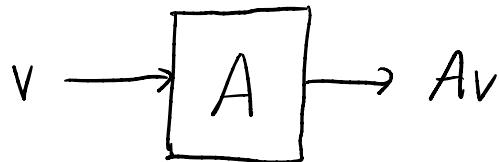
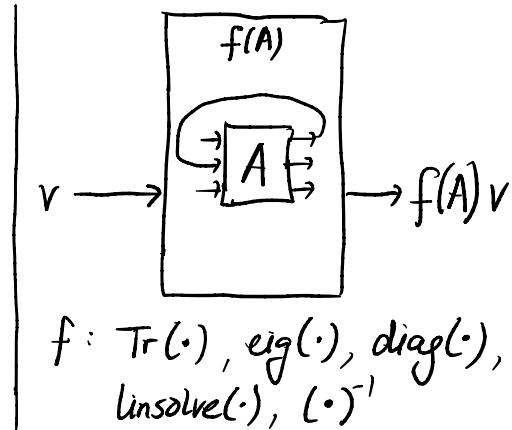


# Large-scale linear algebra with curvature matrices

## Overview: linear operators

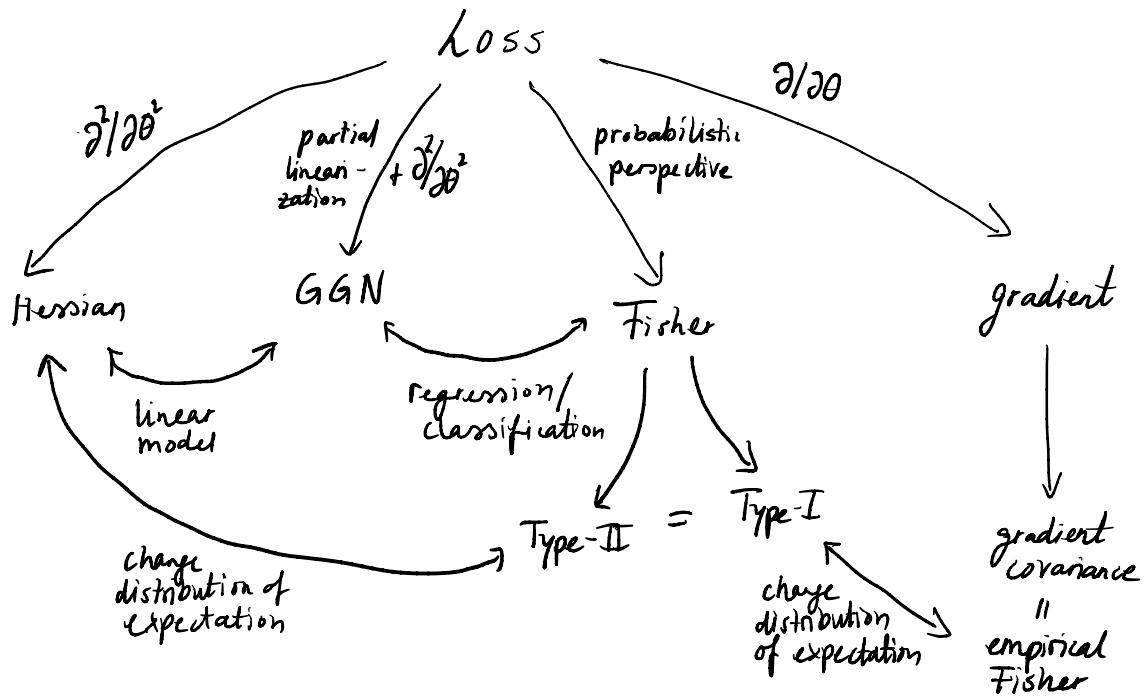


$$\begin{aligned} A(v_1 + v_2) &= Av_1 + Av_2 \\ A(cv) &= cAv \end{aligned}$$



$f: \text{Tr}(\cdot), \text{eig}(\cdot), \text{diag}(\cdot), \text{linsolve}(\cdot), (\cdot)^{-1}$

## Overview: DL curvature matrices and their relation





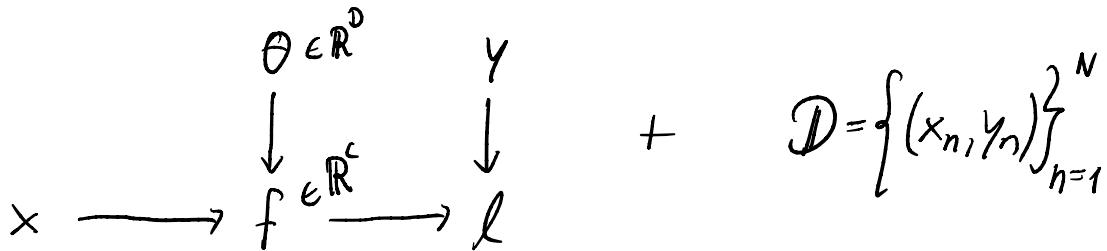
DL

Curvature

matrices



## Loss, gradient, Hessian



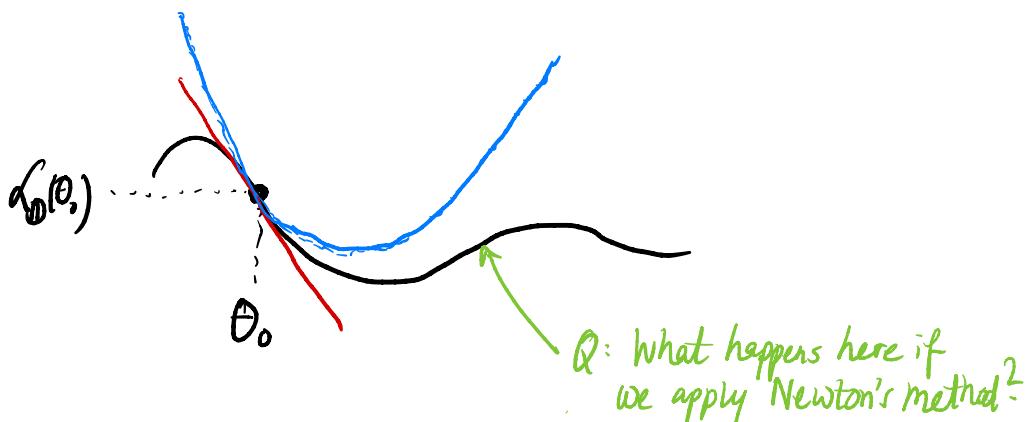
$$\mathcal{L}_{\mathcal{D}}(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(f(x_n, \theta), y_n) \quad \in \mathbb{R}$$

$$g_{\mathcal{D}}(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta} \ell(f(x_n, \theta), y_n) \quad \in \mathbb{R}^D$$

$$H_{\mathcal{D}}(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta}^2 \ell(f(x_n, \theta), y_n) \quad \in \mathbb{R}^{D \times D}$$

Motivation: Taylor expansion

$$\begin{aligned}
 \mathcal{L}_{\mathcal{D}}(\theta - \theta_0) &= \mathcal{L}_{\mathcal{D}}(\theta_0) + g_{\mathcal{D}}(\theta_0)^T (\theta - \theta_0) \\
 &\quad + \frac{1}{2} (\theta - \theta_0)^T H_{\mathcal{D}}(\theta_0) (\theta - \theta_0) + O((\theta - \theta_0)^3)
 \end{aligned}$$



Use cases:

- Optimization (fit local quadratic, jump to minimum)
- Pruning (find weight that least affects loss when removed)
- Sharpness & influence functions (sensitivity of loss w.r.t. perturbations in the data/parameters)
- Laplace approximations (fit a Gaussian to a more complicated distribution)

Pros and cons of the Hessian:

- Cannot touch the full matrix explicitly ( $D \times D$ )
- Can multiply with the Hessian efficiently using first-order autodiff! (double-backward)

$$\left( \frac{\partial^2 L}{\partial \theta^2} \right) v = \frac{\partial}{\partial \theta} \left[ \left( \frac{\partial L}{\partial \theta} \right) v \right] = \frac{\partial}{\partial \theta} [g \cdot v]$$

- Hessian can be indefinite because  $f \circ f$   
is usually non-convex
 

convex

non-convex

## Generalized Gauss-Newton (GGN) matrix

- Compositions of convex functions are convex  
 → can we convexify  $f^2$ ? Yes, through linearization around some anchor  $\theta_0$ :

$$f(\theta) \rightarrow f_{\theta_0}^{\text{lin}}(\theta) = f(\theta_0) + \underbrace{\nabla_{\theta} f(\theta_0)}_{C \times D} (\theta - \theta_0)$$

$$[ \dots ]_{i,j} = \frac{\partial f_i}{\partial \theta_j}$$

$$H_{\text{pl}}(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta}^2 (l_n \circ f_n)(\theta)$$

↓ Hessian of the partially linearized loss

$$G_D(\theta) = \frac{1}{N} \sum_{n=1}^N \left[ \nabla_{\theta}^2 (l_n \circ f_{\theta_0, n}^{\text{lin}})(\theta) \right] \Big|_{\theta_0=\theta}$$

"GGN"

$$= \frac{1}{N} \sum_{n=1}^N \underbrace{\left( \nabla_{\theta} f_n \right)^T}_{D \times C} \underbrace{\left( \nabla_{f_n}^2 l_n \right)}_{C \times C} \underbrace{\left( \nabla_{\theta} f_n \right)}_{C \times D}$$

- Positive semi-definite
- GGN-vector products through  $\nabla_{\theta} f_n$ ,  $\nabla_{f_n}^2 l_n$
- = Hessian if  $f$  is linear in  $\theta$  (e.g. linear regression)

// Prelude to Fisher: Risk minimization

$$X(\theta) = \frac{1}{N} \sum_{n=1}^N l(f(x_n, \theta), y_n) \quad \begin{matrix} \text{(massage into} \\ \text{an expectation)} \end{matrix}$$

Assume  $\exists$  data-generating process  $P_{\text{data}}(x, y)$   
and we want to minimize the expected risk

$$L_{P_{\text{data}}}(\theta) = \mathbb{E}_{P_{\text{data}}(x, y)} [l(f(x, \theta), y)] \quad \text{w.r.t. } \theta.$$

↳  $P_{\text{data}}$  is intractable though. That's why we collected a data set, which approximates

$$P_{\text{data}}(x, y) \approx P_D(x, y) := \frac{1}{N} \sum_{n=1}^N \delta(x - x_n) \delta(y - y_n)$$

↑  
"empirical data distribution", tractable

We then have

$$\mathcal{L}_{\mathbb{D}}(\theta) = \mathcal{L}_{p_{\theta}}(\theta) = \mathbb{E}_{p_{\theta}(x,y)} \left[ l(f(x,\theta), y) \right]$$

for the empirical risk.

// Connection to maximum likelihood

Want to learn a parametrized density  $p_{\theta}(x,y)$  that approximates  $p_{\text{data}}(x,y)$ :

$$\underset{\theta}{\text{minimize}} \quad KL \left( p_{\text{data}}(x,y) \parallel p_{\theta}(x,y) \right)$$

$$\Leftrightarrow \underset{\theta}{\text{minimize}} \quad \mathbb{E}_{p_{\text{data}}(x,y)} \left[ -\log p_{\theta}(x,y) \right]$$

(we will only model  $y|x$  with parameters, so  $p_{\theta}(x,y) = p_{\theta}(y|x) p_{\text{data}}(x)$ , we don't have access to  $p_{\text{data}}$ , so  $p_{\text{data}} \hookrightarrow p_{\mathbb{D}}$ )

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \underbrace{-\log p_{\theta}(y_n | x_n)}_{\text{This looks like } l(f(x_n, \theta), y_n)}$$

This looks like  $l(f(x_n, \theta), y_n)$  from above, where  $l$  is a negative log-likelihood

Rewrite  $p_{\theta}(y|x) = r(y|f(x, \theta))$ , then we can show

$$1) \quad r(y|f(x, \theta)) = N(y | \mu=f(x, \theta), \Sigma=I)$$

$$\Leftrightarrow -\log r = l(\text{MSE Loss})$$

$$2) \quad r(y|f(x, \theta)) = \text{Cat}\left(y \mid \pi = \text{softmax}(f(x, \theta))\right)$$

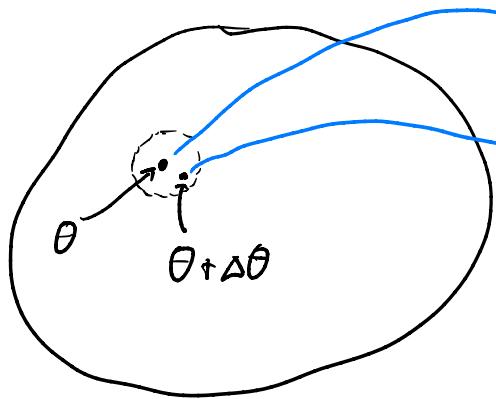
$\in \{1, \dots, C\}$

$$\Leftrightarrow -\log r = l(\text{CrossEntropy Loss})$$

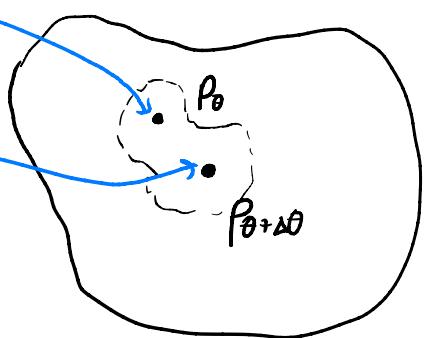
$\Rightarrow$  The NN models a Gaussian/categorical distribution

Probabilistic perspective

parameter space



statistical manifold



$$d(\theta, \theta + \Delta\theta) = \|\Delta\theta\|_2^2$$

$$d(p_{\theta + \Delta\theta}, p_\theta)$$

= (second-order approximation of)  
 $KL(p_{\theta + \Delta\theta} \| p_\theta)$

$$= -\frac{1}{2} E_{p_\theta(x,y)} \left[ \nabla_\theta^2 \log p_\theta(x,y) \right]$$

"Fisher information matrix"

$$(use p_\theta(x,y) = p_\theta(y|x) p_\theta(x))$$

$$F(\theta) = \frac{1}{N} \sum_{n=1}^N E_{p_\theta(y|x_n)} \left[ -\nabla_\theta^2 \log p_\theta(y|x_n) \right]$$

$$= \frac{1}{N} \sum_{n=1}^N E_{p_\theta(y|x_n)} \left[ -\nabla_\theta^2 l(f(x_n, \theta), y) \right]$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \nabla_\theta f_n \right)^T E_{p_\theta(y|x_n)} \left[ \nabla_f^2 l(f_n, y) \right] \nabla_\theta f_n$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \nabla_\theta f_n \right)^T E_{p_\theta(y|x_n)} \left[ D_f l(f_n, y) (\nabla_f l(f_n, y))^T \right] \nabla_\theta f_n$$

looks similar to Hessian

Type-II  
looks like GAN

Type-I  
motivation for ECF

# Similarities and differences ( $\ell = -\log p$ )

$$H(\theta) = \mathbb{E}_{P_{\theta}(x,y)} \left[ \nabla_{\theta}^2 \ell \right] = \mathbb{E}_{P_{\theta}(x)} \mathbb{E}_{P_{\theta}(y|x)} \left[ \nabla_{\theta}^2 \ell \right] \quad \text{Hessian}$$

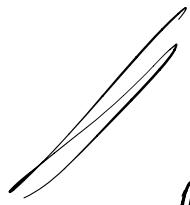
$$\bar{F}(\theta) = \mathbb{E}_{P_{\theta}(x,y)} \left[ \nabla_{\theta}^2 \ell \right] \quad \text{Fisher}$$

$$F_{\text{II}}(\theta) = \mathbb{E}_{P_{\theta}(x)} \left[ \left( \nabla_{\theta} f \right)^T \mathbb{E}_{P_{\theta}(y|x)} \left[ \nabla_f^2 \ell \right] \nabla_{\theta} f \right] \quad \text{Fisher type-II}$$

$$G(\theta) = \mathbb{E}_{P_{\theta}(x)} \left[ \left( \nabla_{\theta} f \right)^T \mathbb{E}_{P_{\theta}(y|x)} \left[ \nabla_f^2 \ell \right] \nabla_{\theta} f \right] \quad \text{GGN}$$

$$\bar{F}_{\text{I}}(\theta) = \mathbb{E}_{P_{\theta}(x)} \left[ \left( \nabla_{\theta} f \right)^T \mathbb{E}_{P_{\theta}(y|x)} \left[ \nabla_f \ell (\nabla_f \ell)^T \right] \nabla_{\theta} f \right] \quad \text{Fisher type-I}$$

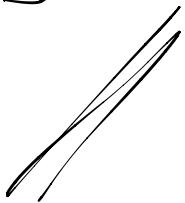
$$EF(\theta) = \mathbb{E}_{P_{\theta}(x)} \left[ \left( \nabla_{\theta} f \right)^T \mathbb{E}_{P_{\theta}(y|x)} \left[ \nabla_f \ell (\nabla_f \ell)^T \right] \nabla_{\theta} f \right] \quad \begin{matrix} \text{empirical} \\ \text{Fisher} \end{matrix}$$



Matrix-free

linear algebra

with linear operators



# Hutchinson trace / diagonal estimation

$$\text{Tr}(A) = \text{Tr}(A \mathbb{I})$$

Trace

$$= \text{Tr}(A \mathbb{E}_{p(v)}[vv^T])$$

$$= \mathbb{E}_{p(v)}[\text{Tr}(Avv^T)]$$

$$= \mathbb{E}_{p(v)}[\text{Tr}(v^T A v)]$$

$$= \mathbb{E}_{p(v)}[v^T A v]$$

$$\approx \frac{1}{S} \sum_{s=1}^S v_s^T A v_s \quad \text{where } v_s \stackrel{i.i.d.}{\sim} p(v)$$

(can also use to estimate  $\|A\|_2^2 = \text{Tr}(A^T A)$ )

$$\text{diag}(A) = \text{diag}(A \mathbb{I})$$

Diagonal

$$= \text{diag}(A \mathbb{E}_{p(v)}[vv^T])$$

$$= \mathbb{E}_{p(v)}[\text{diag}(Avv^T)]$$

$$= \mathbb{E}_{p(v)}[(Av) \odot v]$$

$$\approx \frac{1}{S} \sum_{s=1}^S Av_s \odot v_s$$

where  $v_s \stackrel{i.i.d.}{\sim} p(v)$

## Power iteration

eig

$\lambda_1(A), e_1(A)$

|  $\lambda_1 \leftarrow \infty$   
|  $e_1 \leftarrow v$        $v$  random,  $\|v\|_2 = 1$   
| until converged  
|  $e_1 \leftarrow Ae_1$       matvec  
|  $\lambda_1 \leftarrow \|e_1\|_2$       update eigenval  
|  $e_1 \leftarrow e_1 / \|e_1\|_2$

- (more sophisticated methods, e.g. implicitly restarted Arnoldi iterations (ARPACK), `scipy.sparse.linalg.eigs`)
- Also: Spectral density estimation

$$\rho(\lambda) = \frac{1}{D} \sum_{d=1}^D S(\lambda - \lambda_d)$$

//

## linear systems and inversion

$\text{linsolve}(A, b) = x$  such that  $Ax = b$ .

↪ Iterative solvers like conjugate gradients (CG)  
(details not important, only uses  $[v \mapsto Av]$ ).

// CG inversion: Compute  $A^{-1}b = x$   
 $\Leftrightarrow$  Find  $x : b = Ax$   
 $\Rightarrow x = \text{linsolve}(A, b)$

// Neumann series:

$$A^{-1} = \sum_{k=0}^{\infty} (I - A)^k \quad (\text{if } 0 < \lambda(A) < 2)$$

$$A^{-1}v = v + (I - A)v + (I - A)^2v + \dots \quad (\text{truncate})$$

curvlinops library: offers linear operators for  $H$ ,  
 $\tilde{F}_{\mathbb{II}}, \tilde{F}_{\mathbb{I}}, G, EF$ , and more + functionality to estimate their properties