# PERSONAL NOTES

# PRINCIPLES OF LARGE SCALE ML

Hanfeng Zhai

2023

Scaling $\longrightarrow$ ML
$\uparrow$
principles.

How scale impact performance of ML system.

plu Data $\longrightarrow$ Data $\longrightarrow$ Feature $\longrightarrow$ (Robust).
Collection     cleaning.   Engineering           Training
                          (Embedding).

Train/Valide/Test.
Spliting.

Model selection
(Hyperpara. Opt.)

(Online leaning)          Deployment.

"tilga test loss"

Monitoring. $\longleftarrow$ Inference.

Optimization   Statistics.
        ML
    Systems.

Principle #1: Optimization. Tasks

Solve it using gradient-base algorithms
there are fast & "canned".

backpropagation + numerical linear algebra.

o gradient descent.

o empirical risk min.

Principle #2: X whole dataset, → subsample

o Stochastic gradient d.

  ↳ subsample the component losses.

✩ compute ⟩ faster.

o Cross-validation / train val test split

o bagging.

o kernel sampling.          → subsample with ___
                            in CNN

o dropout.              o data augmentation.
--- VAE. (dimension reduction).

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

o batch size. Selection.

o numerical precision.

o GPU ↪ NN.

o Caching

o data loader / streaming

                ↑
                GPU.
o using fast numerical linear algebra kernels.

parallel / distributed computing.

o neuromorphic computing. (chip → NN).

cs.cornell.edu/courses/9547817.

+ canvas.              + Gradescope → problem.

+ ed discussion.       + CMS - programming a.
                       (python, numpy, torch)

Grading.

- Psets

- PAs.

- 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 as optimization.

✱✱✱ $\downarrow$
continuous optimization.

$\downarrow$

optimizing real numbers.

Data.

$\hookrightarrow$ Embedding steps $\longrightarrow \mathbb{R}^d$ vectors.

e.g. pred. ppl.

• 33

• g. woq.

• 2013 . current job / cornell dun.

• 58,000 / yr

$\hookrightarrow \begin{bmatrix} - \\ - \\ - \end{bmatrix} \Rightarrow \mathbb{R}^4.$

Vector - element in vector space.

$x, y \in V$   $x + y \in V$

+ associated & commutative.

multiplication

~~$\cancel{XXX}$~~ 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 $\mathbb{R}^{n \times n}$. $\Big\}$ can all be
$\mathbb{R}^{n \times n \times p}$. interpreted as
vectors.

Core idea under: $\underline{\text{basis}}$ in linear algebra.

how many numbers does it
take / pinpoint 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 $x_1, x_2, \dots x_n \in B$,
and $a_1, a_2, \dots a_n$.

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

and $\forall x \in B$,

$|B| = $ "Dimension"     $x \notin span(B / \{x\})$

$\downarrow$
terminology confusion

$\swarrow$ $\searrow$ numpy, sympy, ...
Numerical linear algebra

$V = \mathbb{R}^2$.

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

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

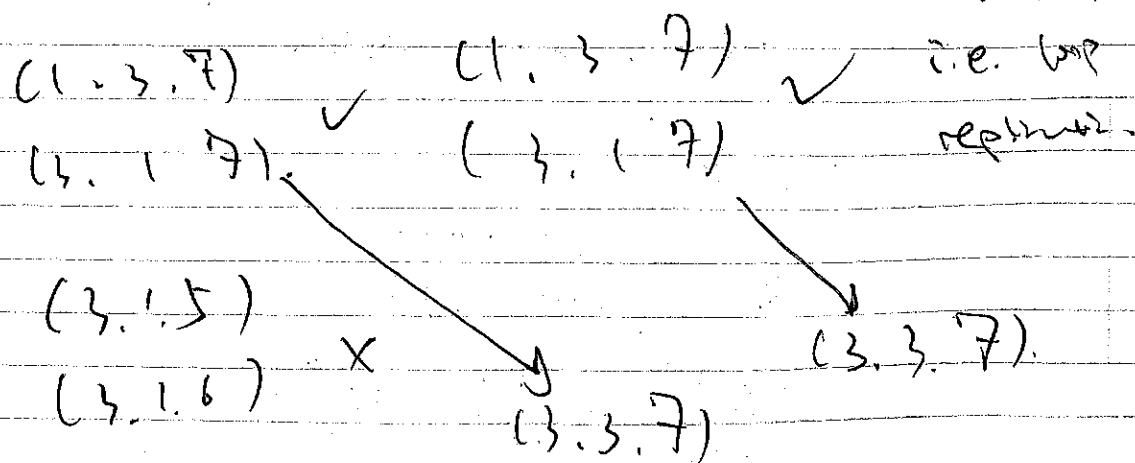numpy faster python.

- multiprocessing
- c++
- cache

* & @ easily masked

↓ multiply per element.  ↓ tensor / matrix multiply

broadcast → faster than manually comp.

(1,3,7) ✓  (1,3,7) ✓  i.e. top replicated
(3,1,7)   (3,1,7)

(3,1,5) ╲  
(3,1,6) ✗ (3,3,7)   (3,3,7)
       (3,3,7)

~~$\sum$~~ numpy broadcasting.

---

Automatic diff.

Definition — $\nabla f(w)^T$ is a linear map such that: $\Delta \in \mathbb{R}^d$

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

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

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

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

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

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

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

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

all elements.

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

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

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

$$f(x + a\Delta) = (x + a\Delta)^T (x + a\Delta)$$
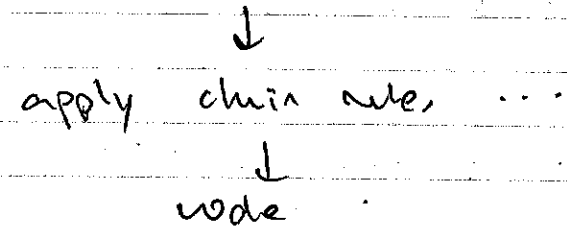
Euclidean norm

$L_1$ 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, ...
↓
code.

→ Problems.

had to do it manually.

○ converting code into mathe   not trivial

○. no guarantee → comput. structure.

why important?

differenting $f_E$ function.

→ Numerical differentiation.

○ high-order different.

○

○ error accumulating

Problems.

- numerical imprecision.
  ↳ number of operations increases.

- f not smooth → wrong results.

- unclear how to setup ε

- for a vector → scalar function,
  you have to compute each partial
  individually, meaning $O(d)$ blowup in cost.

- if ε 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

- Blow-up proportional to $d$.
                          dimension
  ∧
  if we are computing a gradient of f.
  $\mathbb{R}^d → \mathbb{R}$.

Week 2 - 2
Aug 31.

Backpropagation.

Review for Forward Mode AD

$x \in \mathbb{R}$

$\downarrow$

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

⚡ Python supports operator overloading

Reverse - Mode AD

$\rightarrow$ fix one output $\ell$ over $\mathbb{R}$.

$\hookrightarrow$ compute partial derivative $\frac{\partial \ell}{\partial y}$.
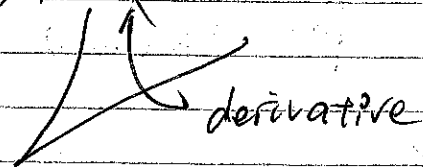
$(y, \nabla_y \ell)$.

Same shape

- Derive backprop. thru chain rule.

$$\ell = f(u), \quad u = g(y)$$

$$\mathcal{L} = P(u_1, u_2, \ldots, u_k), \quad u_i = g_i(y)$$

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

derivative

is a Jacobian

compute gradient w.r.t. $y$

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

process $y$     process $u_i$

$$z = x + y$$

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

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

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

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

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

$$\nabla_x H = ReLU^{-1}$$

initialize

$$\nabla_{U_2} H = \nabla_H H = \nabla_W H = \nabla_x H = \nabla_b H = 0$$

$$\nabla_{U_2} H+ = ReLU'(U_2) \cdot \nabla_H H = 1 \cdot 1 = 1$$

## Problems

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

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

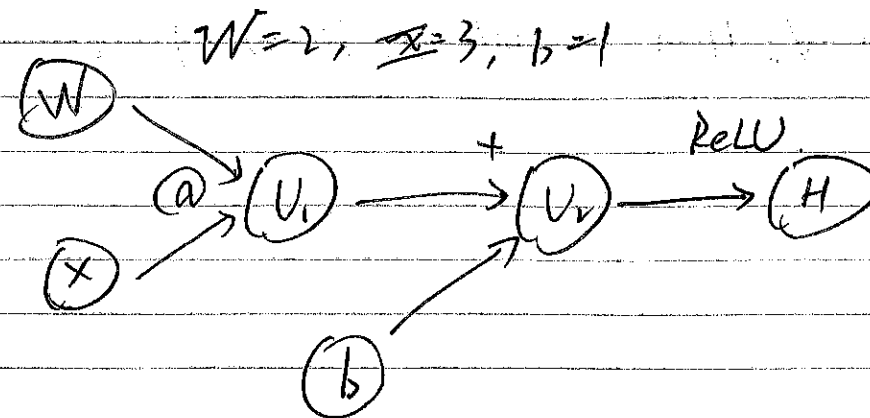$\qquad\downarrow \qquad \bigg) \; why \quad x^T A \Delta = (x^T A \Delta)^T?$

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

d. $\rightarrow$ a Jacobian.

2.5. why $sum(log)$ in np.

Review: Backprop.

$W=2, \; X=3, \; b=1$



$\bullet \; U_1 = W * X \cdot = 2 * 3 = 6$

Forward $\bigg| \; \bullet \; U_2 = U_1 + b \cdot = 7$

$\bullet \; H = ReLU (U_2) = ReLU (WX + b) = ReLU(7)$

$\bullet \; \nabla_H H = 1$

$\bullet \; initialize: \quad \nabla_{U_2} H = \nabla_{U_1} H = \nabla_W H$

$\qquad\qquad\qquad = \nabla_x H = \nabla_b H = 0$

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

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

$\bullet \; \nabla_b H \; += \; \nabla_{U_2} H = 1$

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

$\circ \ \nabla_x H + = W \ \nabla_{u_1} H = 2 * 1 = 2$

$\overline{F}(x) \longrightarrow D F(x)$

$0 \longmapsto$

$A^{-1}$

Machine learning systems

$f : U \longmapsto V$

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

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

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

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

if $F(x) = x^2$

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

$\lim_{\alpha \to 0} \dfrac{F(x + \alpha \Delta) - F(x)}{\alpha}$

$\lim_{\alpha \to 0} \dfrac{x^2 + \alpha \Delta x + \alpha x \Delta + \alpha^2 \Delta^2 - x^2}{\alpha}$

$= \Delta x + x \Delta$

$DF(x) = (\Delta \longmapsto \Delta x + x \Delta)$

or $\left( DF(x) \right)(\Delta) = \Delta x + x \Delta$

Why Backprop.

gradient has 1 output $\ell$.

Backprop manifest gradient

$\nabla_y \ell$ for intermediates $y$.



class / object
↓
knows how
↳ compute itself
↓
knows its dependencies

def __add__(self, other):

two ways:                    (lazy)

EAGER    v.s.    GRAPH (only later)

compute immediately    don't compute, but
                       know the dependencies

Eager                    Graph

1 cutoff                 ○ first computation

1 error                  ○ simplization
                         (transform. ✓)
○ take storage
                         ○ reuse the graph.
○ debugging &            (✓)
development (✗).
                         ○ good for caching
○ if / branch / control flow    (?)
(✓).
                         ○ memory locality
                         (✓).

                         ○ less python
                         overhead (✓)

                         ✦ Deployment.
                         ↓
                         convert it
                         into some form
                         that "inpython"
                         — optimization

TF $\xrightarrow{\text{Default}}$ Graph

PT $\longrightarrow$ Eager

immuterialized

ML framework.

> numerical linear algebra library

> hardware accelerator support.

o backprop engine. (gradient. auto).

> library for writing Deep Nets

PyTorch & TensorFlow

MxNet    Jax.    Flux (Julia)
(c++)

Caffe (ancient) → 1st GPU acc.

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 hardwares.

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.

(1.1)

use graph
search alg.

\* def all_dependencies( ) 
find all the dependencies
(1.2)

\* def backward( ) → Sort the dependencies
in reverse node
backward → find gradients

def grad_fn ( )

---

"Define math operators in backprop"

def add
▷ BA_Add( )

def sub
▷ BA_Sub( )

direct modified
for tensor operation
(2.2)

def mul
▷ BA_Mul( )

def truediv
▷ BA_Div( )

def sum ( )
▷ BA_Sum( ).

Tensor operation
(2.1)
(2.3)

def reshape ( )
▷ BA_Reshape ( )

Compute "grad"

def transpose ( )
▷ BA_Transpose( )

implement for

→ Different classes          scalar - arrays.
(1.3)

compare ad, nd, sd

$(1.5)$

def numerical_diff( )

global funs.

def numerical_grad( )

def backprop_diff( )

$(2.3)$ $(2.4)$ $(2.5)$ $(2.6)$

Store test functions in a class:

class TestFxs( ):

  def {f₁, df₁dx, f₂, ......}

test implementation: Test Fxs. df₃dx $(1.4)$

High d input

if __name__ == "__main__":

Scalar output

write the script to execute

Compare stuff Main functions.

---

for #3:

$$\lim_{\alpha \to 0} \frac{\log(x+\alpha\Delta) - \log(x)}{\alpha} = D(F(x))\Delta$$

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

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

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

$$=$$

$$=$$

Week 4 - Lecture 7.

GD

we compute gradient in large-scale
opt. problem to solve learning task
examples will most be on supervised
learning. but not limited to ...

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

$= \frac{1}{n} \sum_{x \in I} f(w; x)$  size of dataset.

in framework, we strive to:

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

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

Init $w_0 = 0 \rightsquigarrow$ not nec. zero. gradient.

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

✭ How much time does it cost?

$\triangleright \, O(ndk)$.

✭ How much memory required?

$\rightarrow O(nd)$       or      $O(d)$

                    $\downarrow$
           • parallel way

running backprop thru whole thg.
in a NAIVE way $\mapsto O(nd)$.

___

Newton's Method

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

↳ converges faster.

"BUT" more expensive that GD.    Store Hessian matrix.

So, ✭ How much time? & ✭ How m. Mem?

$\triangleright O(nd^2k)$.          $\rightarrow O(d^2)$

(1) I scared away from Newton's method cuz the large memory required, i.e. quadratic cost.

(2) why are we confident that GD will converge?

the value of gradient are computed based on particular weight, if move a lil' bit, it will drastically change grad

we want to make it a bust.

How close GD ~?

$\Rightarrow$ bounds Derivative.
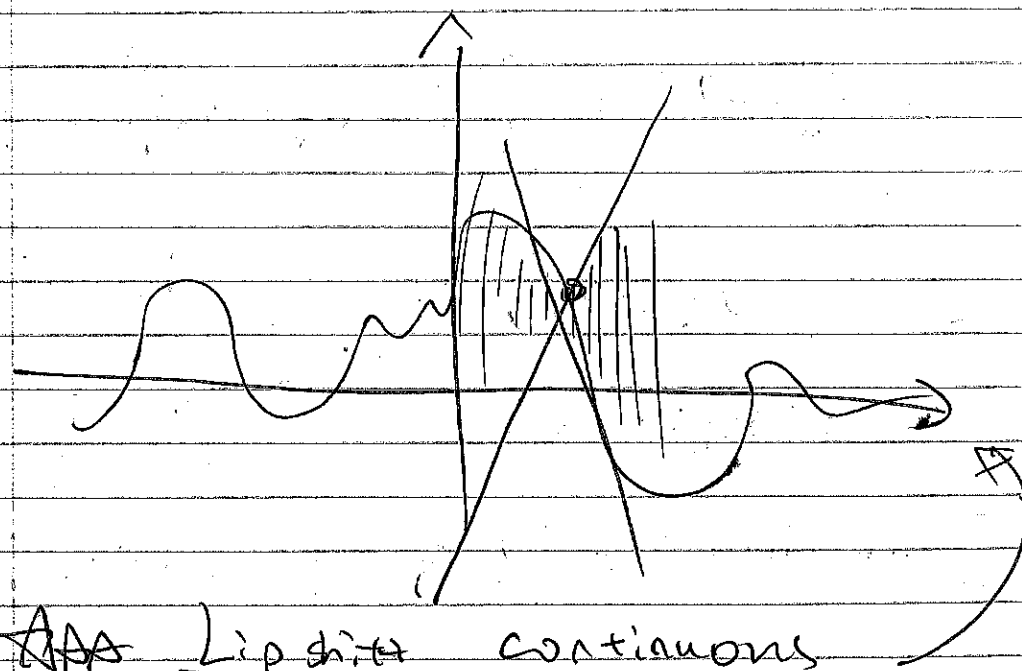
$\quad\hookrightarrow$ Second-order D. is also Bd.

Gradient $\nabla f$ is $L$- Lipnitz continuous

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

Assume $\exists f$ r.t. $f(w) \geq f*$.

$\quad\hookrightarrow \forall w, \forall u, \|u\|=1.$

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



AAA Lipshit continuous

$\quad\hookrightarrow L$- smooth.

How to prove GD converges ???

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

$$\text{fundamental thm. of calc.}$$

$$\Downarrow$$

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

$$\Updownarrow$$

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

$$= f(w_t) + \int_0^\alpha \left(-\nabla f(w_t)\right)^T \nabla f(w_t)$$

$$\eta \nabla f(w_t)) \, d\eta. \qquad \text{how?}$$

$$\hookrightarrow \text{imagine} = 0 \qquad \text{close to grad.}$$

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

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

Expect $\eta$ to be small $\leftarrow$ expect

$$d \cdot \text{small}$$

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

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

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

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

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

$$\Rightarrow \text{close to norm.} \to \text{always}$$

$$\text{positive} \Rightarrow \underline{gd} \text{ works} \to \underline{min}$$

$$\to d \text{ coat.}$$

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

$$+ \int_0^\alpha \left( -\nabla f(w_t) \right)^T \left( \nabla f(w_t - \eta \nabla f'(w_t)) - \nabla f(w_t) \right) d\eta$$

(Apply Cau. theorem:

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

$$\{ \|A \, B\|. \}$$

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

$$+ \int_0^\alpha \| -\nabla f(w_t) \| \, \| \nabla f(w_t - \eta \nabla f(w_t)) - \nabla f(w_t) \| \, d\eta$$

$$\leq f(w_t) + \int_0^\alpha -\|\nabla f(w_t)\| \, d\eta$$

$$+ \int_0^\alpha \| -\nabla f(w_t) \| \cdot L \cdot \| \eta \nabla f(w_t) \| \, d\eta$$

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

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

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

⇒ How does constant L affect how we learn tasks in scales

that's why

assume: $\alpha L < 1$

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

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

$\rightarrow$ Single iteration of GD.

$$\frac{\alpha}{2} \sum_{t=0}^{k} \|\nabla f(w_t)\|^2 \le \sum_{t=0}^{k-1} \left( f(w_t) - f(w_{t+1}) \right)$$

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

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

( assume $\underline{f \text{ is bounded from below}}$ ✱✱✱

$\qquad = \le f(w_0) - f^*.$

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

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

$\rightarrow$ require $\alpha L \le 1$.

LHS run to iteration of GD
return the optimal weights.

RHS upper bound of func.

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

the largest step

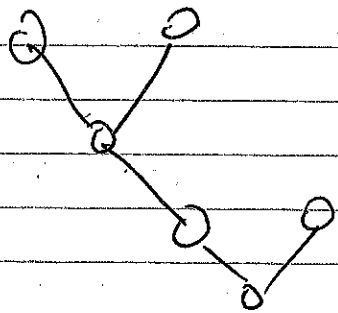$\quad \rightarrow \le \frac{2L(f(w_0) - f^*)}{k}.$

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).

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

confused on the implementation.

② writing a for loop.

while the len

① add empty row/cols. for
extra dim. → new axis.
check the dims are targeting the same
↳ broadcasting.

A: $1 \times 6 \times 5$

B: $\boxed{1} \times 2$

→ iterating over the axis.

1: adding 1.

2: repeating axis (iterating)

helper_func

BA – And

Just recall  += helper (a, b)

---

HW Review:

$\nabla_{\underline{X}} \, \ell$ has same shape as $\underline{X}$

$X$. grad. shape = $X$. data. shape.

* "+=" in numpy mutates the vector

self . X . grad += stuff

self . X . grad = self X grad . + stuff.
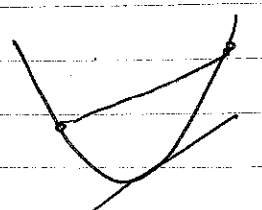
☆ in't both as arrays as all zeros

---

GD continued.

... last time.

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

→ most simple condition – convexity

0$^{th}$ – order convexity

$$f((1-\alpha)x + \alpha y) \le (1-\alpha)f(x) + \alpha f(y)$$

for $\alpha \in [0, 1]$

▽ 1st order def. ──── considering whether func. convex.

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

▽ 2nd order def.

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

"parabola has a const. sec. deriv."

— convex has a unique global optimum.

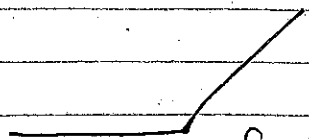"that's why we care 'bout con"

Strongly convex function

"a function that can be bounded from below"

μ - strongly convex function.

↳ scalar positive num.

$$\frac{\partial^2}{\partial \alpha^2} f(x + \alpha u) \geq \mu \quad (\text{for } \|u\| = 1)$$

E.g. "Strongly convex"

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

---

Convexity: vector space ↦ real num.
i.e. scalar.

Strictly convex. vs. Strongly convex

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

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

if larger → strongly convex

parabola

f is μ-strongly convex, iff ∃ g s.t.

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

& g is convex.

✦✦✦✦

Use in a lot of analysis.

Polyak - Lojasiewicz (PL)

··· TBC

$$\|\nabla f(\omega)\|^2 \geq 2\mu \left( f(\omega) - f^* \right)$$

min val.

strongly convex.

$\underset{?}{\downarrow}$ does not imply strongly convex

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

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

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

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

☆ the gradient is only close to zero when its only around global optimum.

satisfies PL yet not strongly cv.

---

$$f(\omega_{t+1}) = f(\omega_t - \alpha \nabla f(\omega_t))$$

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

non-general convex of gd.

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

$$f(\omega_{t+1}) - f^* \leq \left( f(\omega_t) - f^* \right) - 2\mu\alpha$$

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

$$\leq \left[ \frac{1}{2} \mu \alpha \left( 1 - \frac{\alpha L}{2} \right) \right] \left( f(\omega_t) - f^* \right)$$

call it $\gamma$

$\hookrightarrow$

$$\leq \gamma \left( f(\omega_t) - f^* \right) \quad \underset{\text{min}}{\longleftarrow} \quad \text{min is } \gamma$$

$$f(\omega_k) \quad \bigotimes - f^* \leq \gamma^k \left( f(\omega_0) - f^* \right)$$

if $\gamma > 0$.

☆ How should we chose $\alpha$?

to show $1 - \frac{\mu}{L}$ is non-ne...

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

$$1 - \alpha L = 0$$

$$\alpha L = 1$$

$$\alpha = \frac{1}{L}.$$

we got $\gamma = 1 - 2\mu \frac{1}{L}(1 - \frac{\frac{1}{L} L}{2}) = 1 - \frac{\mu}{L}$

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

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

convergence @ a linear rate

even tho is decreasing exponentially
converges with errors a num. of digits.
going to be linear wrt. . . .
"roughly linear with digits wrt. score".

Goal: output $\hat{w}$ with

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

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

∴ suffices to run $T$ iterations of
(num timestep)
(GD), s.t.

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

we can solve this for $T$

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

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

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

hessian
for s.c.
pos. def.
symm.

ratio of
largest eigen
val of hessian
/ the smallest val.

$$\frac{L}{\mu} = K \quad \text{"condition number"}$$

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

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

$$(\text{for } \|u\|=1)$$

eqly. $-LI \leq \underbrace{H f(x)}_{\perp} \leq LI$

$$\nabla^2 f$$

$k$ large $\rightarrow$ lot of steps

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

$k \gg 1 \hookrightarrow k=1$ isotropic
quadratic

Condition number $\propto$ num of steps
to an
to solve the
opt. problem.

---

the other thing impact time o GD.

$\rightarrow$ size o data set.

Stochastic $\leftarrow$ GD goes bad when

GD $\qquad$ DS are large.

Subsampling $\rightarrow$ approx.

Principle #2 of LS ML.

Spoiler: g- faster than GD

GD # steps $\sim \vartheta(K)$.

$$\underset{GD}{\text{SGD}} \text{ \# } \overset{time}{\text{steps}} \sim \vartheta(\log K)$$

batch
size

Lecture 7

Week 3

Subsonptime

GD time $\quad \mathcal{O}\left(\widehat{nK}\right)\log\left(\frac{f_0-f^*}{\varepsilon}\right)$

training set size $\nearrow$

$I$

num. of iterations
of GD.

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

GD $\begin{cases} W_{t+1} = W_t - \alpha \nabla f(W_t) \\ \\ = W_t - \frac{\alpha}{n}\sum_{i=1}^{n} \nabla f_i(W_t) \end{cases}$

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

$$W_{t+1} = W_t - \alpha \nabla f_i(W_t).$$

mini-batch SGD: Sampling $B$

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

batch size

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

Stochastic GD (mini-batch)

$\downarrow$

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

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

$$= \sum_{j=1}^{n} \frac{1}{n} \nabla f_j(W_t).$$

$$= \nabla f(W_t).$$

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

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

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

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

$$f(W_{t+1}) = f(W_t - \alpha g_t) = f(W_t) + \int_0^\alpha \frac{d}{d\eta} f(W_t - \eta g_t)\, d\eta = f(W_t) + \int_0^\alpha \langle \nabla f(W_t - \eta g_t)$$

full loss function $(-g_t)\, d\eta$

$$= f(w_t) - \alpha \nabla f(w_t)^T g_t + \int_0^\alpha [\nabla f(w_t - \eta g_t)$$
$$- \nabla f(w_t))^T (-g_t) \, d\eta$$

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

$$\leq f(w_t) - \alpha \nabla f(w_t)^T g_t + \int_0^\alpha [L\|\eta g_t\| \, \|g_t\|] \, d\eta$$

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

full-batch grad

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

$$+ \frac{\alpha^2 L}{2} \|g_t\|^2 \qquad \nabla f(w_t) \; \text{stochastic grad}$$

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

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

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

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

$$\leq f(w_t) - \alpha \|\nabla f(w_t)\|^2 + \frac{\alpha^2 L}{2} \mathbb{E}[\|g_t\|^2 | w_0]$$

New Assumption   uniform over $w_t$

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

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

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

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

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

$$\frac{1}{B}\left[\underline{\phantom{xxxx}}\right]$$

samples that $b \neq c$   independent each other   $b \neq c$

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

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

same

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

Back to the assumption:

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

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

Easy to bound
$$\|\nabla f(w)\|^2$$

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

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

$$= \|\nabla f(w_t)\|^2 + \frac{\sigma^2}{B}$$

coming back to 2d assumption

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

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

Variance of
the square error of SGD

---

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

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

PL condition

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

$$\mathbb{E}[f(w_{t+1})-f^* \,|\, w_t] \leq \left(1-\alpha\left(1-\frac{\alpha L}{2}\right)2\mu\right)(f(w_t)-f^*) + \frac{\alpha^2\sigma^2 L}{2B}$$

SGD

additional "noise" term

$$\leftarrow \rho_t = \mathbb{E}[f(w_t)-f^*]$$

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

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

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

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

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

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

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

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

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

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

step:

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

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

$$+ \frac{\alpha^2\sigma^2 L}{2\mu B} \longrightarrow \text{noise wall/noise floor}$$

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

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

Lecture 9.    Week 5 -2.

Rev - last time:

$$\mathbb{E}\left[f(w_T) - f^*\right] < \exp(-\alpha \mu T)(f(w_0)$$

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

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

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

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

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

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

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

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

$$\Rightarrow T \geq \frac{1}{\alpha \mu} \cdot \log\left(\frac{2(f(w_0) - f^*)}{\varepsilon}\right)$$

also $\alpha L \leq 1$.

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

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

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

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

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

$$\left. K \cdot \log\left(\frac{1}{\varepsilon}\right)\right)$$

usually, $\varepsilon$ are small

$$\widetilde{\mathcal{O}}\left(\frac{\sigma^2}{B\mu\varepsilon} K\right) \longrightarrow \text{batch } B \text{ cancels}$$

$\hookrightarrow$ ignoring "log" term

runtime of GD. $\vartheta\left(nk\log\left(\frac{1}{\varepsilon}\right)\right)$ ⑥

for GD $n \to \dfrac{\sigma^2}{B\varepsilon}$.

<u>GD</u>　　　　　　<u>SGD</u>

- convex (strongly)　　　∘ non-convex

- Global opt. guaranteed　∘ multi-objective
  (unique).

　　　　　　　　　　　∘ huge dataset

∘ large memory (GPU).　∘ huge $n$
parallel

　　　　　　　　　　　$\varepsilon \to$ large "ish"

∘ really small $\varepsilon$　　or not too small

∘ large $\sigma^2$

How to set batch size $B$
　　　　　　　　　　$\Downarrow$
　　　★ Hardware consideration.

---

⓪

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

$\nabla f_i(w) = x_i \overbrace{(x_i^T w - y_i)}$ → dot  axpy

minibatch $\underline{X} \in \mathbb{R}^{B\times d}$ 　　$y \in \mathbb{R}^B$

$\nabla f_{batch}(w) = \dfrac{1}{B} \underline{X}^T(\underline{X}w - \underline{Y})$

faster　　　　MMPY　$(B,d) @ (d,1)$,
　　　　　　　　　　　　$(d,B) @ (B,1)$

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

minibatch size $\propto$ hardware itself

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

SGD $\longrightarrow$ GD: sampling w/o replacemt.
　　　　max$(n, 1)$ → is GD
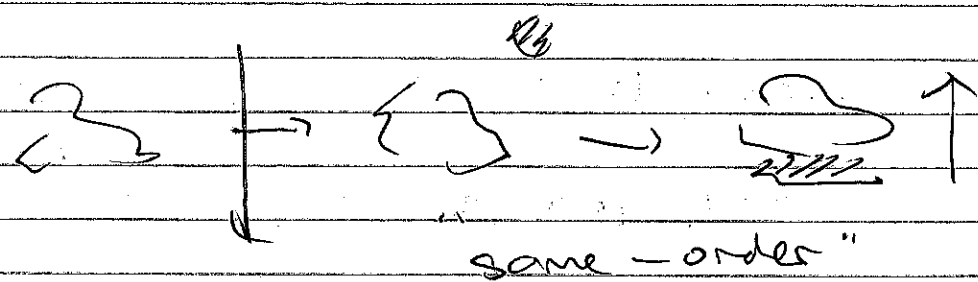　　　larger, slower　　　scenario

<u>random reshuffling</u>

sample w/o rep., only reuse

example after every example

has been used.

epoch $\Rightarrow$ 1 past thru the training set

(e.g. $100 \sim 200$).

<u>shuffle - once</u>



"same - order"

~~Statistically not~~

use a fixed step size.

Diminishing step size (assuming $\alpha_t \leq 1$)

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

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

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

$$\Rightarrow 0 = -\mu \rho_t + \frac{\alpha_t \sigma^2 L}{B}$$

$$\Rightarrow \alpha_t \simeq \frac{\rho_t \mu B}{\sigma^2 L}$$

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

$$\frac{1}{\rho_{t+1}} \geq \frac{1}{\rho_t - \frac{\mu^2 B}{2\sigma^2 L} \rho_t^2} \quad \left( \frac{1}{\rho_t} \right)$$

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

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

$$\Rightarrow \frac{1}{\rho_T} \geq \frac{1}{\rho_0} + \frac{\mu^2 B T}{2\sigma^2 L} \Rightarrow \rho_T = \mathbb{E}\left[ f(w_T) - f^* \right]$$

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

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

$$\simeq \Theta \left( \frac{\sigma^2 K}{\mu BT} \right)$$

Lecture 11    Week 6.

Momentum.

Rev: $K \gtrsim \frac{L}{\mu}$ ; grad. noise $\sim \left( \frac{\sigma^2}{B} \right)$

Announcement:  $\circ$ Pset 3
                       $\circ$ PA 2.

Rev. (GD).

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

$\downarrow$

gap after $T$ epochs

Q. Simplist model with large $K$?

a.  2D  problem $\rightarrow$ Quadratic

e.g. $f(w) = \frac{1}{2} \left( w_1^2 L + w_2^2 \mu \right)$

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

$\nabla^2 f(w)$

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

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

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

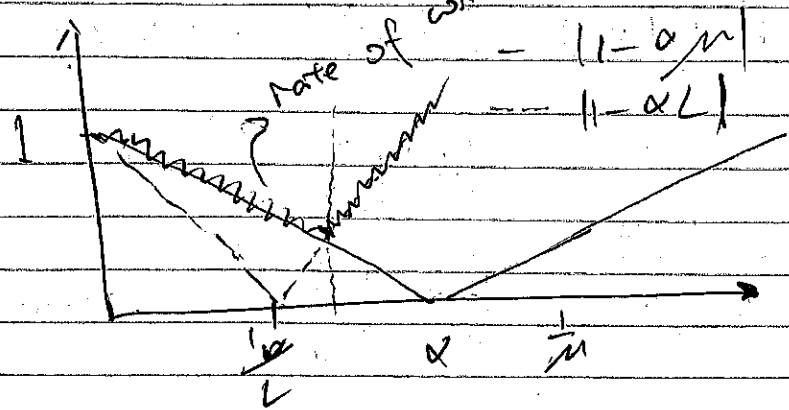$$= w_t - \alpha \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} w_t.$$

$$= \begin{bmatrix} 1-\alpha L & 0 \\ 0 & 1-\alpha\mu \end{bmatrix} w_t \quad \text{, power "T"}$$

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

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

$$f(w_T) = \frac{1}{2}\left( L(1-\alpha L)^{2T}(w_0)_1^2 + \mu(1-\alpha\mu)^{2T}(w_0)_2^2 \right)$$

rate of convergence
$$-|1-\alpha\mu|$$
$$--|1-\alpha L|$$

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



$$1$$

$$\frac{1}{L} \quad \alpha \quad \frac{1}{\mu}$$

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

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

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

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

previous step
(take the more step) step size too large, overshooting the object.

"Self-tune" the step size

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

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

we need
$$0 < \beta < 1$$

particle undergoes

"same" force term

"friction term"
~ or momentum.

How Polyak momentum applies to 2D system?

$$W_{t+1} = W_t - \alpha \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} W_t + \beta(W_t - W_{t-1})$$

$\underbrace{\qquad}_{A}$

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

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

$$\begin{bmatrix} W_{t+1} \\ W_T \end{bmatrix} = \begin{bmatrix} (1+\beta)\mathbb{I} - \alpha A & -\beta\mathbb{I} \\ \mathbb{I} & 0 \end{bmatrix} \begin{bmatrix} W_1 \\ W_0 \end{bmatrix}$$

$\underbrace{\qquad\qquad}$

4×4 block matrix

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

eigenvalue

(2a)

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

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

$\underbrace{\qquad}_{\text{complex vals.}}$

Polyak gives:

$$\alpha = \frac{1+2\beta}{L/\mu}, \qquad \sqrt{\beta} = \frac{\sqrt{K}-1}{\sqrt{K}+1}$$

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

$$T \approx \sqrt{K} \qquad \sqrt{\beta}^T \approx \exp\left(-\frac{LT}{\sqrt{K}+1}\right)$$

$$T \approx \sqrt{K}$$

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

$$|2\sqrt{\beta}| \geq |1+\beta+\alpha\lambda|$$

$$\hookrightarrow 2\sqrt{\beta} \geq 1+\beta+\alpha\lambda \geq -2\sqrt{\beta}$$

$$\mu \le \lambda \le L$$

Conditioning

Rev. — momentum

$$T = \mathcal{O}\left(n\, \underset{\frac{z}{c}}{K}\, \log\left(\tfrac{1}{\varepsilon}\right)\right)$$

condition num.

could be too large

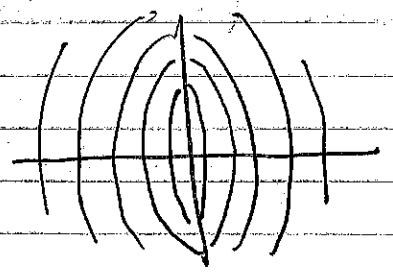too small.

last lec.

$$\hookrightarrow T_{mom} = \mathcal{O}\left(n\sqrt{K}\, \log\left(\tfrac{1}{\varepsilon}\right)\right)$$

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

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

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

$$C = a w_1^2 + b w_2^2$$
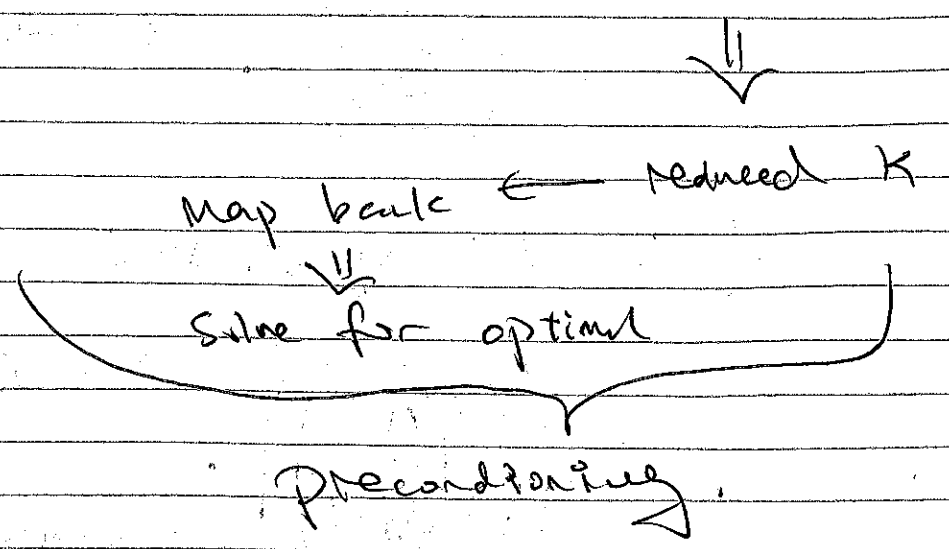
"poorly-conditioned" prob.

$$\overset{\frown}{\omega = Pu}$$

$$\min_{\omega} f(\omega) = \min_{\omega = u} f(P_u)$$

ALA. $|P| \neq 0$

↳ invertible

any matrix, invertible

$$\arg\min_{\omega} f(\omega) = P\left(\arg\min_{u} f(Pu)\right)$$

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

graphically, this is essentially

"stretching" the poorly-conditioned

problem, i.e. apply linear operation

$$\Downarrow$$

map back ⟵ reduced K

$$\Downarrow$$

solve for optimal

preconditioning.

instead of saying minimize $f(\omega)$

over $\omega$ — solve minimize $f(Pu)$

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

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

- how to set $P$? what cond. can

we get? (PD, $A > 0$, all eigenvel.

$A$ are positive)

$$\overset{\frown}{P^T A P = I}^{\text{D. symmetric}}$$

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

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

suppose $P$ symmetric, invertible

- $PAP = I$

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

replace eigenvel.

$$A = P^{-2} \Rightarrow P = A^{-\frac{1}{2}}$$

in transformed space.

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

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

assume $W = Pu$ ; $W_t = Pu_t$

$$Pu_{t+1} = Pu_t - \alpha PP^T \nabla f(Pu_t).$$

$$W_{t+1} = W_t - \alpha \underset{\underset{R}{\downarrow}}{PP^T} \nabla f(W_t) \qquad \text{Preconditioner}$$

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

positive invertible - definite
symmetric.

Preconditioned GD

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

$$x^T R x = x^T P P^T x$$
$$= (P^T x)^T (P^T x)$$
$$= \|P^T x\|^2 > 0$$

chose $R$ to be diagonal preconditioner

$R$ diagonal $\longrightarrow$ stretch along the
coordinates'
axis

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

time multiply $\mathcal{O}(d)$.

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

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

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

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

— run GD in implicit scaled space.

AdaGrad

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

$$r_i \leftarrow r_i + g_i^2$$

$$\left( (r_i)_t = \sum_{k=0}^{t} (g_i^2)_t \right)$$

$$W_i \leftarrow W_i \theta - \frac{\alpha}{\sqrt{r_i}} g_i \quad \left( \alpha_i = \frac{\alpha}{\sqrt{r_i}} \right)$$

– AdaGrad.

$\hookrightarrow$ "Adaptive Optimization" Alg.

---

Lecture 12.     Week 7.

AdaGrad.

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

loop:
get example grad $g \in \mathbb{R}^d$

(depends too much on the history)

$g = \nabla f(w; x)$ minibatch etc

$\uparrow$ large example

(infinite memory.) $\Leftarrow$

update $r \leftarrow r + g^2$

$\quad$ update $w \leftarrow w - \frac{\alpha}{\sqrt{r} + \xi} \cdot g$

– done in element–wise implicitly.

small number.

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

mul. num. of steps.

$\hookrightarrow \approx T(\mathbb{E}(g_i)^2 + Var(g_i))$

could be dominant

$\rightarrow$ fix this : RMSProp.

RMSProp

hyperparameter $\alpha > 0$, $\rho \in (0,1)$,

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

loop:

get gradient sample $g$ at $w$.

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

update $w \leftarrow w_0 - \frac{\alpha}{\sqrt{r} + \varepsilon} g$

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

typical $\rho = 0.999$

Combining momentum + RMSProp +
Correction for $0$-init

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

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

momentum:

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

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

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

$\alpha, \rho_1, \rho_2$ hyperparam.

loop:

get gradient sample $g$ @ $w$.

$$s \leftarrow \rho_1 s + (1-\rho_1) g$$

$$r \leftarrow \rho_2 r + (1-\rho_2) g^2$$

Update:

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

$$\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$$

$$\hat{r} \leftarrow \frac{r}{1 - \rho_2^t} \qquad t = \# \text{ of steps}$$

if $g_t$ denotes the $g$ drawn at

step $t$, $\quad s_t = \sum_{k=0}^{t-1} \rho_1^k (1-\rho_1) g_{t-k}$

$s$ after $t$ steps

$$\sum_{k=0}^{t-1} \rho_1^k (1-\rho_1)$$

$$= (1-\rho_1)\left(\frac{1 - \rho_1^t}{1 - \rho_1}\right)$$

SI: exponential MA:
sequence $x_0, x_1, x_2, \ldots$ } recurrence

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

today — parameter:

subsampling — training number n.

momentum — k.

variance of the gradient

               ↳ avoid having

                    example

## Variance Reduction

→ SVRG

    ↳ strive to behave GD

    — yet SGD each steps.

↳ Large-Scale convex optimization

               → great !!

for DL — Not so good.

  GD        SVRG.

$$\mathcal{O}\left(nk\log\left(\frac{1}{\varepsilon}\right)\right) \quad \mathcal{O}\left((n+k)\log\left(\frac{1}{\varepsilon}\right)\right)$$

Polyak Averaging

— Output an average:

if $W_t$ is parameter vector after

t steps

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

Stochastic Weight averaging

          (SWA).

Rev...

$$\mathcal{O}\left(n|K \cdot \log\left(\frac{1}{\varepsilon}\right)\right) \quad - \quad GD.$$

$$\mathcal{O}\left(K/\varepsilon^2\right) \quad -SGD$$

Momentum: $\mathcal{O}\left(n\sqrt{K} \log\left(\frac{1}{\varepsilon}\right)\right) \quad - MSGD$

$$\mathcal{O}\left((n+K)\log\left(\frac{1}{\varepsilon}\right)\right) \quad SVRC$$

$\downarrow$ Polyak  averaging / Stochastic
                    weight averaging

$\Rightarrow \mathcal{O}\left(nK \cdot d \cdot \log\left(\frac{1}{\varepsilon}\right)\right) \Longleftrightarrow \mathcal{O}\left(K \sigma^2 d/\varepsilon^2\right)$

$$\mathcal{O}\left(n\sqrt{K} d \log\left(\frac{1}{\varepsilon}\right)\right)$$

linear model : $w \in \mathbb{R}^d$.

$$\ell(w) = \frac{1}{n} \sum_{x,y \in \mathbb{D}} \mathcal{L}(w^T x, y).$$

time to compute a example grad

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

---

PCA $\to$ covariance matrix.    $(d \times d)$.

                    computationally expensive

Random Projection. (Johnson – Lindenstrauss

$X \in \mathbb{R}^{D} \Rightarrow Ax \in \mathbb{R}^d$  (Transform)

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

length

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

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

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

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

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

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

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

$$\mathbb{E}\left[\|A(x_i - x_j)\|^2\right] = \mathbb{E}\left[(x_i - x_j)^T A^T A (x_i - x_j)\right]$$

$$= (x_i - x_j)^T \, \mathbb{E}[A^T A] \, (x_i - x_j)$$

$\hookrightarrow$ Gaussian matrix.

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

$$= \begin{cases} 0 & \text{if } i \neq j \\ \sum_{k=1}^{d} \mathbb{E}[A_{ki}^2] & \text{if } i = j. \end{cases}$$

$\sigma^2$

$$\mathbb{E}[A^T A] = \sigma^2 \, d I.$$

$$\hookrightarrow = (x_i - x_j)^T (\sigma^2 d I)(x_i - x_j)$$

$$= \sigma^2 \, d \, \|x_i - x_j\|^2 \overset{\text{if}}{=} \|x_i - x_j\|^2$$

$\hookrightarrow \sigma^2 = \frac{1}{d}.$

__Auto Encoders__



$\downarrow$ downstream.

---

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

$$e.g. \; \frac{1}{\|D\| \times D} \sum_i \|x - \hat{x}\|_2^2$$

$\uparrow$ unsupervised learning.

__Sparsity__

$$\text{density of } x = \frac{\# \text{ of nonzero entries } x_i}{\text{size of } x \, d.}$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 2 \\ 0 \\ 0 \\ 3 \end{bmatrix} = \{(4,2,0),(7,3,0)\} = \begin{matrix} - \text{indices} = [4,7] \\ - \text{values} = [2.0, 3.0] \end{matrix}$$

$16 \quad$ ish

$\swarrow 8$

__Sparse matrix.__

$COO$

$\begin{cases} \text{expensive to comp.} \\ \text{easy to write} \end{cases}$

$$\begin{bmatrix} 0 & 5 & 0 & 0.3 \\ 0 & 10 & 0 & 0.0 \end{bmatrix} \mapsto$$

· row idx. $[1,1,2]$
· col. idx. $[3,6,2]$
· values $[5, 3, 1]$

$CSR \xrightarrow{tr.} CSC$

- row offsets. $[0,2]$ $[3,6]$ $[2] \Rightarrow [3.6.2]$
- col idxs. $[5.0. 3.2]$ $[1.0]$ $[5.0. 3.0 1.2]$
- values.

# Week 8. Lecture 15/16

## Deep learning

convex opt.

## Neural Networks

- MLP.

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

$x \rightarrow \boxed{\text{linear}} \rightarrow O_1 \rightarrow \boxed{\sigma(\cdot)} \rightarrow a_1 \rightarrow O_2$

$\downarrow$
$\sigma$
$\downarrow$
$a_2$
$\downarrow$
...

$O_1 = W_1 x + b_1$

$\begin{cases} W_1 \in \mathbb{R}^{d_1 \times d_0} \\ b \in \mathbb{R}^{d_1} \\ O_1 \in \mathbb{R}^{d_1} \end{cases}$

$a_1 = \sigma(O_1)$

Nonlinear function acts element-wise as the activation function.

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

---

universal: for continuous input → output, there exists a function th approx. with bounded parameters (good accuracy)

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

★ NN tend not to overfit!
$\rightarrow$ larger models generalize better.

( ▷ double-descent. )

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

$|D| = n.$   $n = d.$

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

larger model generalize better

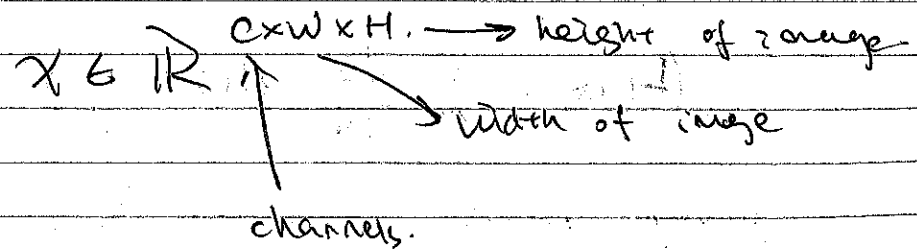Spectrum for SGD, Adam. or sth similar



test "double-descent" phenomenon.

complexity
↳ data-dependent.

$a_\ell, o_\ell \in \mathbb{R}^{d_\ell}$

$a_\ell = W_\ell \, o_{\ell-1} + b_\ell \quad \rightarrow \quad \mathcal{O}(\quad )$

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

$\mathbb{R}^{d_\ell} \quad [\ \cdots\ ] \in \mathbb{R}^{d_\ell}$

$o_\ell = o_\ell(a_\ell) \quad \rightarrow \quad \mathcal{O}(\quad d_1 )$

$W_1 \in \mathbb{R}^{d_1 \times d_0}$

$b_1 \in \mathbb{R}^{d_1}$

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

Image processing $\leftarrow$ CN

$x \in \mathbb{R}^{C \times W \times H} \longrightarrow$ height of image

$\searrow$ width of image

$\uparrow$ channels.

1) convolutional neural network

— restricted subset

convolutional layer in place of the linear layers.

---

Convolution in 1D:



padding

$x_1 w_1 + x_2 w_2 + x_3 w_3$

filter: arbitrarily set of parameters.
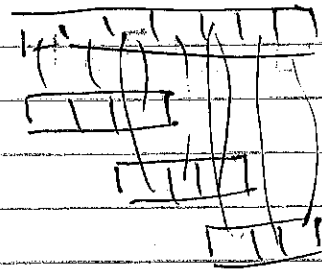
$x_{k+1} w_1 + x_{k+2} w_2 + x_{k+3} w_3$
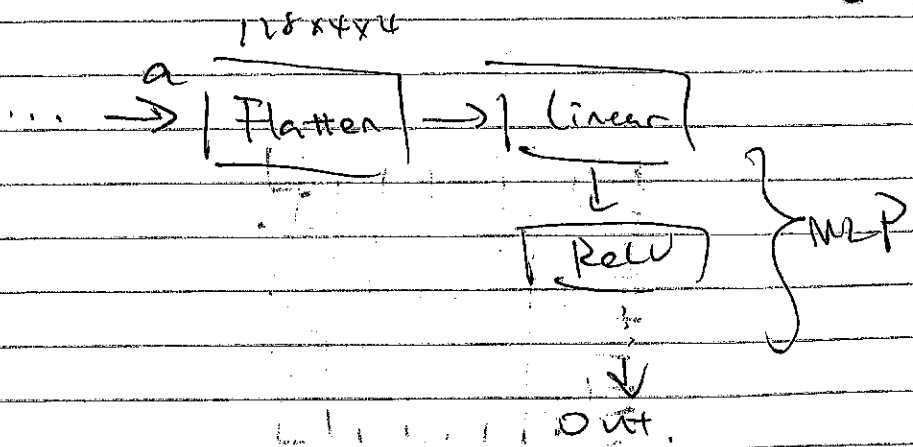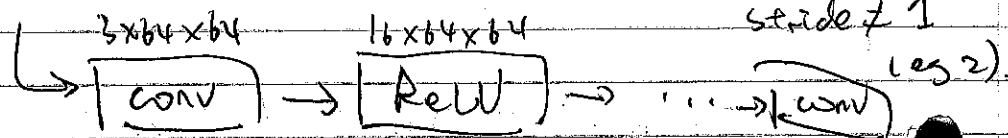
in 2D



Conv. layer              diff. channels

input $x \in \mathbb{R}^{C_0 \times W_0 \times H_0}$   output $a \in \mathbb{R}^{C_1 \times W_1 \times H}$

parameters $W \in \mathbb{R}^{C_1 \times C_0 \times k_w \times k_H}$

$a_i = \sum_{j=1}^{C_0} \text{conv}(W_{ij}, x_j) \qquad W \in \mathbb{R}^{C_1 \times C_0 \times k_w \times k_H}$

$b \in \mathbb{R}^{C_1}$

· stride



outputs 3

input image:                                          32x32x32

$3 \times 64 \times 64$      $16 \times 64 \times 64$      stride ≠ 1

→ [conv] → [ReLU] → ⋯ → [conv]      (eg 2)

$128 \times 4 \times 4$

⋯ → a [Flatten] → [Linear]
                              ↓
                          [ReLU]  } MLP
                              ↓
                          Out.

Week 9.    lecture 17.

Overfitting / Underfitting ← low capacity

↳ Model has high capacity
(overparametrized).

DNNs high ← representational
any weights.

weights coming ← effective.

from the learning algorithm.

Dropout
    ↳ increase the representation capacity

During training → randomly drop
(independently each step) neurons.

    ↓
forces the learn to employ
most neurons.

## Batch Norm.

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

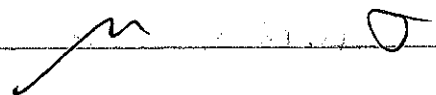$$u \longrightarrow \frac{u - \mu}{\sqrt{\sigma^2}} = \frac{u - \mu}{\sigma}$$

↳ instead. look at $u \in \mathbb{R}^B$

— conv. ⟹ image → per channel!
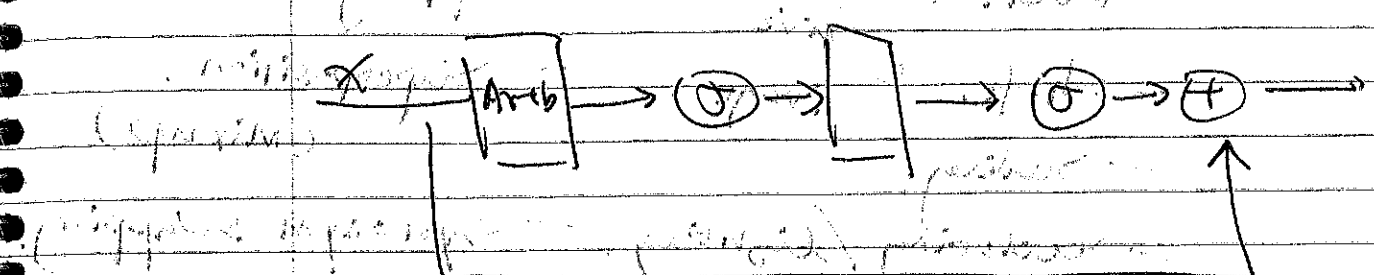
— saves a running [exp.] average, of the mean & variance

$$\mu \qquad \sigma$$

Model. train ( )  ⟹ turning on/off
Model. eval ( )       batch norm.

$$u \longrightarrow \frac{u - \mu}{\sigma} \cdot \gamma + \beta$$

because scalar learned param (SGD)

---

Residual connections / skip connections



Res Net.

Dense Net

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

$$(x_i, y_i) \in D \subseteq X \times Y$$

⟹ sparse data (prob.)

$n$ too small.

⟹ Data Augmentation.

Aug. function $A \cdot x \xrightarrow{[0,1]} x \cdot$ ← random noise

↑ randomness

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

example chosen @ t

IMAGES:

— Translation ⎫          — Noise / Arti fact
— Rotation    ⎬ mix       Addition
— Shift    ⎭ image       — Cropping
— Scaling                 + Superposition
— recoloring/lighting          (MIXUP)
                          — synonym swapping
                          — deletion
                          — backtrans.

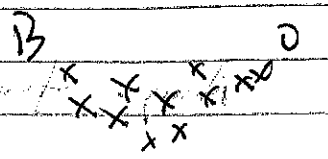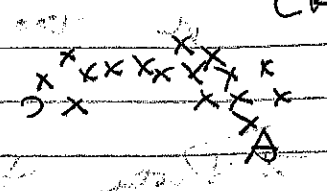Semi-supervised learning

"data - cheap"  "label - expens."

lots of data, few labels
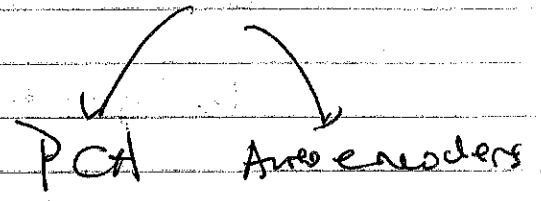
assump.
   Similar inputs have similar outputs.
                    ↳ (kNN)



Clustering → data is clustered
             labels by cluster.

Manifold → data lie on low-dimension
           manifold

           PCA    Autoencoders

Week 10.  lecture 19.

Sequence models.

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

↳ different per example.

discussion:

takes sequence as input.

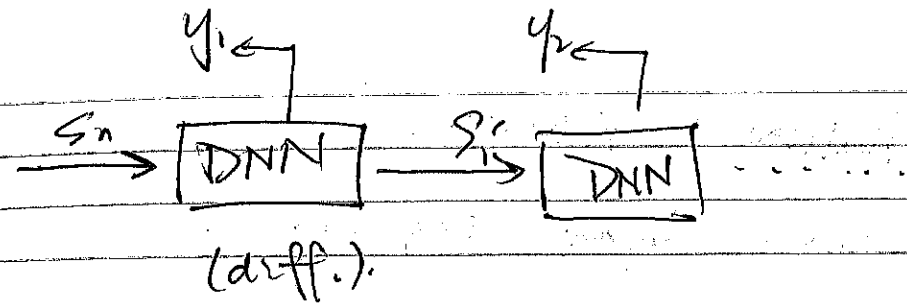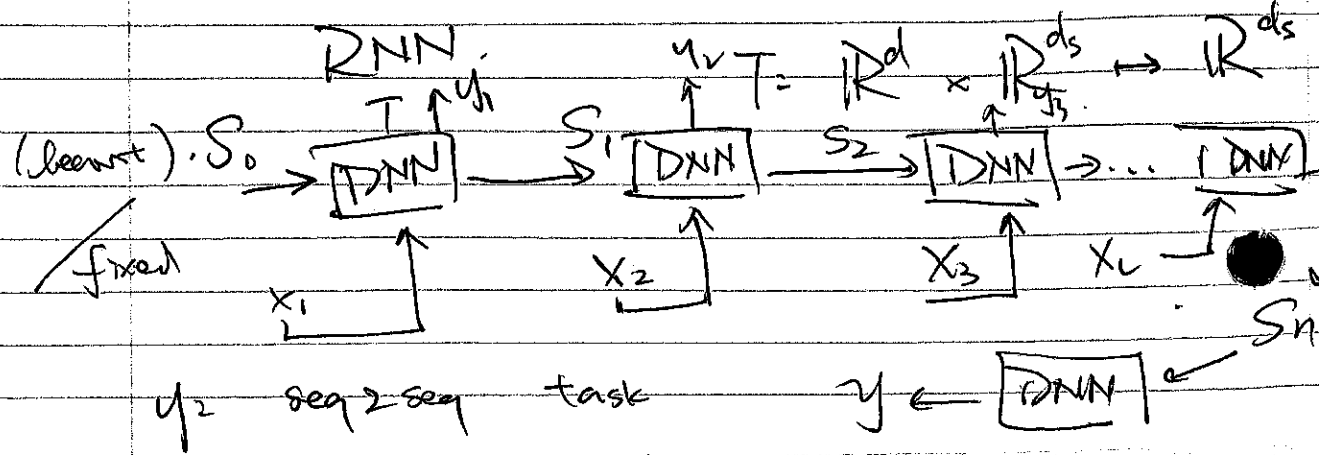↳ generate sequence output.

▽ handle sequence input/output?
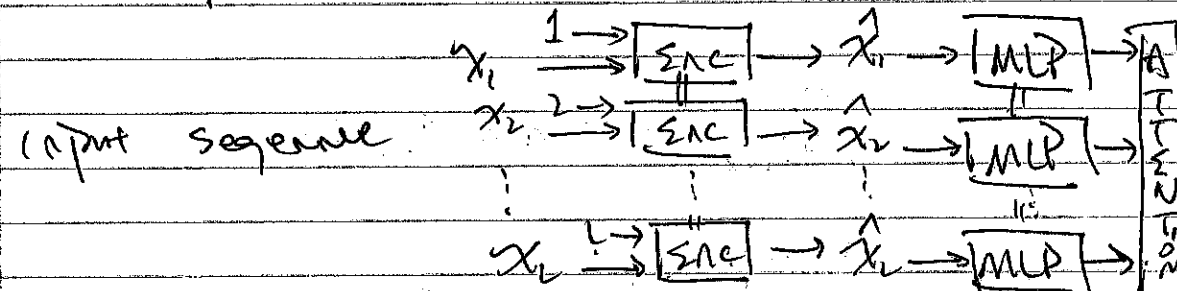
— padding → max length

— counts/sum

— RNN ≈ sequential

— Transformers. ≈ parallel.

DFA:   $S_{next} = T(S_{current}, X_i)$.      $X_i \mathbb{R}^c$

RNN.

$T: \mathbb{R}^d \times \mathbb{R}^{ds} \rightarrow \mathbb{R}^{ds}$

(learnt) $S_0$ →[DNN]→ $S_1$ →[DNN]→ $S_2$ →[DNN]→... [DNN]

fixed        $X_1$         $X_2$         $X_3$   $X_L$
                                                      $S_n$
$y_2$  seq 2 seq  task        $y$ ←[DNN]←

$S_n$ →[DNN]— $S_i$ →[DNN] ······
        $y_1$            $y_k$
(diff.)

• LSTM.

→ Transformers.

input sequence
$X_1 \xrightarrow{1}$ [Enc] → $\hat{X}_1$ →[MLP]→ A
$X_2 \xrightarrow{2}$ [Enc] → $\hat{X}_2$ →[MLP]→ T
                                              T
                                              E
$X_L \xrightarrow{} $ [Enc] → $\hat{X}_L$ →[MLP]→ N

$X_i = Enc(X_i, \cdot)$.

Mathematically attention layer:

$X_1 \xrightarrow{128}$ [NN] ⇒ $k_1$, $Q_1$, $V_1$
$X_2 \xrightarrow{128}$ [NN] ⇒ $k_2$, $Q_2$, $V_2$
$X_3 \xrightarrow{128}$ [NN] ⇒ ...
$X_L \xrightarrow{128}$ [NN] ⇒ $k_L$, $Q_L$, $V_L$

$\left\{ \begin{array}{l} k_i \in \mathbb{R}^{d_k} \\ V_i \in \mathbb{R}^{d_v} \\ Q_i \in \mathbb{R}^{d_k} \end{array} \right.$

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

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

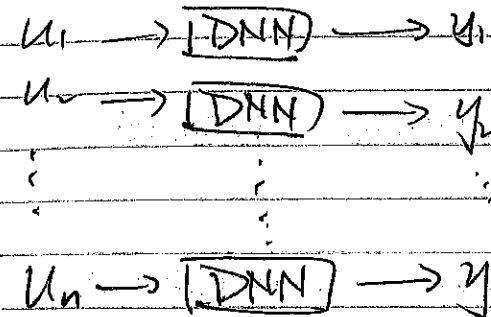$Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V$
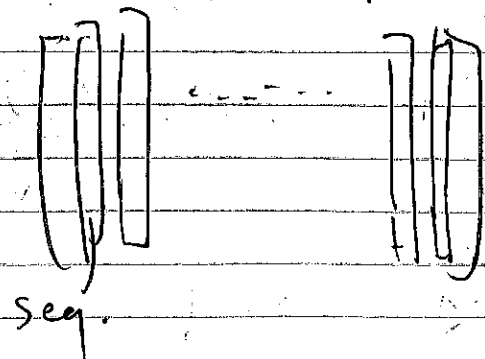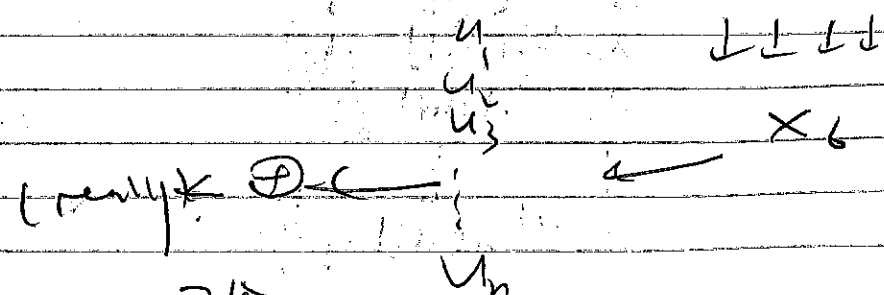
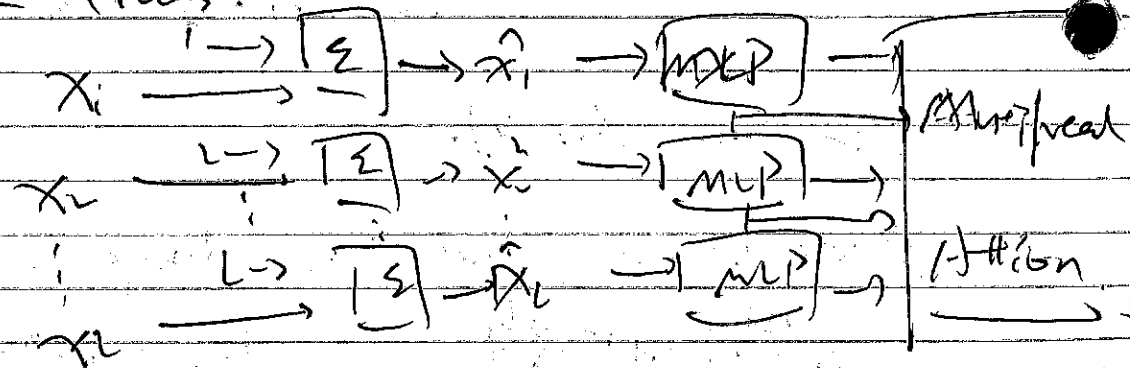approximate table lookup:       Attention block.

Multihead Attention

Same input sequence.

↓

Same no. attention blocks

↓

diff weights. → concatenate output

$$\hat{x}_{ins} = Enc(x_i, v).$$

- Trans.

$$x_1 \longrightarrow \boxed{\Sigma} \longrightarrow \hat{x}_1 \longrightarrow \boxed{MLP}$$
$$x_2 \longrightarrow \boxed{\Sigma} \longrightarrow \hat{x}_2 \longrightarrow \boxed{MLP}$$
$$x_L \longrightarrow \boxed{\Sigma} \longrightarrow \hat{x}_L \longrightarrow \boxed{MLP}$$

Multihead Attention

$u_1$
$u_2$
$u_3$

(really) ⊕

$X_6$

$u_n$

seq.

$u_1 \longrightarrow \boxed{DNN} \longrightarrow y_1$
$u_2 \longrightarrow \boxed{DNN} \longrightarrow y_2$
$\vdots$
$u_n \longrightarrow \boxed{DNN} \longrightarrow y_n$

instead, we do

$u_1$
$u_2 \longrightarrow ⊕ \longrightarrow \boxed{DNN}$  $y_1$
$\vdots$
$u_n$       1

# Lecture 20.

## Kernels.

opposite of dimension reduction.

map into higher dimensional space.

$$\phi(x) = \mathbb{R}^d \mapsto \mathbb{V}$$

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

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

## 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

$$k(x_1, x_2) = (1 + x_1^2 x_2)^P.$$

## Kernel definition.

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

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

$\rightarrow$ 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 semidefinite

$k$ 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 maps $\phi_1$ & $\phi_2$, then $k = k_1 + k_2$

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

$$k(u,v) = k_1(u,v) + k_2(v)$$
$\leftarrow$ is a kernel.

$$k(u,v) = c\, k_1(u,v) \quad \text{if } c > 0.$$
$$\phi(x) = \sqrt{c}\ \phi_1(x).$$

$$k(u,v) = k_1(u,v)\, k_2(u,v).$$
$$\phi(x) = \phi_1(x) \otimes \phi_2(x).$$
$$k(u,v) = f(u)\, k_1(u,v)\, f(v).$$
$$\phi(x) = f(x)\, \phi(x).$$

$$f: \mathcal{X} \to \mathbb{R}.$$

component loss $\in$
$$f(w, x, y) = \ell(w^T \phi(x), y)$$
$$= \ell(\langle w, \phi(x) \rangle, y)$$

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

$$\nabla_w f(w, x, y) = \ell'(w^T \phi(x), y)\, \phi(x).$$
$$\Rightarrow \text{for GD / SGD.}$$
$$w_k \in \text{span}\{\phi(x_1), \phi(x_2), \dots \phi(x_n)\}$$
$$w_k = \sum_{i=1}^{n} a_i \phi(x_i)$$

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

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

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

$$\left(\text{let } K_{ij} = k(x_i, x_j)\right)$$
$$= \frac{1}{n} \sum_{j=1}^{n} \ell((Ka)_j, y).$$

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

② precompute & store $\phi(x_i)$ before training.

— train in the transformed space

③ compute $k$ on-the-fly as needed

④ precompute $K$ & store.

Lecture 21.

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

ways to do learning:

① Learn in the transformed
space, compute $\phi$ on-the-fly.

② learn " " , precompute $\phi$.

③ learn in the "kernel" space
while computing $k$ on-the-fly.

④ learn in "kernel trick" space,
precompute $K_{ij} = k(x_i, x_j)$

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

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

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 | $\mathcal{O}(Td)$ | $\mathcal{O}(d)$ (not inc. train.s) |
| ② | | |
| ③ | | |
| ④ | | |

— learn in the transformed space

| | | | |
|---|---|---|---|
| ① | $\mathcal{O}(Tmd)$. precompute | $\mathcal{O}(m)$. | "tmd" |
| ② | $\mathcal{O}(Tm + nmd)$ | $\mathcal{O}(nm)$. | (+training d) |
| ③ | $\mathcal{O}(Tnd)$. | $\mathcal{O}(n)$. | |
| ④ | $\mathcal{O}(Tn + n^2d)$ | $\mathcal{O}(n^2)$. | |

$\left( \, \ell'\left(\sum_{j=1}^{n} u_j \, k(x_i, x_j), y_i\right) \phi(x_i) \right.$

$\nabla f_i(w) = \ell'\left(\phi(x_i)^T \sum_{j=1}^{n} u_j \phi(x_j), y_i\right) \phi(x_i)$

an SGD step on compute $i$

$$u_i \leftarrow u_i - \alpha \, l' \left( \sum_{j=1}^{n} u_j k(x_i, x_j) ; y \right)$$

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.

$\vdots$ Random fourier features.

if kernel : $k(x_i, x_j) = k(x_i - x_j)$.

then $k(x_i, x_j) = \mathbb{E} \left[ 2 \cos(w^T x_i + b), \right.$

$$\left. \cos(w^T x_j + b) \right].$$

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

$w \sim \mathcal{F}(k)$ fourier transform of $k$.

"Reg" $k$

RBF kernel:

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

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

$$\mathbb{E} \left[ 2 \cos(w^T x_i + b) \cdot \cos(w^T x_j + b) \right]$$

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

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

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

Set $\psi_i(x) = \frac{1}{\sqrt{D}} \sqrt{\frac{2}{D}} \cos(w_i^T x + b_i)$

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

$\vdots$ (TBC)

⑤ learn w/ approx. feature map comp.
$\psi$ on-the-fly.

⑥ learn w/ approx. f.m. compute & cache

Vanilla linear model.

compute.            Memory.

Ⓐ  Ⓗ $(TD\alpha)$.              Ⓗ $(D)$

Ⓑ  Ⓗ $(TD + nD\alpha)$   Ⓗ $(nD)$.
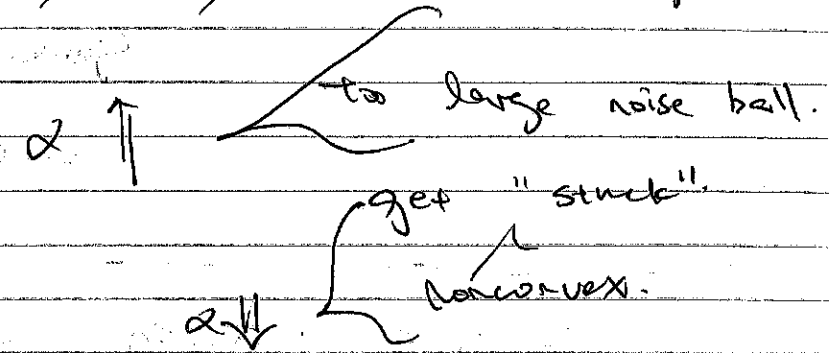
large $D$ - high "accuracy" in repr. $k$.

— high compute & memory cost.

---

Lecture 22.

▷ Hyperparameter optimization

— step size / lr / $\alpha$  → diverge.

$\alpha \uparrow$ {  too large noise ball.

  get "stuck"
$\alpha \downarrow$ {  nonconvex.

— $\lambda$ for regularization
        ↓
    over/under regularize

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

— $\beta$ momentum.  → if $\beta \notin (0, 1)$
                            ‼
                            bad.

— $B$ batch size. → too small can
                affect noise ball.
          too large → much like GA.
performance  matching parallel capabilities
                        of HW.

— choice of architecture: DNN

$\downarrow$

Generalize poorly.

systems implications.

— — — — — — — — Non-DNN — —

— tree depth $\rightarrow$ weak learner

overfit, & lots of mem.

— feature vector length for

Random features.

$\rightarrow$ poor fidelity,

high compute/memory

cost.

— # of epochs $\rightarrow$ high loss/error

compute cost.

— dimension $\rightarrow$ poor accuracy.

high compute cost.

— k in KNN & k-means $\rightarrow$ under/overfit

More compute for larger k.

the task for assigning hp:

$\Rightarrow$ HPO — Any process assign hps

Standard: $\beta$  0.9  0.99

$(\beta_1, \beta_2)$  $(0.9, 0.999)$

B   power of 2, e.g. 256.

— Grid Search

$\hookrightarrow$ CoD

— Random Search

less reproducible than GS.

Parallelism.

Early stopping

— dropping bad hyperparameters

that perform poor each epoch

goal: minimize $F(\alpha, \beta, B, k, \alpha)$.

"derivative - free" optimization. (DFO)

Encode our knowledge about the space

& the problem $\implies$ HPO

Mean / Median  $f_0$ $f_1$ $f_3$ $f_4$ ...

HPO

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

look at: $P(f(x_*) \mid f(x_1) = y_1, f(x_2) = y_2, \cdots)$

$f \in \mathbb{R}^d$, $f \sim \mathcal{N}(0, \Sigma_{\downarrow})$

$\underset{\text{mean}}{\downarrow}$ $\underset{\text{covariance}}{\downarrow}$

☆ Fact: if $(x, y)$ are

jointly Gaussian.
$\wedge$
multivariate

$\{$ Set by also
$\{$ according to
$\{$ intuition.

the $x \mid y = y$ is also Gaussian.

☆ Fact: " $\sim$ ", then $X$ is Gaussian

$f(x_*) \mid f(x_1) = y_1, f(x_2) = y_2, \cdots$

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

How to pick the next point?

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

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

(lower confidence bound)

$\longrightarrow a$ - acquisition function

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

(Probability of Improvement)

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

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

accumulated distribution function

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

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

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

$$= P(\sigma u + \mu \leq f_{best})$$

for $U \sim \mathcal{N}(0, 1)$

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

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

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

Expected improvement

by convention

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