time is GD.

$\Rightarrow$  size of dataset.

$\frac{\partial^2}{\partial x^2} f(x)$

Stochastic  $\leftarrow$  GD goes well when

GD

DS are large.

Subsampling.  $\rightarrow$  approx.

Principle #2 of LS ML.

SGD:  $g$  faster than GD

GD # steps.  $\sim \mathcal{O}(k)$ .

SGD + step.  $\sim \mathcal{O}(\log k)$

GD

Lecture 7

Week 5

GD time

Sampling set  $\mathcal{X}$

Size

num. of iterations

of GD

$$f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$$

$$\text{GD: } \begin{cases} w_{t+1} = w_t - \alpha \nabla f(w_t) \\ = w_t - \alpha \sum_{i=1}^n \nabla f_i(w_t) \end{cases}$$

Stochastic GD (mini-batch)

$$E [w_{t+1} | w_t] = w_t - \alpha \nabla f(w_t)$$

$$E [\nabla f(w_t) | w_t] = \sum_{j=1}^n P(j=j) \nabla f_j(w_t)$$

$$= \sum_{j=1}^n \bar{\eta} \nabla f_j(w_t)$$

$$= \nabla f(w_t)$$

Assume  $\exists L > 0$  s.t.

$$\|\nabla f(w) - \nabla f(v)\| \leq L \|w - v\|$$

SGD: sample  $\eta$  random  $\eta \in \{1, n\}$ .

$$w_{t+1} = w_t - \alpha \nabla f_\eta(w_t)$$

mini-batch SGD: Sampling

$$w_{t+1} = w_t - \frac{\alpha}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t)$$

batch size

sample is uniformly from  $\{1, n\}$  for each  $b \in \{1, B\}$

Assume  $f(w) \geq f^*$

let  $g_t = g_t$  denote  $\frac{1}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t) - g_t$

$$f(w_{t+1}) = f(w_t - \alpha g_t) = f(w_t) + \int \frac{\partial}{\partial w} f$$

$$w_t - \eta g_t) dw = f(w_t) + \int (\nabla f(w_t - \eta g_t)) (-g_t) dw$$

$$= f(w_t) - \alpha \nabla f(w_t)^T g_t + \int_0^{\alpha} (\nabla f(w_t + \eta) g_t)$$

$$- \nabla f(w_t))^T (-g_t) d\eta.$$

$$\leq f(w_t) - \alpha \nabla f(w_t)^T g_t + \int_0^{\alpha} (\|\nabla f(w_t + \eta) g_t\| d\eta)$$

$$\leq f(w_t) - \alpha \nabla f(w_t)^T g_t + \int_0^{\alpha} \|\eta g_t\| \|g_t\| d\eta$$

$$\leq f(w_t) - \alpha \nabla f(w_t)^T g_t + \frac{\alpha^2 L}{2} \|g_t\|^2$$

$$f(w_{t+1}) \leq f(w_t) - \alpha \nabla f(w_t)^T g_t$$

full batch  
grad

$$+ \frac{\alpha^2 L}{2} \|g_t\|^2$$

stochastic  
grad

$$\mathbb{E}[f(w_{t+1}) | w_t] \leq f(w_t) - \alpha \nabla f(w_t)^T \mathbb{E}[g_t | w_t]$$

$$+ \frac{\alpha^2 L}{2} \mathbb{E}[\|g_t\|^2 | w_t]$$

$$\leq f(w_t) - \alpha \nabla f(w_t)^T \nabla f(w_t)$$

$$+ \frac{\alpha^2 L}{2} \mathbb{E}[\|g_t\|^2 | w_t]$$

$$\leq f(w_t) - \alpha \|\nabla f(w_t)\|^2 + \frac{\alpha^2 L}{2} \mathbb{E}[\|g_t\|^2 | w_t]$$

New Assumption

$$\frac{1}{B} \sum_{b=1}^B \|\nabla f_i(w_t) - \nabla f(w_t)\|^2 \leq \sigma^2$$

for any  $w \in \mathbb{R}^d$ .

$$\mathbb{E}[\|g_t\|^2 | w_t] = \mathbb{E}\left[\left(\frac{1}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t)\right)^T \left(\frac{1}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t)\right)\right] | w_t$$

$$= \mathbb{E}\left[\left(\frac{1}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t)\right)^T \left(\frac{1}{B} \sum_{b=1}^B \nabla f_{i_b}(w_t)\right)\right] | w_t$$

$$= \frac{1}{B^2} \sum_{b=1}^B \sum_{c=1}^B \mathbb{E}[\nabla f_{i_b}(w_t)^T \nabla f_{i_c}(w_t) | w_t]$$

samples that are independent each other

$b \neq c$

$$\begin{pmatrix} & \\ & \end{pmatrix}$$

$$\begin{pmatrix} & \\ & \end{pmatrix}$$

$$= \frac{1}{B^2} (B^2 - B) \|\nabla f(w_t)\|^2 + \frac{1}{B^2} \sum_{b=1}^B \mathbb{E}[\|$$

$$\left(\|\nabla f_{i_b}(w_t)\|\right)^2 | w_t]$$

same

$$E[Y] = E[E[Y|X]]$$

$$= \frac{B-1}{B} \|\nabla f(w_e)\|^2 + \frac{1}{B} \sum_{i=1}^n \|\nabla f_i(w_i)\|^2$$

From  $\rightarrow$  the assumption:

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(w_i)\|^2 = \frac{2}{n} \sum_{i=1}^n \nabla f_i(w)^T \nabla f(w)$$

$$\begin{aligned} &+ \frac{1}{n} \sum_{i=1}^n \|\nabla f(w)\|^2 \leq \sigma^2 \quad \rightarrow \|\nabla f(w)\| \|\nabla f(w)\| \\ &\text{Easy to bound } \|\nabla f(w)\|^2 \end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(w)\|^2 - \|\nabla f(w)\|^2 \leq \sigma^2$$

$$\leq \frac{B-1}{B} \|f(w_e)\|^2 + \frac{1}{B} (\sigma^2 + \|\nabla f(w)\|^2)$$

$$= \|\nabla f(w_e)\|^2 + \frac{\alpha^2}{B}$$

Coming back to 2d assumption

$$E[f(w_{e+1})|w_e] \leq f(w_e) - \alpha \|\nabla f(w_e)\|^2$$

$$+ \frac{\alpha^2}{2} \|\nabla f(w_e)\|^2 + \frac{\alpha^2 \sigma^2}{2B}$$

Variance of  
the square error of SGD

Also assume  $\|\nabla f(w)\|^2 \geq 2\mu(f(w) - f^*)$

PL condition

$$n \leq f(w_e) - \alpha \left(1 - \frac{\alpha L}{2}\right)^2 \mu (f(w_e) - f^*) + \frac{\alpha^2 \sigma^2}{2B}$$

$$E[f(w_e)] - f^*$$

$$E[f(w_e)] \leq \left(1 - \alpha \left(1 - \frac{\alpha L}{2}\right)^2 \mu\right) (f(w_e) - f^*) + \frac{\alpha^2 \sigma^2}{2B}$$

SGD

additional "noise" term.

$$P_t = E[f(w_e) - f^*]$$

$$P_{t+1} \leq \left(1 - \alpha \left(1 - \frac{\alpha L}{2}\right)^2 \mu\right) P_t + \frac{\alpha^2 \sigma^2}{2B}$$

$\Rightarrow$  Assume  $\alpha L \leq 1$ .

$$\leq (1 - \alpha/\mu) P_t + \frac{\alpha^2 \sigma^2}{2B}$$

$$\rho_\alpha = (1 - \alpha/\mu) \rho_\infty + \frac{\alpha^2 \sigma^2 L}{2B}$$

$$\alpha/\mu \rho_\infty = \frac{\alpha \sigma^2 L}{2B}$$

$$f_\alpha = \frac{\alpha \sigma^2 L}{2\mu B} = \frac{\alpha \sigma^2 K}{2B}$$

$$\rho_{t+1} - \frac{\alpha \sigma^2 L}{2\mu B} \leq (1 - \alpha/\mu) \rho_t - \frac{\alpha \sigma^2 L}{2\mu B}$$

$$+ \frac{\alpha^2 \sigma^2 L}{2B}$$

$$\leq (1 - \alpha/\mu) \left( \rho_t - \frac{\alpha \sigma^2 L}{2\mu B} \right)$$

$$\rho_T - \frac{\alpha \sigma^2 L}{2\mu B} \leq (1 - \alpha/\mu)^T \left( \rho_0 - \frac{\alpha \sigma^2 L}{2\mu B} \right)$$

Step:

$$\leq (1 - \alpha/\mu)^T \rho_0$$

$$\mathbb{E}[f(w_T) - f^*] \leq (1 - \alpha/\mu)^T (f(w_0) - f^*)$$

$$+ \frac{\alpha^2 \sigma^2 L}{2\mu B}$$

noise wall  
noise floor

$$\Rightarrow 1 - x \leq \exp(-x) \quad (\text{plug in})$$

LHS

$$\leq \exp(-\alpha \mu T) (f(w_0) - f^*) + \frac{\alpha \sigma^2 L}{2\mu B}$$

lecture 9. Week 5 - 2.

Rev - last time.

$$\mathbb{E}[f(w_t) - f^*] \leq \exp(-\alpha \mu T)(f(w_0)$$

$$- f^*) + \frac{\alpha L \sigma^2}{2B\mu}.$$

$$\text{Goal: } \mathbb{E}[f(w_{\text{out}}) - f^*] \leq \varepsilon.$$

$$\text{suffice for } \varepsilon \geq \exp(-\alpha \mu T)(f(w_0) - f^*)$$

$$+ \frac{\alpha L \sigma^2}{2B\mu}.$$

$$\text{also suffice } \frac{\varepsilon}{2} \geq \exp(-\alpha \mu T)(f(w_0) - f^*)$$

$$\text{and } \frac{\varepsilon}{2} \geq \frac{\alpha L \sigma^2}{2B\mu}.$$

$$\Rightarrow \alpha \leq \frac{B\mu\varepsilon}{L\sigma^2}$$

$$\log\left(\frac{\varepsilon}{2(f(w) - f^*)}\right) \geq -\alpha \mu T.$$

$$\Rightarrow T \geq \frac{1}{\alpha \mu} \cdot \log\left(\frac{2(f(w) - f^*)}{\varepsilon}\right)$$

also  $\alpha L \leq 1$ .

$$\alpha \leq \min\left(\frac{B\mu\varepsilon}{L^2}, 1\right) \cdot \frac{1}{L}$$

$$T \geq \max\left(\frac{\sigma^2}{B\mu\varepsilon}, 1\right) \frac{L}{\mu} \log( \dots )$$

$$\frac{2(f(w_0) - f^*)}{\varepsilon}$$

$$\geq \max\left(\frac{\sigma^2}{B\mu\varepsilon}, 1\right) \cdot K \cdot \log(n)$$

(runtime of SGD):  $\mathcal{O}\left(\max\left(\frac{\sigma^2}{B\mu\varepsilon}, 1\right) \cdot \frac{B}{\varepsilon} \cdot K \cdot \log\left(\frac{1}{\varepsilon}\right)\right)$

usually,  $\varepsilon$  are small  $\rightarrow$

$$\mathcal{O}\left(\frac{\sigma^2}{B\mu\varepsilon} K\right)$$

$\rightarrow$  batch  $B$  cancels

$\rightarrow$  ignoring "log" term.

runtime of GD.  $\mathcal{O}(NK \log(\frac{1}{\epsilon}))$

for GD  $n \rightarrow \frac{\nabla^2}{B\epsilon}$ .

GD

SGD

- convex (strongly).
- o Non-convex

o Global opt. guaranteed  
(unique).

o Large memory. (GPU).

parallel

o really small  $\Sigma$

o large  $\Sigma$

How to set batch size  $B$

→ Hardware consideration.

$$\text{e.g. } f_i(w) = \frac{1}{2} (x_i^T w - y_i)^2$$

$$\nabla f_i(w) = x_i (x_i^T w - y_i)$$

minibatch:  $X \in \mathbb{R}^{B \times d}$ ,  $y \in \mathbb{R}^B$ .

$$\nabla f_{\text{batch}}(w) = \frac{1}{B} X^T (Xw - Y)$$

faster

MMPI  $(B, d) @ (d, 1)$ ,  
 $(d, B) @ (B, 1)$

$(1024, 1024) @ (1024, 1024)$  a lot  
 $(1024, 1024) @ (1024, 1)$  less than  $1024 \times$

minibatch size  $\propto$  hardware itself

e.g.  $B=256, 64, \dots$

SGD → GD: Sampling w/o replacement.

$\max(n, 1)$  is GJ  
longer, slower scenario

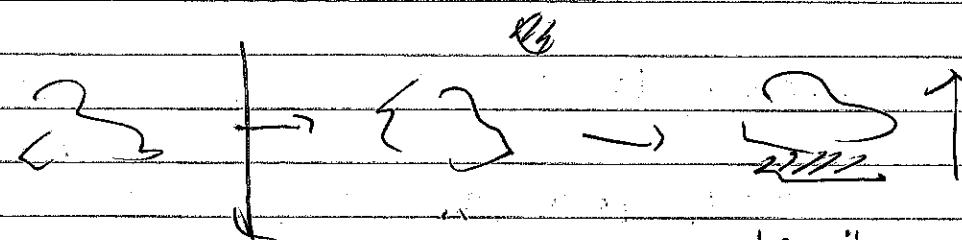
Random reshuffling

sample w/o repl., only reuse  
example after every example

has been used.

epoch  $\Rightarrow$  1 pass thru the training set  
(e.g. 100,  $n \geq 100$ ).

shuffle - once



Statistically not

use a fixed step size.

Decreasing step size (assuming  $\alpha_t \in I$ )

$$E[f(w_{t+1}) - f^*] \leq (1 - \alpha_t \mu)$$

$$E[f(w_{t+1}) - f^*] \leq \frac{\alpha_t \sigma^2 L}{2B}$$

$$P_{t+1} \leq (1 - \alpha_t \mu) P_t + \frac{\alpha_t^2 \sigma^2 L}{2B}$$

$$\Rightarrow \dot{P} = -\mu P + \frac{\alpha \sigma^2 L}{B}$$

$$\Rightarrow \Delta t \approx -\frac{P_t / \mu B}{\sigma^2 L}$$

$$P_{t+1} \leq P_t - \frac{\mu^2 B}{2\sigma^2 L} P_t^2$$

$$\overset{\downarrow}{P_{t+1}} \geq \overset{\downarrow}{P_t} - \frac{\mu^2 B}{2\sigma^2 L} \overset{\circ}{P_t^2}$$

$$= \frac{1}{P_t} \left( 1 - \frac{\mu^2 B}{2\sigma^2 L} P_t^2 \right)^{-1}$$

$$\geq \frac{1}{P_t} \left( 1 + \frac{\mu^2 B}{2\sigma^2 L} P_t \right) = \frac{1}{P_t} + \frac{\mu^2 B}{2\sigma^2 L}$$

$$\Rightarrow \frac{1}{P_t} > \frac{1}{P_0} + \frac{\mu^2 B T}{2\sigma^2 L} \Rightarrow P_t = E[f(w_t) - f^*]$$

Week 11 Week 6.

$$\leq \left( \frac{1}{f(w_0) - f^*} + \frac{n^2 BT}{2\sigma^2 L} \right)^{-1}$$

$$= \Theta\left(\frac{1}{T}\right)$$

$$\approx \mathcal{O}\left(\frac{\sigma^2 K}{nBT}\right)$$

Momentum:

Rev:  $\kappa \sim \frac{h}{n}$ ; great. noise  $\mathcal{N}(\frac{\alpha}{n})$

- Arrangement:
  - PSet 3
  - PA 2.

Rev: G1

$$\ell \rightarrow f(w) - f^* \leq \exp(-\frac{T}{K})(f(w_0) - f^*)$$

gap after  $T$  epochs

Q: Simplist model with large  $K$ ?

a 2D problem  $\rightarrow$  Quadratic

$$\text{e.g. } f(w) = \frac{1}{2}(w^T L w + w^T \mu)$$

$$= \frac{1}{2} w^T \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} w$$

$$\nabla^2 f(w)$$

$$\Rightarrow H(w) = \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix}, \text{ suppose } L > \mu > 0$$

$$\nabla f(w) = \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} w$$

$$w_{t+1} = w_t - \alpha \nabla f(w_t)$$

$$= w_t - \alpha \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} w_t$$

$$= \begin{bmatrix} 1 - \alpha L & 0 \\ 0 & 1 - \alpha \mu \end{bmatrix} w_t$$

↑ power " $T^k$ "

$$\Rightarrow w_T = \begin{bmatrix} 1 - \alpha L & 0 \\ 0 & 1 - \alpha \mu \end{bmatrix}^T w_0$$

$$= w_T - \begin{bmatrix} (1 - \alpha L)^T & 0 \\ 0 & (1 - \alpha \mu)^T \end{bmatrix} w_0$$

$$f(w_T) = \frac{1}{2} \left( L((1 - \alpha L)^T (w_0))^2 + \mu ((1 - \alpha \mu)^T (w_0))^2 \right)$$

convergence

$$\text{Rate of convergence} = \frac{1 - \alpha \mu}{1 - \alpha L}$$

$$\alpha = \frac{2}{L + \mu}$$

$$= \frac{1 - \mu}{L + \mu} =$$

$$= 1 - \frac{2}{k+1}$$

$$f(w_T) = \sigma \left( \left( 1 - \frac{2}{k+1} \right)^{2T} \right) = \sigma \left( \exp \left( - \frac{4T}{k+1} \right) \right)$$

$$w_{t+1} = w_t - \alpha \nabla f(w_t) + \beta (w_T - w_{t+1})$$

Previous step  
use the next step, step size  $\rightarrow$   
large, overshooting  
the object.

"Self-tune" the step size

$$w_{t+1} - w_t = -\alpha \nabla f(w_t)$$

$$\beta (w_t - w_{t+1}) - \alpha \nabla f(w_t)$$

we need

$$0 < \beta < 1$$

"some" force term

"friction" term  
or momentum

How Polyak's momentum applies to

$$\Rightarrow \text{sgd} \xrightarrow{\text{?}}$$

$$w_{t+1} = w_t - \alpha \begin{bmatrix} L & \\ & \mu \end{bmatrix} w_t + \beta (w_t - w_{t-1})$$

$\curvearrowleft$

$$\begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} = \begin{bmatrix} w_t - \alpha A w_t + \beta (w_t - w_{t-1}) \\ w_t \end{bmatrix}$$

$$= \begin{bmatrix} (1+\beta)I - \alpha A & -\beta \\ I & 0 \end{bmatrix} \begin{bmatrix} w_t \\ w_{t-1} \end{bmatrix}$$

$$\begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} = \begin{bmatrix} (1+\beta)I - \alpha A & -\beta I \\ I & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix}$$

$\curvearrowleft$

$$\boxed{\begin{bmatrix} (1+\beta)-\alpha L-\alpha L & -\beta \\ 1 & 0 & & \\ & 0 & (1+\beta)-\alpha L-\beta \\ & & 0 & 1 \end{bmatrix}}$$

4x4 block matrix

Eigenvalue

use  $\lambda$  to find eigenvalues.

$$\begin{bmatrix} ((1+\beta)-\alpha L-\beta) - \lambda & -\beta \\ 1 & 0 + \lambda \end{bmatrix} = \lambda^2 - (1+\beta-\alpha L) \lambda + \beta$$

$$\lambda = \frac{1+\beta-\alpha L \pm \sqrt{(1+\beta-\alpha L)^2 - 4\beta}}{2}$$

Complex vals.

Polyadic gives:

$$\lambda = \frac{L+2\beta}{L/\mu}, \quad \sqrt{\beta} = \frac{\sqrt{L}-1}{\sqrt{L}+1}$$

$$|\lambda|^2 = \frac{1}{4} [(1+\beta-\alpha L)^2 + 4\beta - (1+\beta-\alpha L)^2]$$

$$T \approx \sqrt{L} \cdot \sqrt{\beta} \approx \exp\left(-\frac{L}{\sqrt{L}+1}\right)$$

Works when  $L(\beta \geq (1+\beta-\alpha L)^2)$

$$\begin{aligned} |2\sqrt{\beta}| &\geq |1+\beta-\alpha L| \\ 2\sqrt{\beta} &\geq 1+\beta-\alpha L \geq -2\sqrt{\beta} \end{aligned}$$

Lesson 11. - Week 6.

$$\mu \leq \lambda \leq L$$

Conditioning.

Rev. - Momentum.

$$T = \mathcal{O}\left(nK \log\left(\frac{1}{\epsilon}\right)\right)$$

condition run.

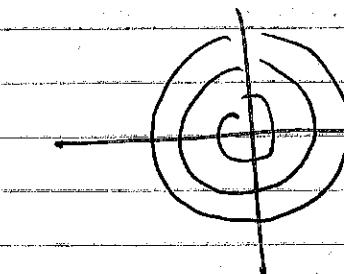
last see.

Condition be too large

/ too small.

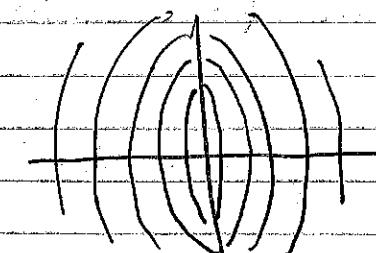
$$T_{\text{mom}} = \mathcal{O}\left(\sqrt{nK} \log\left(\frac{1}{\epsilon}\right)\right)$$

$$f(w_1, w_2) = \frac{1}{2}w_1^2 + \frac{1}{2}w_2^2 \quad (K=1)$$



$$C = w_1^2 + w_2^2$$

$$f(w_1, w_2) = \frac{10}{2}w_1^2 + \frac{1}{2}w_2^2$$



$$C = aw_1^2 + bw_2^2$$

"poorly-conditioned" prob.

$$w = P u$$
$$\min_w f(w) = \min_{P u} f(P u)$$

ALG.  $|P| \neq 0$

$\downarrow$

invertible  $\rightarrow$  any matrix, invertible

$$\arg\min_w f(w) = P (\arg\min_u f(P u))$$

let  $g(u) = f(P u)$ .

graphically, this is essentially

"stretching" the poorly-conditioned

problem, i.e. apply linear operation



Map back  $\leftarrow$  reduced K

$\Downarrow$   
Solve for optimal

Preconditioning.

instead of solving minimize  $f(w)$

over  $w$  solve minimize  $f(P u)$

over  $u$  (for clever  $P$ ).

e.g.  $f(w) = \frac{1}{2} w^T A w$ .  $A^T = A$

- how to set  $P$ ? what cond. can we get? (PD,  $A > 0$ , all eigenvalues

$A$  are positive)

$P^T A P = I$   $\rightarrow$  symmetric.

$$2f(P u) = (P u)^T A (P u) = u^T P^T A P u$$

$$= u^T u = \|u\|^2$$

Suppose  $P$  symmetric, invertible

-  $P^T A P = I$

-  $P^{-1} P A P P^T = P^{-1} I P^{-1}$

$A = P^{-2} = P = A^{-2}$

replace  
eigenval.

in transformed space.

$$g(w) = f(Pw) \rightarrow \nabla g(w) = P^T \nabla f(Pw)$$

$$U_{t+1} = U_t - \alpha \nabla g(U_t) = U_t - \alpha P^T \nabla f(PU_t)$$

assert:  $w = Pu$ ;  $U_t = Pv_t$

$$PU_{t+1} = PU_t - \alpha PP^T \nabla f(PU_t)$$

$$W_{t+1} = W_t - \alpha P^T \nabla f(W_t)$$

$P$  Preconditioner

$$W_{t+1} = W_t - \alpha R \nabla f(W_t)$$

positive-definite  
symmetric.

Preconditioned GM

$$R = A^{-1}$$

$$x^T R x = x^T P P^T x$$

$$= (P^T x)^T (P^T x)$$

$$= \|P^T x\|^2 > 0$$

choose  $R$  to be diagonal preconditioner

$R$  diagonal  $\rightarrow$  stretch along the coordinates' axes  
memory  $O(d)$ .

time multiply  $O(d)$ .

$$(W_{t+1})_i = (W_t)_i - \alpha R_{ii} (\nabla f(W_t))_i$$

let  $\alpha_i = \alpha R_{ii}$ .

$$(W_{t+1})_i = (W_t)_i - \alpha_i (\nabla f(W_t))_i$$

$$W_{t+1} = W_t - \alpha \odot \nabla f(W_t)$$

element-wise

run GM in implicit scaled space.

AdaGrad

loop

$$g \leftarrow \nabla f_i(w), \text{ random } i$$

$$\tau_i \leftarrow \tau_i + g_i^2$$

## Lecture 12. Week 7.

$$(\mathcal{F}_i)_t = \sum_{k=0}^t (g_i^2)_t$$

$$w_i \leftarrow w_i - \frac{\alpha}{\sqrt{\mathcal{F}_i}} g_i \quad (\alpha_i = \frac{\alpha}{\sqrt{r_i}})$$

- AdaGrad

(G) "Adaptive Optimization" Alg.

AdaGrad.

init  $w \in \mathbb{R}^d$ ,  $r = 0 \in \mathbb{R}^d$ .

loop:

depends too much on the history get example grad  $g \in \mathbb{R}^d$

infinite

memory.

large example update  $r \leftarrow r + g^2$

) update  $w \leftarrow w - \frac{\alpha}{\sqrt{r}} \cdot g$

- done in element-wise

impractically.

$$r_i \approx \mathbb{E}[g_i^2]$$

small number

num. num. of steps.

$$\mathcal{G} \approx r (\mathbb{E}(g_i^2) + \text{Var}(g_i))$$

could be dominant

fix this: RMSProp.

### RMSProp

hyperparameters  $\alpha > 0$ ,  $p \in (0, 1)$ ,

init  $w \in \mathbb{R}^d$ ,  $r = 0 \in \mathbb{R}^d$

loop:

get gradient sample  $g$  at  $w$ .

update  $r \leftarrow p r + (1-p) g^2$

update  $w \leftarrow w - \frac{\alpha}{\sqrt{r} + \epsilon} g$

$$P_i \approx \mathbb{E}[g_i^2]$$

typical  $p = 0.99$

Combining momentum + RMSProp +

Correction for  $0 = \text{init}$

init  $w \in \mathbb{R}^d$ ,  $r, \sigma = 0 \in \mathbb{R}^d$

$$S = 0 \in \mathbb{R}^d$$

Momentum:

$$w_{t+1} = w_t - \alpha g_t + \beta (w_t - w_{t+1})$$

$$(w_{t+1} - w_t) = \beta (w_t - w_{t+1}) - \alpha g_t$$

$$\begin{cases} v_{t+1} = \beta v_t - \alpha g_t = \beta v_t - (1-\beta) \frac{\alpha}{(1-\beta)} g_t \\ w_{t+1} = w_t + v_t \end{cases}$$

$\alpha, p_1, p_2$ , hyperpara.

loop:

get gradient sample  $g$  @  $w$ .

$$s \leftarrow p_1 s + (1-p_1) g^2$$

$$r \leftarrow p_2 r + (1-p_2) g^2$$

Update:

$$w \leftarrow w - \frac{\alpha}{\sqrt{r} + \epsilon} s$$

$$s \leftarrow \frac{s}{1-p_1}$$

$$r \leftarrow \frac{r}{1-p_2}$$

$t = \# \text{ of steps}$

if  $g_t$  denotes the  $g$  drawn at

$$\text{Step } t, S_t = \sum_{k=0}^{t-1} p_1^k (1-p_1) g_{t+k}$$

$S$  after  $t$  steps

$$\sum_{k=0}^{t-1} p_1^k (1-p_1)$$

$$= (1-p_1) \left( \frac{1-p_1^t}{1-p_1} \right)$$

SI: superexponential MA:  
sequence  $x_0, x_1, x_2, \dots$  reverse

$$S_0 = 0; S_{t+1} = p S_t + (1-p) x_t$$

- today - parameters:

subsampling - binary number  $n$ .

momentum -  $k$ .

Variance of the gradient.

↳ avoid having

sample

### Variance Reduction

→ SVRG

↳ strive to reduce  $G_D$

- yet SGD each steps.

↳ Large-Scale convex optimization

→ great !!

for DL - Not so good.

$G_D$

SVRG.

$$G(nk \log(\frac{1}{\epsilon})) \quad O((n+k) \log(\frac{1}{\epsilon}))$$

### Polyak Averaging

- Output an average:

if  $W_t$  is parameters vector after

$t$  steps

$$\hookrightarrow \frac{1}{T} \sum_{t=0}^{T-1} W_t$$

Stochastic weight averaging

(SWA).

(wed.)

Wednesday 15/14. Week 7.

few ...

$$\mathcal{O}(n+k \log(\frac{1}{\epsilon})) - \text{GI}$$

$$\mathcal{O}(K \epsilon^2) - \text{SGI}$$

Momentum:  $\mathcal{O}(n\sqrt{K} \log(\frac{1}{\epsilon})) - \text{MSGI}$

$$\mathcal{O}((n+k) \log(\frac{1}{\epsilon})) - \text{SRC}$$

$\rightarrow$  Polyak averaging / Stochastic

weight averaging

$$\mathcal{O}(nK(d \cdot \log(\frac{1}{\epsilon}))) \Leftrightarrow \mathcal{O}(K \sigma_d^2 d / \epsilon^2)$$

$$\text{linear model: } w \in \mathbb{R}^d, \mathcal{O}(n\sqrt{K} d \log(\frac{1}{\epsilon}))$$

$$J(w) = \frac{1}{n} \sum_{x_i, y_i} \mathcal{L}(w^T x_i, y_i).$$

time to compute an example grad

$$\mathcal{O}(d)$$

$P(A) \Rightarrow$  covariance matrix.  $(d \times d)$ .

computationally expensive

Random Projection: (Johnson-Lindenstrauss)

$$x \in \mathbb{R}^d \Rightarrow Ax \in \mathbb{R}^d \text{ (trans form)}$$

feature  $A \in \mathbb{R}^{d \times D}$  random length

$$A_{ij} \sim N(0, \sigma^2)$$

$$\leq \|Ax - Ay\|^2 \leq$$

$$(1-\epsilon) \|x_i - x_j\|^2 \leq \|Ax_i - Ax_j\|^2 \leq (1+\epsilon) \|x_i - x_j\|^2$$

$$\forall x_i, x_j \in \mathcal{D}$$

$$d \geq \frac{8 \log(n)}{\epsilon^2}$$

Why? choose  $\sigma^2$ . s.t.

$$\mathbb{E} [\|Ax_i - Ax_j\|^2] = \|x_i - x_j\|^2$$

$$(\mathbb{E} [\|(Ax_i - Ax_j)\|^2]) - (\mathbb{E} [(x_i - x_j)^T A^T (x_i - x_j)])$$

$$= (x_i - \bar{x})^T [E[A^T A]] (x_i - \bar{x})$$

↳ Gaussian matrix.

$$[E[A^T A]]_{ij} = \sum_{k=1}^d E[A_{ki} A_{kj}]$$

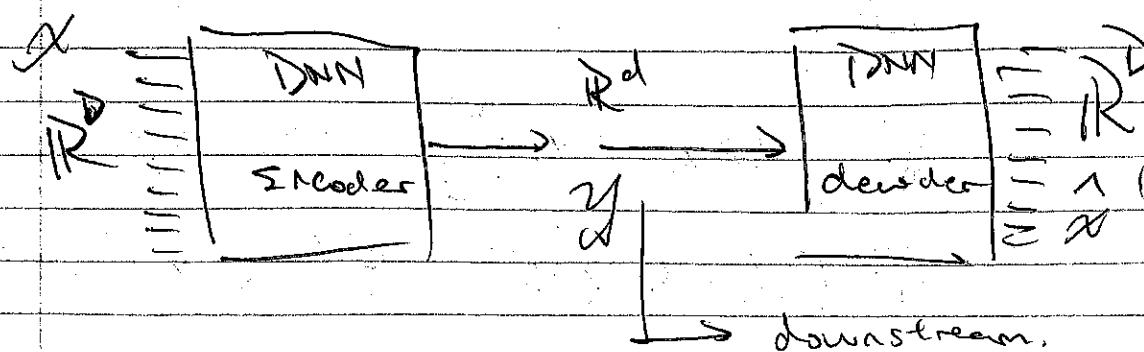
$$\begin{aligned} &= \left\{ \begin{array}{ll} 0 & \text{if } i \neq j \\ \sum_{k=1}^d E[A_{ki}^2] & \text{if } i=j. \end{array} \right. \\ [E[A^T A]] &= \sigma^2 d I. \end{aligned}$$

$$\Rightarrow = (x_i - \bar{x})^T (\sigma^2 d I) (x_i - \bar{x})$$

$$= \sigma^2 d \|x_i - \bar{x}\|^2 \stackrel{\text{if}}{=} \|x_i - \bar{x}\|^2$$

(since  $\sigma^2 = \frac{1}{d}$ )

Autoencoders



$$\ell(x, \hat{x})$$

$$\text{e.g. } \sum_{i \in \text{ID}} \|x - \hat{x}\|_2^2$$

↳ unsupervised learning.

Sparcity

density of  $X = \frac{\# \text{ of non-zero entries}}{\text{size of } X \cdot d}$

$$\begin{bmatrix} 0 \\ 0 \\ 2 \\ 0 \\ 0 \\ 0 \\ 3 \end{bmatrix} = \{(4, 2, 0), (7, 1, 0)\} = \text{- indices.} = [4, 7]$$

(6. pos)

v8

Sparse matrix.

COO

→ expensive to comp.  
→ easy to write

$$\begin{bmatrix} 0 & 5 & 0 & 3 \\ 0 & 10 & 0 & 0 \end{bmatrix}$$

→ row idx. [1, 1, 2]

→ col. idx. [3, 6, 2]

values

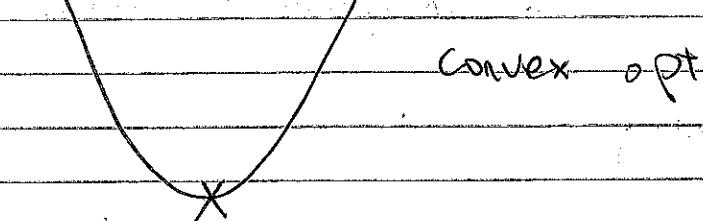
[5, 3, 1]

CSC → tr. CSC  
→ row offsets. [0, 2] [3, 6] [2] → [3, 6, 2]

→ col cols. [5, 0, 3, 0] [1, 0] [5, 0, 3, 0, 1, 0]

Q Week 8. Lecture 15/16

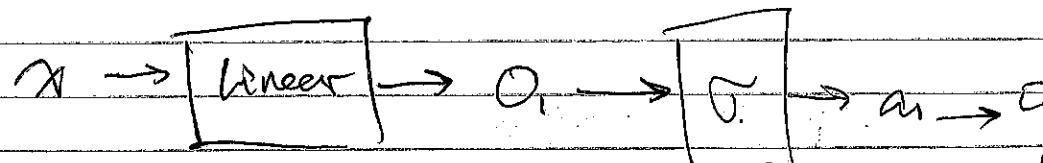
## Deep learning



### Neural Networks

• MLP

$$x \in \mathbb{R}^{d_0}$$



$$o_1 = w_1 x + b_1$$

$$w_i \in \mathbb{R}^{d_{i-1} \times d_i}$$

$$a_1 = \sigma(o_1)$$

$$\left\{ \begin{array}{l} w_i \in \mathbb{R}^{d_{i-1} \times d_i} \\ b_i \in \mathbb{R}^{d_i} \end{array} \right.$$

$$o_i \in \mathbb{R}^{d_i}$$

Nonlinear: function onto element-wise

as the activation function.

$$\sigma(x) = \text{ReLU}(x) = \max(x, 0).$$

universal: for continuous input  $\rightarrow$  output.

there exists a function  $f$  approx.

with bounded parameters (good accuracy).

dependency  $\xrightarrow{\text{imply}}$  overfit.

the NN tend not to overfit!

$\rightarrow$  larger models generalize better.

▷ double-descent.

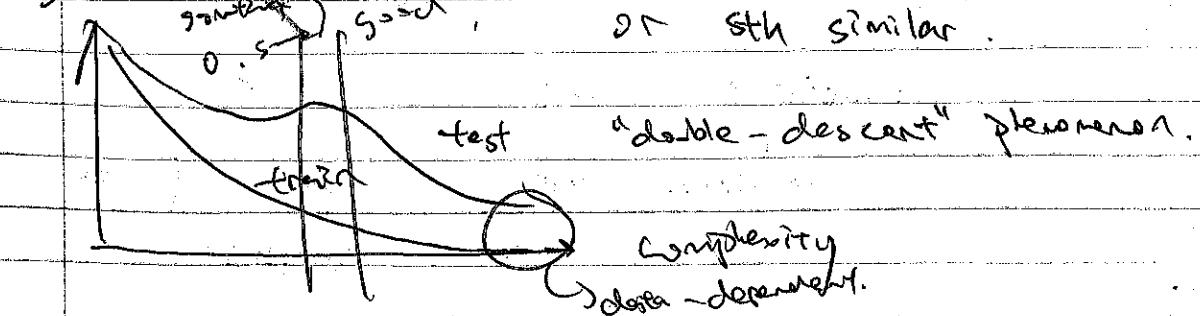
$$(x, y) \in \mathbb{R}^d \times \mathbb{R}.$$

$$|\mathcal{D}| = n. \quad n=d.$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n (x_i^T w - y_i)^2 + \text{regularization}$$

larger model generalize better

Specifying for SGD, Adam.  
loss function gradient good, or 5th similar.



$a_e, o_e \in \mathbb{R}^{d_e}$

$$a_e = W_e o_e + b_e \rightarrow J( \text{ } \overset{d_e}{\underset{d_1}{\text{---}}})$$

$\beta \cdot [ \text{ } \overset{d_1}{\underset{d_1}{\text{---}}} ] \in \mathbb{R}^{d_1}$

$$o_e = o_e(a_e) \rightarrow J( \text{ } \overset{d_1}{\underset{d_1}{\text{---}}})$$

$w_i \in \mathbb{R}^{d_1 \times d_0}$

$b_i \in \mathbb{R}^{d_1}$

$x \in \mathbb{R}^{d_0}$

Image processing ( $\leftarrow c_w$ )

$x \in \mathbb{R}^{c_w \times h_w}$   $\rightarrow$  height of image

width of image

channels.

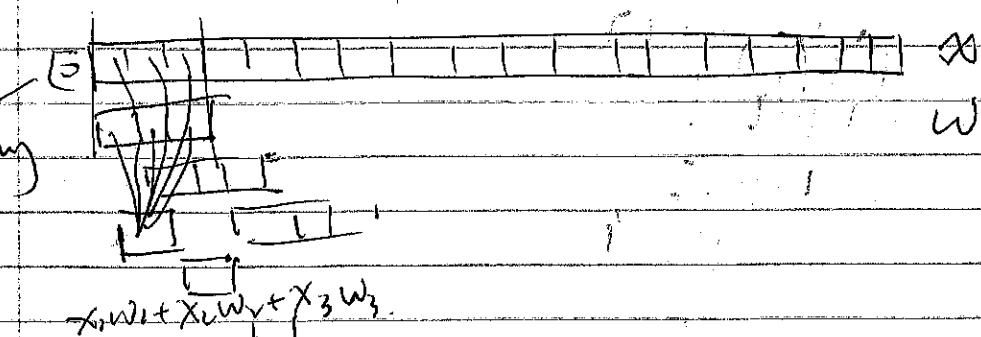
1) convolutional neural network

- restricted subset

convolutional layers. in place of the linear layers.

$d_i(d_0 + d_0 - 1) + d_1$

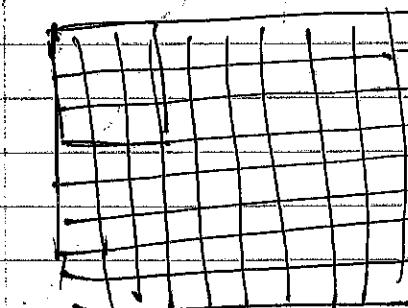
Convolution in 1D:



fixed arbitrary set of parameters.

$$x_{k+1}w_1 + x_kw_2 + x_{k-1}w_3$$

$\hookrightarrow 2D$



diff. channels

Conv. layer

input  $x \in \mathbb{R}^{c_w \times h_w}$  output  $a \in \mathbb{R}^{c_o \times h_o \times w_o}$

parameters  $W \in \mathbb{R}^{c_o \times c_w \times h_w \times w_h}$

$$a_i = \sum_{j=1}^{c_o} \text{conv.}(W_j, x_j) \quad W \in \mathbb{R}^{c_o \times c_w \times h_w \times w_h} \quad b \in \mathbb{R}^{c_o}$$

Gorkaray

Week 9. lecture 17.

Overfitting / Underfitting  $\leftarrow$  capacity

$\leftarrow$  model has high capacity.

(Overparametrized)

DNNs with  $\leftarrow$  representational  
easy weights.

weights coming  $\leftarrow$  effective

from the learning algorithm.

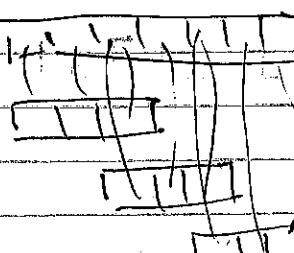
Dropout

$\hookrightarrow$  increase the representation capacity.

During training  $\rightarrow$  randomly drop

(independently each step) neurons.

forces the learn to employ  
more neurons.



Output: 3

Input image:

$3 \times 64 \times 64 \rightarrow$  conv  $\rightarrow$  ReLU  $\rightarrow \dots \rightarrow$  (eg 2).  
 $16 \times 64 \times 64$  stride 1  
 $32 \times 32 \times 32$

$\dots \rightarrow$  Flatten  $\rightarrow$  (linear)  
 $a$   
 $\downarrow$   
ReLU  $\rightarrow$  MLP  
 $\downarrow$   
Out.

### Batch Norm.

let  $u \in \mathbb{R}^n$

$$u \rightarrow \frac{u - \mu}{\sqrt{\sigma^2}} = \frac{u - \mu}{\sigma}$$

Instead, look at  $u \in \mathbb{R}^B$

- conv.  $\rightarrow$  merge  $\rightarrow$  per channel!

- saves a running exp.  
average,

of the mean & variance

$$\mu \text{ and } \sigma$$

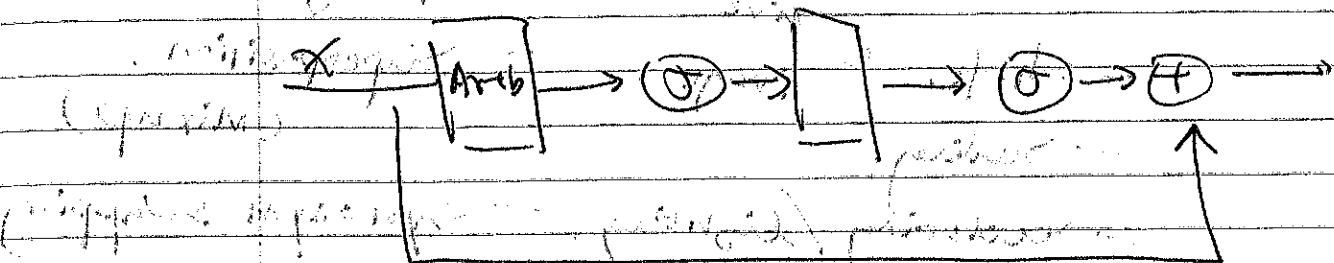
Model.train()  $\rightarrow$  turning on/off

Model.eval()  
batch norm.

$$u \rightarrow \frac{u - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

batchnorm scalar learned param (SGD)

Residual connections / skip connections.



### ResNet

### DenseNet

$$\min_w l(w) = \frac{1}{n} \sum_{i=1}^n l(w, x_i, y_i)$$

$$(x_i, y_i) \in D \subseteq \mathbb{R}^n \times \mathbb{R}$$

$\triangleright$  sparse data ( $D$  sub)

$n$  too small.

randomness.

$\star$  Data Augmentation.

Aug. function  $A: \mathcal{X} \xrightarrow{[0,1]} \mathcal{X}$ .

$$w_{t+1} = w_t - \alpha \nabla l(w_t; A(x_t; \eta_t), y_t)$$

example chosen at  $t$ .

IMAGES:

- Translation → Noise / Artif fact
- Rotation
- Shift, scale → Cropping, Addition
- Sealing → Superposition (mix up)
- rendering / lighting, → synonym swapping.
- deletion
- backgrounds.

clustering → data is clustered

labels by cluster.

Manifold → data lie on low-dimensional manifold

PCA  dimensionality reduction

Semi-supervised learning

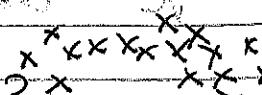
"data - cheap" "label - expensive."

lots of data, few labels

assump.

Similar inputs have similar outputs.

↳ (kNN)



B

0

↳ kNN finds the nearest neighbor  
and assigns it the same label.

# Week 10, lecture 19.

## Sequence models.

$$x \in \mathbb{R}^{d \times n}$$

different per example.

Discussion:

take sequence as input.

↳ generate sequence output.

↳ handle sequence input / output?

- padding to max length.

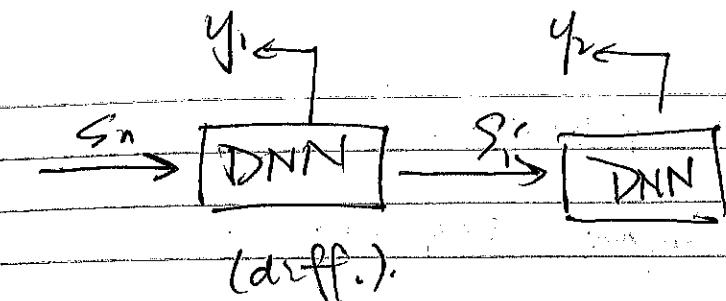
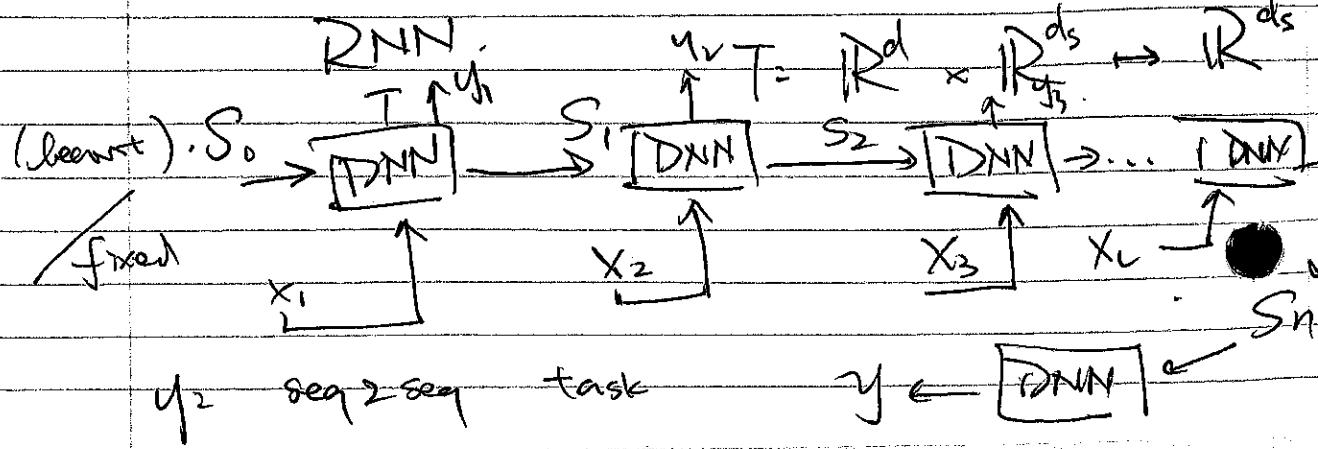
- counts / sum

- RNN  $\approx$  sequential

- Transformers  $\approx$  parallel.

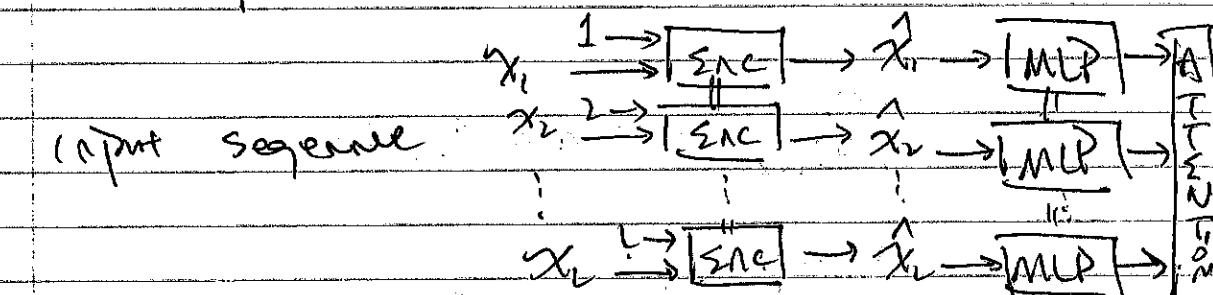
$$\text{DFA: } S_{\text{next}} = T(S_{\text{current}}, x_i).$$

$$x \in \mathbb{R}^c$$



↳ LSTM.

↳ Transformers.



$$x_i = \text{Enc}(x_i, \cdot)$$

Mathematically, attention layer:

$$x_i \xrightarrow{\text{FF}} k_i$$

$$x_i \xrightarrow{\text{FF}} v_i$$

$$x_i \xrightarrow{\text{FF}} v_i$$

$$\vdots$$

$$x_i \xrightarrow{\text{FF}} k_i$$

$$x_i \xrightarrow{\text{FF}} v_i$$

$$\vdots$$

$$Q_i \in \mathbb{R}^{d_k}$$

$k_i \in \mathbb{R}^{d_k}$

$v_i \in \mathbb{R}^{d_v}$

$$Q, K \in \mathbb{R}^{n \times d_k}$$

$$V \in \mathbb{R}^{n \times d_v}$$

$$\text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V$$

approximate table lookup:

Attention  
block.

## Multihed Attention

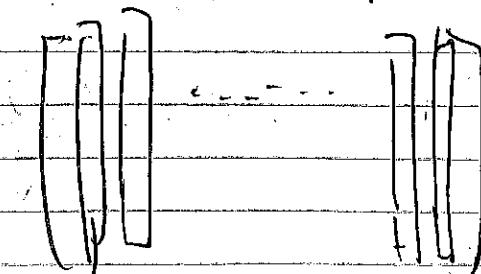
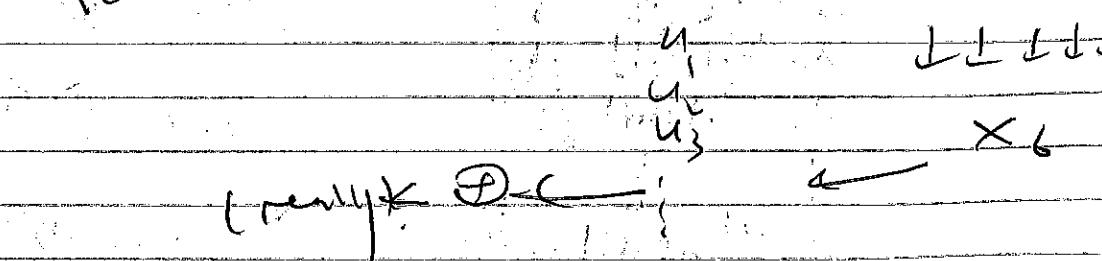
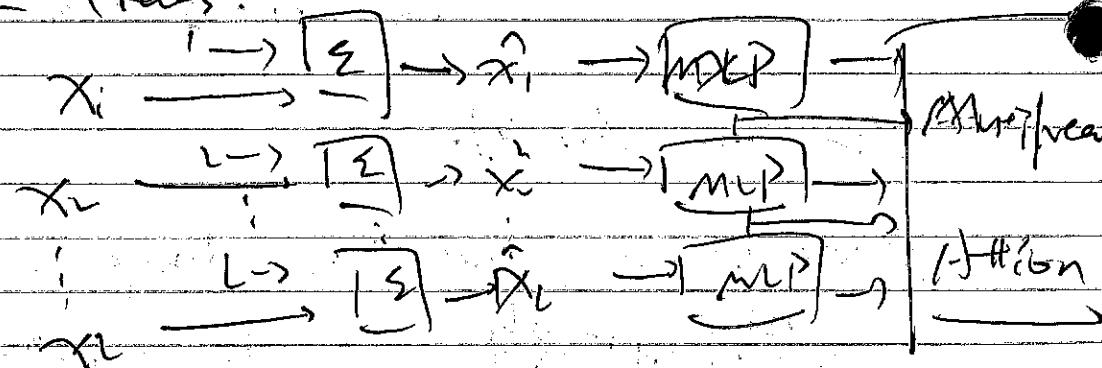
Some input sequence:

↓  
Some no. attention blocks

↓  
diff. weights. → concatenated output

$$\hat{x}_{i,18} = \text{Enc}(x_i, v)$$

- Trans.



$$u_1 \rightarrow [\text{DNN}] \rightarrow y_1$$

$$u_2 \rightarrow [\text{DNN}] \rightarrow y_2$$

$$\vdots \quad \vdots \quad \vdots$$

$$u_n \rightarrow [\text{DNN}] \rightarrow y_n$$

instead, we do:

$$u_1 \quad y_1$$

$$u_2 \rightarrow \oplus \rightarrow [\text{DNN}]$$

$$\vdots \quad \vdots$$

$$u_n \quad y_n$$

## Lecture 20

### Kernels.

opposite of dimension reduction.

map into higher dimensional space.

$$\phi(x) : \mathbb{R}^d \rightarrow \mathbb{V}$$

$$\langle \phi(x_1), \phi(x_2) \rangle = k(x_1, x_2).$$

$$k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

### Examples

#### - Gaussian Kernels

$$k(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|^2).$$

a.k.a. RBF kernel

#### - linear kernel

$$k(x_1, x_2) = x_1^T x_2$$

#### - exponential kernel (laplacian)

$$k(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|)$$

#### - Polynomial kernels

$$k(x_1, x_2) = (1 + x_1^T x_2)^P$$

#### Kernel definition.

$$k : \mathbb{V} \times \mathbb{V} \rightarrow \mathbb{R}$$

$$k(u, v) = k(v, u)$$

→ kernel must be symmetric.

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0.$$

if  $K_{ij} = k(x_i, x_j) \Rightarrow K$  is positive definite

$k$  is a kernel implies:

$\exists \mathbb{V}$  vector space &  $\phi$  s.t.

$$\langle \phi(x_1), \phi(x_2) \rangle = k(x_1, x_2)$$

if  $k_1$  &  $k_2$  are kernels w/ feature

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}$$

$$k(u, v) = k_1(u, v) + k_2(u, v)$$

is a kernel.

$$w_k = \sum_{i=1}^n a_i \phi(x_i).$$

$$k(u, v) = c k_1(u, v) \quad \text{if } c > 0.$$

$$\phi(x) = \sqrt{c} \phi_1(x).$$

$$k(u, v) = k_1(u, v) k_2(u, v).$$

$$\phi(x) = \phi_1(x) \otimes \phi_2(x).$$

$$k(u, v) = f_u(u) k_1(u, v) f_v(v).$$

$$\phi(x) = f(x) \phi(x).$$

$$f: x \rightarrow \mathbb{R}.$$

component loss  $\in f(w, x, y) = l(w^\top \phi(x), y)$ .

$$= l(w^\top \phi(x), y).$$

$$f(w) = \frac{1}{n} \sum_{i=1}^n f(w, x_i, y_i).$$

$$= \frac{1}{n} \sum_{i=1}^n l(w^\top \phi(x_i), y_i).$$

$$\nabla_w f(w, x, y) = l'(w^\top \phi(x), y) \phi(x).$$

$\Rightarrow$  for GD / SGD,

$$w_k \in \text{span}\{\phi(x_1), \phi(x_2), \dots, \phi(x_n)\}.$$

$$w_k = \sum_{i=1}^n \alpha_i \phi(x_i)$$

$$F(a, x, y) = l\left(\left\langle \sum_{i=1}^n a_i \phi(x_i), \phi(x) \right\rangle, y\right)$$

$$= l\left(\sum_{i=1}^n a_i \langle \phi(x_i), \phi(x) \rangle, y\right)$$

$$= l\left(\sum_{i=1}^n a_i k(x_i, x), y\right)$$

$$f(a) = \frac{1}{n} \sum_{j=1}^n l\left(\sum_{i=1}^n a_i k(x_i, x_j), y\right).$$

let  $K_{ij} = k(x_i, x_j)$

$$= \frac{1}{n} \sum_{j=1}^n l(Ka_j, y),$$

① compute  $\phi(x_i)$  on-the-fly as needed.

② precompute & store  $\phi(x_i)$  before training.  $\checkmark$  train in the transformed space

③ compute  $K$  on-the-fly as needed

④ precompute  $K$  & store.

lecture 21

Kernels: function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ .

ways to do learning:

① Learn in the transformed

space, compute  $\phi$  on-the-fly.

② Learn in "kernel trick" space, precompute  $\phi$ .

③ Learn in the "kernel" space

while computing  $k$  on-the-fly.

④ Learn in "kernel trick" space,

precompute  $K_{ij} = k(x_i, x_j)$

$$f_i(w) = l(w^T \phi(x_i), y_i)$$

$$\nabla f_i(w) = l'(w^T \phi(x_i), y_i) \phi(x_i)$$

Suppose  $n$  examples,  $x_i \in \mathbb{R}^d$ ,

and  $k$  takes  $\mathcal{O}(d)$  to compute.

$$\phi(x_i) \in \mathbb{R}^m$$

run for  $T$  steps

Cost to compute  $\phi \Rightarrow \mathcal{O}(md)$ .

Vanilla non-kernel linear model:  $\mathbb{R}^d$

Compute Memory

①  $\mathbb{H}(Td)$

$\mathbb{H}(d)$  (not inc.  
train.s)

②

③

④

Learn in the transformed space

①  $\mathbb{H}(Td)$ , precompute

$\mathbb{H}(M)$ , "trn"

②  $\mathbb{H}(Tm + nTd)$

$\mathbb{H}(nm)$ , (train,d.)

③  $\mathbb{H}(Td)$

$\mathbb{H}(n)$

④  $\mathbb{H}(Tn + n^2d)$

$\mathbb{H}(n^2)$

$$\rightarrow l'(\sum_{j=1}^n w_j \langle c(x_i, x_j), y_i \rangle \phi(x_j))$$

$$\nabla f_i(w) = l'(\phi(x_i)^T \sum_{j=1}^n w_j \phi(x_j), y_i) \phi(x_i)$$

an SGD step is compute i.

$$u_i \leftarrow u_i - \alpha l' \left( \sum_j u_j k(x_i, x_j); y \right)$$

For kernels - subsampling.

$$k(x_i, x_j) \approx \langle \psi(x_i), \psi(x_j) \rangle.$$

$$\psi(x) \in \mathbb{R}^D.$$

Approximate feature map.

Random Fourier features.

$$\text{If kernel: } k(x_i, x_j) = k(x_i - x_j).$$

$$\text{then } k(x_i, x_j) = \mathbb{E} [2 \cos(\omega^T x_i + b),$$

$$\cos(\omega^T x_j + b)]$$

$$b \sim \text{Unif}([0, 2\pi])$$

$\omega \sim f(k)$  Fourier transform of  $k$ .

$$\mathcal{F}[k]$$

RBF kernel:

$$\omega \sim \mathcal{N}(0, 2\gamma I),$$

$$b \sim \text{Uniform}([0, 2\pi]).$$

$$\begin{aligned} & \mathbb{E} [2 \cos(\omega^T x_i + b), \cos(\omega^T x_j + b)] \\ &= \exp(-\gamma \|x_i - x_j\|^2). \end{aligned}$$

$$\text{draw } w_1, w_2, \dots, w_D \sim \mathcal{N}(0, 2\gamma I)$$

$$b_1, b_2, \dots, b_D \sim \text{Unif}([0, 2\pi])$$

$$\text{Set } \Psi_i(x) = \sqrt{\frac{2}{D}} \sum_{j=1}^D w_j (w_j^T x + b_j)$$

$$\Rightarrow \mathbb{E} [\langle \psi(x_i), \psi(x_j) \rangle] = \exp(-\gamma \|x_i - x_j\|^2)$$

i (TBC)

⑤ Learn w/ approx. feature map comp.

• w/ on-the-fly.

⑥ Learn w/ approx. f. m. compare & code

Vanilla linear model.

compute

Memory

(5)  $\hat{H}(\text{TDd})$ .

( $\hat{H}(\text{D})$ )

(6)  $\hat{H}(\text{TD}) + \text{nDd}$  ( $\hat{H}(\text{nD})$ ).

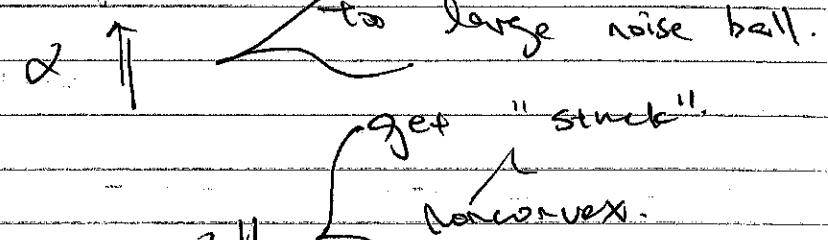
large  $D \Rightarrow$  high "accuracy" in repr.  $K$ .

- high compute & memory cost.

Lesson 22.

▷ Hyperparameter optimization.

Step size / lr /  $\alpha \rightarrow$  diverge.



-  $\lambda$  for regularization

over/under regularize.

converge too slowly. wrong level of expressivity.  
poor generalization.

-  $\beta$  momentum.  $\rightarrow$  if  $\beta \notin (0, 1)$ .

$\beta$   
bad.

-  $B$  batch size.  $\rightarrow$  too small can

affect noise ball.

too large  $\rightarrow$  mem like GA.

performance matching parallel capabilities  
of HW.

choice of architecture: DNN



Generalize poorly.

systems implications,

- - - - - Non-DNN - -

- tree depth → weak learner

overfit & lots of mem.

- feature vector length for

random features.

→ poor fidelity,

high compute/memory

cost.

- # of epochs → high loss/error

compute cost.

- dimension → poor accuracy.

high compute cost

- k in kNN & k-means → under/overfit

more weights for larger k.

the task for assigning hp:

→ HPO - Any process assign hp

Standard:  $\beta$  0.9 0.99

$(\rho_1, \rho_2)$  (0.9, 0.999)

B

power of 2, e.g. 256.

- Grid Search

↔ C,D

Random Search

less reproducible than GS.

Parallelism

Early Stopping

- dropping bad hyperparameters

that perform poor each epoch.

goal: minimize  $F(\alpha, \beta, \gamma, k, \sigma)$

"derivative-free" optimization. (DFO).

Encode our knowledge abt the spec  
of the problem  $\rightarrow$  HPO

Bayesian Opt - Continued.

Mean / Median

$f_0, f_1, f_2, f_3, f_4, \dots$

HPO

$$P(f(x_1) = y_1, f(x_2) = y_2, \dots)$$

$$\text{look at: } P(f(x^*) | f(x_1) = y_1, f(x_2) = y_2, \dots)$$

$$f \in \mathbb{R}^d, f \sim \mathcal{N}(\mu, \Sigma)$$

↓ ↓

mean covariance

\* Fact: if  $(x, y)$  are

jointly Gaussian.

multivariate

{ Set by user

{ According to

intuition.

the  $x | y = y$  is also Gaussian.

\* Fact: "a", then  $X$  is Gaussian.

$$f(x^*) | f(x_1) = y_1, f(x_2) = y_2, \dots$$

$$\sim \mathcal{N}(\mu, \sigma^2)$$

How to pick the next point?

$$x_{\text{next}} = \arg \min_{x^*} a(\mu(x^*), \sigma^2(x^*))$$

$$a(\mu, \sigma) = \mu - k\sigma.$$

(lower confidence bound)

→  $a$  - acquisition function.

$$a(\mu, \sigma) = P(f(x^*) \leq f_{\text{best}}).$$

(Probability of Improvement)

$$f(x^*) \sim N(\mu, \sigma^2)$$

$$= -\Phi\left(\frac{f_{\text{best}} - \mu}{\sigma}\right)$$

↓  
accumulated distribution

$$\rightarrow \Phi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u e^{-\frac{z^2}{2}} dz$$

$$P(f(x^*) \leq f_{\text{best}} | \text{obs}) = P(Z \leq f_{\text{best}})$$

for  $Z \sim N(\mu, \sigma^2)$

$$= P(\text{out}_n \leq f_{\text{best}})$$

for  $U \sim N(0, 1)$

$$= P\left(U \leq \frac{f_{\text{best}} - \mu}{\sigma}\right) = \int_{-\infty}^{\frac{f_{\text{best}} - \mu}{\sigma}} P(u) du$$

$$a(\mu, \sigma) = \mathbb{E}[f(x^*) - f_{\text{best}}$$

$$\min(f(x^*), f_{\text{best}}) | \text{obs}]$$

Expected Improvement.

by convention:

$$= \mathbb{E}[\min(f(x^*) - f_{\text{best}}, 0) | \text{obs}]$$

Thm 3: (Stekely max min principle for heat eqn.)

Assume  $\mu$  is  $C_0^2$  in space &  $C_0^1$  in time, all within  $\bar{U}_T$  & it is continuous on  $\bar{U}_T$ , & it satisfies the heat eqn. in  $U_T$ .

$$\text{then } \underset{\bar{U}_T}{\text{max}} \mu = \underset{\partial U_T}{\text{max}} \mu$$

(7) if  $\Omega$  is connected, & there exists  $\exists (x_0, t_0) \in U_T$ , s.t.  $\mu(x_0, t_0) = \underset{\bar{U}_T}{\text{max}} \mu$

then  $\mu$  is constant in  $\bar{U}_{t_0}$ .

If  $\mu$  attains a max-min @ an interior pt. within  $U_T$ , then  $\mu$  must have been constant for all earlier times.

$\Rightarrow$  If  $\mu$  satisfies  $\mu_t - \Delta \mu = 0$  in  $U_T$ ,

B.C.  $\mu = 0$ , on  $\partial\Omega \times [0, T]$ .

T.C.  $\mu = g$ , on  $\Omega \times \{T\}$ .

where  $g \geq 0$ , then  $\mu(x, t) \geq 0$ . if  $g(x) \geq 0$  for some  $x \in \Omega$ .

Another example of infinite propagation speed of  $\nabla \phi$ . within ht. eqn.

Thm 5 . (Uniqueness in Poisson's Eqn.).

If  $g \in C'(\partial\Omega)$  &  $f \in C'(\Omega)$ ,

then  $\exists$  sol'n  $u \in C^2(\bar{\Omega}) \cap C'(\bar{\Omega})$ .

to

$$\Delta u = f \text{ in } \Omega$$

$$u = g \text{ on } \partial\Omega$$

Thm 6 . (Regularity in Laplace's eqn).

If  $u \in C^1(\bar{\Omega})$  & also satisfies

$$\int_{\partial B_r(x)} u ds = \int_{B_r(x)} u dy \text{ for any}$$

ball  $B_r(x) \subset \Omega$  then  $u \in C^\infty(\bar{\Omega})$

guarantees nothing about what happens on  $\partial\Omega$

Thm 7 (Harnack's inequality)

Let  $\Delta \subset \Omega$  be compact, connected set, &

$\forall v \in \Delta$  a positive constant,  $C$ , for each such  $v$ , where  $\sup_{\Delta} u \leq C \inf_{\Delta} u$

Holds for all non-negative harmonic functions

$v$ , defined on  $\Omega$ .

$\Rightarrow$  The values of a non-negative harmonic function, within an arbitrary  $\mathcal{V}$ , can't be that different from one another: Laplace's equation smooths things out away from the boundary  $\partial\Omega$ .

The point of all these things is to illustrate that, once we have a mathematical model, we can analyze it to understand its properties, & the properties of its solution... and potentially this transfers back to the real world & its properties.

Also, we need to understand the properties of our continuous solutions before we can properly discretize things for the computer.

A good for code debugging.

lets go back to Poisson eqn.

Lemma 2)

Let's recall Poisson eqn.  $\Delta u(x) = f(x)$  & use it to come up with a general viewpoint for thinking about classical solutions to linear PDE

- The Green's function (GF).

$$G(x, \tilde{x})$$

The whole idea comes down to finding a solution for a special instance of, e.g. the Poisson eqn.:

$$\Delta G(x, \tilde{x}) = \delta(x, \tilde{x}).$$

Dirac Delta

The Dirac Delta is kind of linear functional, known as a distribution.

This means  $\delta$  likes to operate on some smooth function to produce a real number. The "operation" involves an integral.

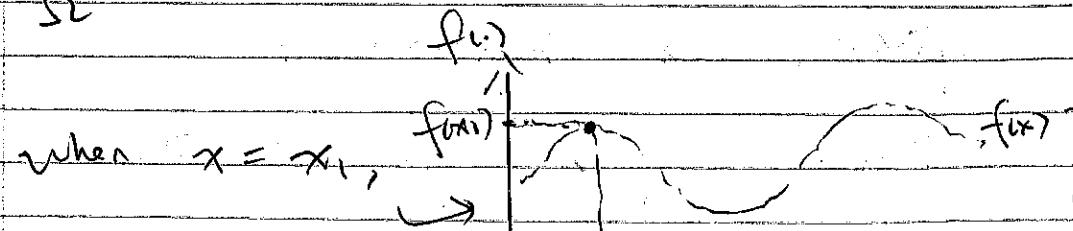
$$\delta(x) = \begin{cases} \infty & x=0 \\ 0 & x \neq 0 \end{cases} \text{ s.t. } \int_{-\infty}^{\infty} \delta(x) dx = 1$$

We can use these interesting properties.

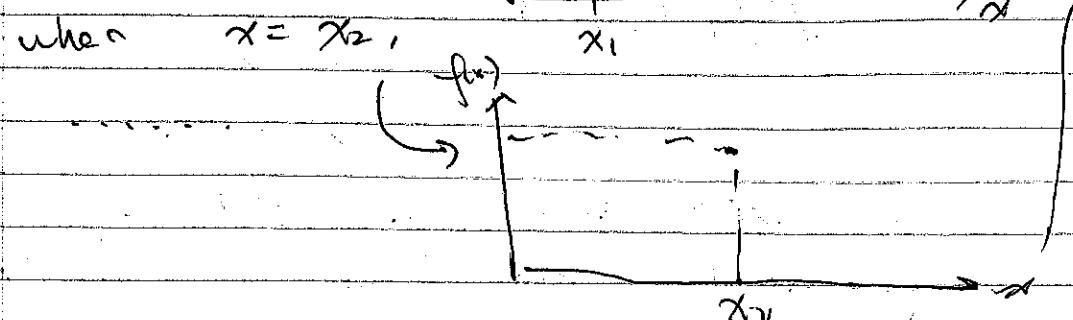
of  $\delta$  to "break up" our source term,

$f(x)$ , into "impulses" meaning:

$$\int_{\Omega} \delta(x, \tilde{x}) f(\tilde{x}) d\tilde{x} = f(x) \quad (+)$$



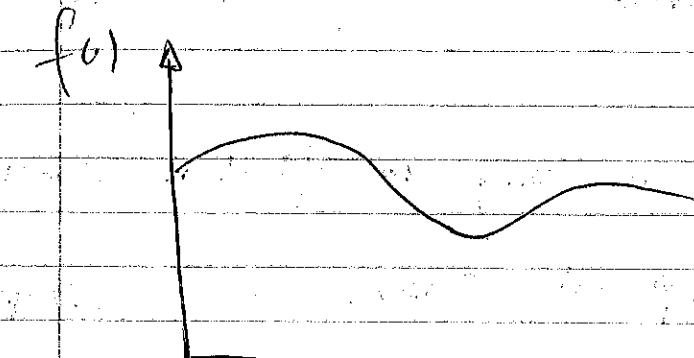
when  $x = x_1$ ,  $f(x)$



when  $x = x_2$ ,  $f(x)$

"taking" an infinity of such pts,

we get a smooth curve"



but we also know that if we substitute a Green's function in  $\Delta(\cdot)$ , we get the Dirac Delta:  $\Delta(G(x, \tilde{x}))$

$$= \delta(x, \tilde{x})$$

$$\text{we rewrite (+): } \int_{\Omega} \Delta(G(x, \tilde{x})) f(\tilde{x}) d\tilde{x} = f(x)$$

Integration is a linear differential operator

$\Omega$  is fixed

$$\Delta \left( \int_{\Omega} G(x, \tilde{x}) f(\tilde{x}) d\tilde{x} \right) = f(x).$$

is a soln to P.S., since it satisfies the PDE  $\Delta(u(x)) = f(x)$

$$\text{thus } u(x) = \int_{\Omega} G(x, \tilde{x}) f(\tilde{x}) d\tilde{x}$$

" $\Omega$  is fixed &  $f$  is nice".

G.F. not so easy to find... available

only for special cases. w/ simple geometries.

e.g. The Green's function for Poisson's

e.g. on the unit dist. w/ homogeneity

Dirichlet B.C.s. i.e.  $(u(x)=0, \forall x \in \partial\Omega)$ , is

$$G(x, \tilde{x}) = \frac{1}{2\pi} (\ln|x - \tilde{x}| - \ln|x^* - \tilde{x}^*| -$$

$$\ln|\tilde{x}|)$$

Heat Eqn. (Diffusion).

$$\frac{\partial u}{\partial t} - \Delta u = 0$$

if there is a source term within the domain, eq. becomes  $\frac{\partial u}{\partial t} - \Delta u = f$ .

where  $t > 0$  &  $x \in \Omega \subset \mathbb{R}^n$ , the solution  $u(x, t)$  is the function satisfying the PDE & I.C. & B.C.

such that  $u: \Omega \times [0, \infty) \rightarrow \mathbb{R}$ , i.e.

our sol'n is a scalar field over  $\Omega$ .

But where did that PDE come from?

Let's begin from the case with no flow:

flowing: First we need some terminology.

Thermal Energy = instantaneous molecular kinetic energy

Heat = Thermal Energy in transit.

Temperature = Average of thermal energy

Physics Theory = Conservation of Energy

Rate of change of Heat Generated by  
in Thermal Energy = sources within  $\Omega$ : f  
within  $\Omega$   
taking arbitrary subd.  $\Omega' \subset \Omega$  + Flow of heat into  $\Omega$ ,  
from outside

Note: Thermal Energy contained in a

body is tied to "heat capacity",  $C$ ,

$C$  = amount of heat generated per unit mass, per rise in temp.,  $\mu$ .

$\rho$  = mass density.

The total thermal energy within  $\Omega'$  @ a given temp.  $t$  is:

$$\int_{\Omega'} C(x) \rho(x) \mu(x, t) dx$$

temp. field

If  $f(x, t)$  is our heat source, specified per unit volume, per unit time, then the

total heat generated by the source within

$$\Omega' \text{ is: } \int_{\Omega'} f(x, t) dx$$

Lastly, heat flux thru  $\partial\Omega'$  is  $q(x, t)$ ; thus

the total heat flux, the body @ time "t", is given by:

$$\int_{\partial\Omega'} q(x, t) \downarrow \text{unit outward normal to } \partial\Omega' ds.$$

Rate of change

in thermal energy

within  $\Omega'$

considering arbitrary sub-domain  $\Omega'$  + flow of heat into domain

$$- \int_{\partial\Omega'} q \cdot \nabla \cdot \mathbf{V} ds$$

↓ divergence

$$\frac{d}{dt} \int_{\Omega'} C(x) \rho(x) \mu(x, t) dx$$

$$= \int_{\Omega'} f(x, t) dx - \int_{\Omega'} \nabla \cdot q dx.$$

$$\int_{\Omega'} \left[ C(x) \rho(x) \frac{\partial \mu(x, t)}{\partial t} - f(x, t) + \nabla \cdot q \right] dx = 0$$

not the same.

$$\Rightarrow Cp \frac{\partial \mu}{\partial t} + \nabla \cdot q = f \quad (\#)$$

we now need

a constitutive model (Fourier's law), relating  
heat flux w/ gradients in temp.

$$q = -k \nabla \mu \quad (*)$$

thermal conductivity.

(\*)  $\rightarrow$  (#)

$$\frac{\partial \mu}{\partial t} - \left( \frac{k}{Cp} \right) \nabla \cdot \nabla \mu = f$$

$$\frac{\partial \mu}{\partial t} - \alpha \nabla^2 \mu = f$$

thermal diffusivity

(RHS=0)

The fundamental sol'n for homogeneous  
heat eqn.

$$\phi(x, t) := \begin{cases} \frac{1}{(4\pi t)^{1/2}} e^{-\frac{|x|^2}{4t}} & (x \in \mathbb{R}^n, t > 0) \\ 0 & (x \in \mathbb{R}^n, t \leq 0) \end{cases}$$

Note singularity @ pt. (0, 0)

$$|x| = (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2}$$

(another radial sol'n)

Sol'n is normalized s.t.

$$\int_{\mathbb{R}^n} \phi(x, t) dx = 1$$

let's consider IVP heat equation:

$$\mu_t - \Delta \mu = 0 \quad \text{in } \mathbb{R}^n \times (0, \infty)$$

$$\mu = g \quad \text{on } \mathbb{R}^n \times \{t=0\}$$

we note that the fundamental sol'n,  $\phi$ , solves

the homogeneous heat eqn. away from  
the singularity @ (0, 0), as it does for any  
shift in spatial coordinate, as in  $\phi(x-y, t)$

for each fixed  $y \in \mathbb{R}^n$ , As a result,

the following convolution is a solution to our

Thm 1: (Initial value heat eqn.)

Assume  $g \in C(\mathbb{R}^n) \cap L^\infty(\mathbb{R}^n)$ , and  
define " $\mu$ " as,  $\mu(x, t) = \int_{\mathbb{R}^n} \phi(x-y, t) g(y) dy$

$$= \frac{1}{(4\pi t)^{n/2}} \int_{\mathbb{R}^n} e^{-\frac{|x-y|^2}{4t}} g(y) dy.$$

$(x \in \mathbb{R}^n, t > 0)$

then,

1)  $\mu \in C^\infty(\mathbb{R}^n \times (0, \infty))$  c-sol'n & smooth

2)  $\mu_t(x, t) - \nabla^2 \mu(x, t) = 0, \quad x \in \mathbb{R}^n, t > 0$ .

&

3)  $\lim_{(x,t) \rightarrow (x^0, 0)} \mu(x, t) = g(x^0)$  for each  $x^0 \in \mathbb{R}^n$ .

Since  $g$  is bounded & continuous, if we

also make it positive semi-definite, then

$\mu(x, t)$  is positive for all Pts.  $x \in \mathbb{R}^n$

& times  $t > 0$ .

implies prop. speed for  
info. regarding I.C.

$$\mu_t - \Delta \mu = f \quad \text{in } \mathbb{R}^n \times (0, \infty)$$

$\mu = 0 \quad \text{on } \mathbb{R}^n \times \{0\}$

Thm 2 - (forced heat eqn.)

Let  $\mu(x, t)$  be of the form

$$\mu(x, t) = \iint_0^t \int_{\mathbb{R}^n} \phi(x-y, t-s) f(y, s) dy ds.$$

$$= \int_0^t \frac{1}{(4\pi(t-s))^{n/2}} \int_{\mathbb{R}^n} e^{-\frac{|x-y|^2}{4(t-s)}} f(y, s) dy ds$$

treating our forcing as I.C. @ "s"

then (1)  $\mu \in C_c^2(\mathbb{R}^n), C_c([0, \infty))$ .

(2)  $\mu_t - \Delta \mu = f \quad (x \in \mathbb{R}^n, t > 0)$ .

and

(3)  $\lim_{t \rightarrow 0} (\mu(x, t)) = g(x)$ .

(b)  $\mu(x, t) = 0$  for each  $x \in \mathbb{R}^n$

Combining Thms 1 & 2, we can solve

$$u_t - \Delta u = f \text{ in } \mathbb{R}^n \times (0, \infty)$$

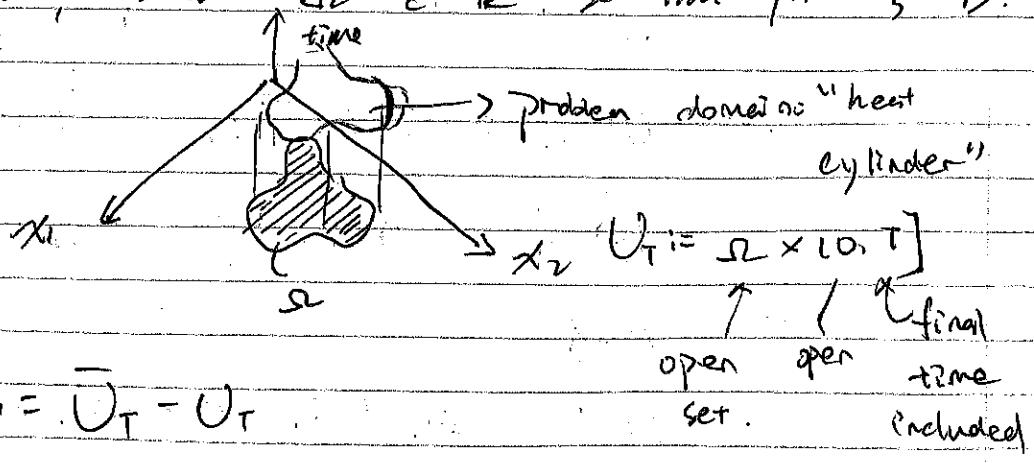
$$u = g \text{ on } \mathbb{R}^n \times \{0\}$$

with  $u(x, t) = \int_{\mathbb{R}^n} \phi(x-y, t) g(y) dy$

$$+ \iint_{\mathbb{R}^n} \phi(x-y, t-s) f(y, s) dy ds$$

moving away from all  $\mathbb{R}^n$ , now, let's define an open region of space serving as our spatial domain:  $\Omega \subset \mathbb{R}^n$ .

For now, draw  $\Omega \subset \mathbb{R}^n$  & time for 3'.



$$\partial U_T := \bar{U}_T - U_T$$