

(2) Differential
Derivative

$$\begin{cases} f: U \subset \mathbb{R} \rightarrow \mathbb{R} \\ U \text{ open set} \end{cases}$$

f admits a derivative at a point x of U

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = a$$

where a is finite

$a = f'(x)$ derivative of f at x

Equivalently,

$$f(x+h) = f(x) + f'(x)h + \epsilon(h)$$

$$\text{where } \epsilon \lim_{h \rightarrow 0} \epsilon(h) = 0$$

this means that the function f is approximated by the affine function $L(h) = f'(x)h$ in the neighbourhood of x

$f(x+h) \approx f(x) + L(h)$ where h is small geometrically

Let E, F be two vectorial spaces and

$$f: E \rightarrow F$$

Goal: extend the notion of derivative for f

Def: Let U be an open set of E and $x \in U$.

We say that f is differentiable at x if there exists a continuous and linear mapping L from E into F

such that

$$f(x+h) = f(x) + L(h) + \|h\|_E \epsilon(h)$$

$$\forall h \in E \text{ where } \lim_{h \rightarrow 0} \|\epsilon(h)\|_F = 0$$

When it is the case, L is called the differential of f at x
 we will denote L by df_x

$$\begin{cases} df_x : E \rightarrow F \\ \text{linear, continuous} \end{cases}$$

$df_x \in \mathcal{L}(E, F)$

space of linear and continuous
functions from E to F .

if h is small, then

$$f(x+h) \underset{\substack{\text{affine approximation} \\ \text{of } f \text{ on the neighborhood of } x.}}{\approx} f(x) + df_x(h)$$

Link with the derivative

$$E = F = \mathbb{R}$$

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

If f has a derivative $f'(n)$ at n then

$$f(n+h) = f(n) + f'(n)h + h\epsilon(h)$$

$$\text{let } L(h) = f'(n)h$$

L is linear and continuous. From $\mathbb{R} \rightarrow \mathbb{R}$

$$\text{Moreover } \lim_{h \rightarrow 0} \epsilon(h) = 0$$

$$\text{so } L = df_n$$

$$df_x: \mathbb{R} \rightarrow \mathbb{R}$$

$$h \mapsto df_x(h) = f'(n)h$$

Examples:

① $f: E \rightarrow F$

Let $f(x) = c$

where $c \in F$

$$f(x+h) - f(x) = 0$$

Let $L(h) = 0$ if h

$$\epsilon = 0$$

$$f(x+h) = f(x) + L(h) + (\epsilon) \epsilon(h)$$

$df_x : E \rightarrow F$

$$h \mapsto 0$$

$$df_x = 0$$

$$df_x(h) = 0 \quad \forall h$$

Q) Let f be continuous and linear from E to F .

$$f(x+h) - f(x) = f(x) + f(h) - f(x) = f(h)$$

f is linear and continuous

$$\text{so } df_x = f$$

② Let f be a bilinear continuous mapping from $E \times F$ into

II Link with partial derivatives.

Let E be a finite nonempty vectorial space equipped with a basis (e_1, \dots, e_n)

Let $f : E \rightarrow \mathbb{R}$

Def! (partial derivative)

Let U be an open set of E and $x \in U$

For $i = 1, \dots, n$ let

$$g_i : \mathbb{R} \xrightarrow{\sim} \mathbb{R} \quad g_i(x+e_i)$$

If g_i admits a derivative at 0 then we say that f has a partial derivative at x in the direction e_i . We denote $\partial_i f(x)$ this derivative.

$$\begin{aligned}\partial_i f(x) &= \lim_{h \rightarrow 0} \frac{g_i(h) - g_i(0)}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(x+he_i) - f(x)}{h}\end{aligned}$$

Assume that f is differentiable at a point x of U

then, for any h

$$f(x+h) = f(x) + df_x(h) + \varepsilon(h)$$

$$\text{where } \frac{\varepsilon(h)}{|h|} \rightarrow 0, |h| \rightarrow 0$$

Let $h = re_i$ with $r \in \mathbb{R}$

$$f(x+re_i) = f(x) + df_x(re_i) + \varepsilon(re_i)$$

but by linearity

$$df_x(re_i) = r df_x(e_i)$$

$$\text{so that } \frac{f(x+re_i) - f(x)}{r} = df_x(e_i) + \frac{\varepsilon(re_i)}{r}$$

Moreover

$$\lim_{r \rightarrow 0} \frac{\varepsilon(re_i)}{r} = \lim_{r \rightarrow 0} \frac{\varepsilon(re_i)}{|re_i|} = 0$$

Then

$$\lim_{r \rightarrow 0} \frac{f(x+re_i) - f(x)}{r} = df_x(e_i)$$

Prop: if f is differentiable at x then it admits partial derivative in any direction e_i and

$$\partial_i f(x) = df_x(e_i)$$

Moreover

$$df_x(h) = \sum_{i=1}^n \partial_i f(x) h_i$$

Where h_i are the coordinates of h in basis (e_1, \dots, e_n)

Indeed

$$\begin{aligned} df_x(h) &= df_x\left(\sum_{i=1}^n h_i e_i\right) = \sum_{i=1}^n h_i df_x(e_i) \\ &= \sum_{i=1}^n h_i \partial_i f(x) \end{aligned}$$

Now assume that

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

where E is finite dimensional with basis (e_1, \dots, e_n)

We can describe f through its components :

$$j = 1, \dots, m \quad f_j: E \rightarrow \mathbb{R}$$

$\xrightarrow{\text{a b}} f_j(x)$

We can compute partial direction of f_j if they exist

Dif. let $x \in U$ open set of E and $f_j: E \rightarrow \mathbb{R}^m$.

if the components of f admit partial derivative in all direction

then we can define the Jacobian matrix, denoted $Df(x)$

in the special case where $m = n$
 we call jacobian the determinant of this matrix
 let us assume that f is differentiable at x

$$df_x \in \mathcal{L}(E, \mathbb{R}^m)$$

$$\begin{pmatrix} df_x(h) \\ df_x^{(1)}(h) \\ \vdots \\ df_x^{(n)}(h) \end{pmatrix} = \begin{pmatrix} \sum_i \partial_i f(x) h_i \\ \vdots \\ \sum_i \partial_i f(x) h_i \end{pmatrix}$$

Prop:

$$df_x(h) = Df(x) h$$

III - 6 radians

let E be an Hilbert space equipped with an inner product $\langle \cdot, \cdot \rangle$
 and $f: E \rightarrow \mathbb{R}$

if f is differentiable at a point x of an open set U of

E then.

$$\{df_x \in \mathcal{L}(E, \mathbb{R}) : = \mathcal{L}(E)\}$$

(df_x is a linear continuous form on E).

In an Hilbert space a linear form l on E is uniquely determined an element of E' :

$$l(h) = \langle a, h \rangle \quad (\text{Riesz representation theorem})$$

Def: we call gradient of f at x the element of E' , denoted $Df(x)$, such that

$$\langle Df_n(h), h \rangle$$

Let E be finite dimensional with a basis (e_1, \dots, e_n)

Assume that $f: E \rightarrow \mathbb{R}$ is differentiable at x

then $Df_x(h) = Df(x)h$

$Df(x)$ is a matrix of size $(1, n)$

$$Df_x(h) = \langle Df(x)^T, h \rangle$$

By identification

$$Df(x) = Df(x)^T$$

In other words

$$(Df(x))_i = \partial_i f(x).$$

Ex-G-1

Exercise 6.1: First and Second Differentials of a Bilinear Map

Let f be a continuous bilinear map from E^2 to F . We equip the product space E^2 with the norm

$$\|(h, k)\|_{E^2} = \sqrt{|h|^2_E + |k|^2_E}, \quad \forall (h, k) \in E^2.$$

1. Show that f is differentiable on E^2 and compute its differential.
2. Show that f is twice differentiable and compute its second differential.
3. Suppose E is finite-dimensional and $F = \mathbb{R}$. Let (e_1, \dots, e_p) be an orthonormal basis for E . We equip E with the inner product

$$\langle x, y \rangle = \sum_{i=1}^p x_i y_i,$$

where $(x_i)_{i=1}^p$ and $(y_i)_{i=1}^p$ are the coordinates of x and y in this basis, respectively.

- (a) Specify the gradient and Hessian of f in this basis of E .
- (b) Write the first and second differentials in terms of the gradient and Hessian, respectively.

$$1) \quad f(x+h) = f(x) + L(h) + R(h)$$

with L : linear and continuous

$$\frac{R(h)}{\|h\|_E} \xrightarrow{\|h\|_E \rightarrow 0} 0$$

$$x = (x_1, x_2) \text{ with } x_1 \in E, x_2 \in F$$

$$f(x) = f(x_1, x_2)$$

$$x = (x_1, x_2) \quad h = (h_1, h_2)$$

$$f(x+h) - f(x) = f(x_1 + h_1, x_2 + h_2) - f(x_1, x_2)$$

Since f is bilinear, we can write

$$\begin{aligned} f(x+h) - f(x) &= f(x_1, x_2 + h_2) + f(h_1, x_2 + h_2) \\ &\rightarrow f(x_1, x_2) \end{aligned}$$

$$= f(x_1, h_2) + f(h_1, x_2) + f(h_1, h_2)$$

$$\text{let } L(h) = f(x_1, h_2) + f(h_1, x_2)$$

L is linear $\lambda, \mu \in \mathbb{R}, h, k \in \mathbb{R}^2$

$$L(h+k)$$

$$= f(x_1, h_2 + k_2) + f(h_1 + k_1, x_2)$$

$$= f(x_1, h_2) + f(x_1, k_2) + f(h_1, x_2) + f(k_1, x_2)$$

$$= f(x_1, h_2) + f(h_1, x_2) + f(h_1, k_2) + f(k_1, x_2)$$

$$= L(h) + L(k)$$

$$L(\lambda h) = f(x_1, \lambda h_2) + f(\lambda h_1, x_2)$$

$$= \lambda f(x_1, h_2) + \lambda f(h_1, x_2)$$

$$= \lambda (f(x_1, h_2) + f(h_1, x_2))$$

$$= \lambda L(h)$$

L is also continuous as it is a sum of continuous functions

$$h_1 \rightarrow f(h_1, x_2)$$

$$\text{and } h_2 \rightarrow f(x_1, h_2)$$

$$\text{let } E(h) = f(h_1, h_2)$$

we study the limit of

$$R(h) = \frac{\|E(h_1, h_2)\|_F}{\|h\|_{E \times E}}$$

since f is linear continuous, there exists a positive constant K such that

$$\|E(h_1, h_2)\|_F \leq K \|h_1\|_E \|h_2\|_E$$

$$\text{so } R(h) \leq K \frac{\|h_1\|_E \|h_2\|_E}{\|h\|_{E \times E}}$$

$$\text{with } \|h\|_{E \times E}^2 = \|h_1\|_E^2 + \|h_2\|_E^2$$

$$\begin{aligned} \text{Noticing that } (a-b)^2 &= a^2 + b^2 - 2ab \geq 0 \\ &= ab \leq \frac{1}{2}(a^2 + b^2) \end{aligned}$$

we obtain

$$R(h) \leq \frac{K}{2} \|h\|_{E \times E}$$

$$R(h) \xrightarrow[\|h\|_{E \times E} \rightarrow 0]{} 0$$

We concluded that f is differentiable at α and

$$\begin{aligned} df_w &\xrightarrow[E \times E]{} F \\ (h_1, h_2) &\mapsto f(x_1, h_1) + f(h_2, x_2) \end{aligned}$$

Exercise 6.4: Linear Model

Consider the generalized linear model

$$Y = X\theta + \epsilon,$$

where the design matrix X is of size $n \times p$ and rank p , the parameters θ are in \mathbb{R}^p , and the errors ϵ are distributed according to $\mathcal{N}(0, \Sigma)$ with a positive definite covariance matrix Σ .

We define the generalized least squares criterion used to estimate

$$J(\theta) = (Y - X\theta)^T \Sigma^{-1} (Y - X\theta).$$

1. Justify that this criterion is of class C^∞ .
2. Determine the first and second-order differentials of this criterion.
3. Equip \mathbb{R}^p with the usual inner product. Specify the gradient and the Hessian matrix of J in the canonical basis of \mathbb{R}^p .
4. Verify that the Hessian matrix of J is positive definite.

Linear model:

$$Y = X\theta + \varepsilon$$

where

$$\begin{cases} Y \text{ is in } \mathbb{R} \\ X \in \mathbb{M}_{n \times p} \\ \theta \in \mathbb{R}^p \\ \varepsilon \text{ is in } \mathbb{R}^n \end{cases}$$

$$\varepsilon \sim N(0, \sigma^2 I)$$

MLE of θ minimizes

$$\left\{ \begin{array}{l} J(\theta) = \frac{1}{2} \|Y - X\theta\|^2 \\ \text{least square criterion} \end{array} \right.$$

(1) if we assume that $\varepsilon \sim N(0, \Sigma)$

where Σ is a positive definite matrix

then the MLE of θ minimizes the generalized least squares criterion

$$J(\theta) = (Y - X\theta)^T \Sigma^{-1} (Y - X\theta)$$

Let us show that J is differentiable $\theta \in \mathbb{R}^n$

$$\begin{aligned} J(\theta + h) &= \frac{1}{2} \|Y - X(\theta + h)\|^2 \\ &= \frac{1}{2} \|Y - X\theta - Xh\|^2 \\ &= \frac{1}{2} \langle Y - X\theta - Xh, Y - X\theta - Xh \rangle \\ &= \frac{1}{2} (\|Y - X\theta\|^2 + \|Xh\|^2 - 2 \langle Y - X\theta, Xh \rangle) \end{aligned}$$

$$\text{then } J(\theta + h) - J(\theta) = \langle Y - X\theta, Xh \rangle + \frac{1}{2} \|Xh\|^2$$

$$\text{Let } L(h) = \langle X\theta - Y, Xh \rangle = (X\theta - Y)^T Xh$$

$$\begin{aligned}
 &= (X + (X\theta - Y))^\top h \\
 &= \langle X^\top(X\theta - Y), h \rangle \\
 &= \langle X^\top X\theta - X^\top Y, h \rangle
 \end{aligned}$$

L is linear and since it is defined on a finite dimensional space it is continuous.

Let us study the limit

$$R(R) = \frac{1}{2} \frac{\|Xh\|^2}{\|Rh\|}$$

Since $R \in M^{n \times n}$.

as $R \rightarrow 0$

$$\begin{aligned}
 \|Xh\|^2 &= \langle Xh, Xh \rangle \\
 &= h^\top X^\top X h
 \end{aligned}$$

II- Higher order differentials

E, F

$$f: E \rightarrow F$$

Def

Let U be an open set of E . If f is differentiable on every point x of U , we say that f is differentiable on U .

If f is differentiable on U , we can define on U the mapping

$$df: U \subset E \rightarrow \mathcal{L}(E, F)$$

$$x \mapsto df_x$$

Def! Let df be defined on $U \subset E$, and $x \in U$.

If df is differentiable at x , then we can define

$d(df)$

In this case, we say that f is twice differentiable at x ,
and we call $d(df)_x$ the second order differentiability of f at x .

We denote

$$d^2f_x = d(df)_x$$

Remark:

$$d^2f_x \in \mathcal{L}(E, \mathcal{L}(E, F))$$

More specifically

let $h \in E$

then $d^2f_x(h) \in \mathcal{L}(E, F)$

let $\ell \in E$

then $d^2f_x(h)(\ell) \in F$

$$\begin{aligned} g: (E, E) &\rightarrow F \\ (h, \ell) &\mapsto d^2f_x(h)(\ell) \end{aligned}$$

g is linear with respect to each argument.

it is also continuous on $E \times E$

we can identify d^2f with a bilinear continuous mapping

$$\mathcal{L}(E, \mathcal{L}(E, F)) = \mathcal{L}^2(E, F)$$

$\approx \mathcal{L}_2(E, F)$ space of

So we may use a notation above

By a recursion we can define differential of an arbitrary order n

let $d^n f^x$ denotes the derivative of order n at x
and $d^n f^x \in \mathcal{L}^n(E, F)$

Seite exercise

$$J_1(\theta) = \frac{1}{2} \|y - X\theta\|^2$$

$$J_2(\theta) = \frac{1}{2} (y - X\theta)^T \sum (y - X\theta)$$

$$J_n(\theta + h) = \frac{1}{2} \|y - X(\theta + h)\|^2$$

$$= \frac{1}{2} \|y - X\theta + Xh\|^2$$

$$= \frac{1}{2} \langle y - X\theta + Xh, y - X\theta + Xh \rangle \text{ fail over}$$

$$= \langle X^T X \theta - X^T y, h \rangle + \frac{1}{2} \|Xh\|^2$$

$$\lim \frac{R(h)}{\|h\|^2} = 0 ?$$

$$\text{Let } A = X^T X$$

A is symmetric

A is definite positive since $h^T A h = h^T X^T X h = \|Xh\|^2 \geq 0$

A is diagonalizable in a orthonormal basis (u_1, \dots, u_n)

and its eigenvalues $\lambda_1, \dots, \lambda_n$ are non negative

$$\text{So let } h = \sum_{i=1}^n \alpha_i u_i$$

$$Ah = A \left(\sum_{i=1}^n \alpha_i u_i \right) = \sum_{i=1}^n \alpha_i A u_i = \sum_{i=1}^n \alpha_i \lambda_i u_i$$

$$h^T Ah = \left(\sum_{i=1}^n \alpha_i u_i \right)^T \left(\sum_{i=1}^n \alpha_i \lambda_i u_i \right)$$

$$= \left(\sum_{j=1}^n R_j u_j \right)^\top \left(\sum_{i=1}^n R_i \lambda_i u_i \right)$$

$$= \sum_{i=1}^n \sum_{j=1}^n R_i R_j \lambda_i u_j^\top u_i$$

but $u_j^\top u_i = \langle u_j, u_i \rangle = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$

then $R^\top A R = \sum_{i=1}^n R_i^2 \lambda_i$

But $\lambda_i \leq \lambda_p$ if $i = 1, \dots, p$

so $R^\top A R \leq \lambda_p \sum_{i=1}^p R_i^2$

but $\sum_{i=1}^p R_i^2 = \|R\|^2$

$$R^\top A R \leq \lambda_p \|R\|^2$$

$$R(R) = \frac{1}{2} R^\top A R \quad \text{with } A = X^\top X$$

$$\frac{R(R)}{\|R\|} \leq \frac{\lambda_p \|R\|}{2} \quad \text{where } \lambda_p = \text{largest eigenvalue of } A$$

$$\lim \frac{R(R)}{\|R\|} = 0$$

Conclusion J is differentiable at θ and $\frac{dJ_\theta(R)}{dR} = \langle X^\top X \theta - X^\top Y, R \rangle$
if $R \in \mathbb{R}^p$.

Remark 1 The gradient of J is:

$$\nabla J(\theta) = X^\top X \theta - X^\top Y \in \mathbb{R}^p$$

$$dJ_\theta(R) = \langle X^\top X \theta - X^\top Y, R \rangle$$

$$\begin{aligned} d\bar{J}_{\theta+\eta}(R) - d\bar{J}_\theta(R) &= \langle X^\top X(\theta+\eta) - X^\top Y, R \rangle - \langle X^\top X \theta - X^\top Y, R \rangle \\ &= \langle X^\top X \theta + X^\top X \eta - X^\top Y, R \rangle - \langle X^\top X \theta - X^\top Y, R \rangle \\ &= \langle X^\top X \theta - X^\top Y, R \rangle + \langle X^\top X \eta, R \rangle - \langle X^\top X \theta - X^\top Y, R \rangle \\ &= \langle X^\top X \eta, R \rangle \end{aligned}$$

L: $R \rightarrow (R \mapsto \langle X^\top X \theta - X^\top Y, R \rangle)$ is linear, continuous with values

in $\mathcal{L}(\mathbb{R}^p)$

so L is the differential of dJ

Hence $dJ_\theta : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$(\mathbf{x}, \mathbf{R}) \mapsto \langle \mathbf{x}^T \mathbf{x}, \mathbf{R}, \mathbf{R} \rangle$$

Condition for the $J(\theta)$ to be definite positive

- 1) $\mathbf{x}^T \mathbf{x}$ has only positive eigenvalues
- 2) $\det(\mathbf{x}^T \mathbf{x}) > 0$
- 3) $\text{ker}(\mathbf{x}^T \mathbf{x}) = \{\mathbf{0}\}$
- 4) $\text{rank}(\mathbf{x}) = p$

Exercise 6.5: Logistic Regression

Consider a sample consisting of n independent random variables Y_i following a Bernoulli distribution with parameter p_i . We assume that for each i , the parameter p_i depends on explanatory variables $x_i \in \mathbb{R}^p$ and unknown parameters $\theta \in \mathbb{R}^p$ through the relation $p_i = \sigma(x_i^T \theta)$, where $\sigma(t) = \frac{1}{1+e^{-t}}$.

1. Show that

$$\sigma(t) + \sigma(-t) = 1$$

and that

$$(\log_e \circ \sigma)'(t) = \sigma(-t).$$

2. Determine the likelihood $\mathcal{L}(\theta; y_1, \dots, y_n)$ of the model.

3. For fixed y_i , set $\ell(\theta) = -\log_e(\mathcal{L}(\theta; y_1, \dots, y_n))$, the logistic loss function. Show that ℓ is of class C^∞ .

4. Show that the gradient $\nabla \ell(\theta)$ and the Hessian matrix $\nabla^2 \ell(\theta)$ of ℓ are given by:

$$\nabla \ell(\theta) = -\sum_{i=1}^n x_i(y_i - \sigma(x_i^T \theta)) \quad \text{and} \quad \nabla^2 \ell(\theta) = \sum_{i=1}^n \sigma(x_i^T \theta)(1 - \sigma(x_i^T \theta))x_i x_i^T.$$

5. Under what condition(s) is $\nabla^2 \ell(\theta)$ a positive definite matrix?

Logistic regression

$$Y_i \in \{0, 1\}$$

$$x_i \in \mathbb{R}^p$$

$$\begin{cases} P(Y_i=1) = \exp(\sigma(\langle x_i, \theta \rangle)) \\ \sigma(t) = \frac{1}{1+e^{-t}} \end{cases}$$

is the logit function

Assume Y_i are independent

$$\begin{aligned} \text{Likelihood of } (y_1, \dots, y_n) : \mathcal{L}(y_1, \dots, y_n, \theta) &= \prod_{i=1}^n P(Y_i=1)^{y_i} P(Y_i=0)^{1-y_i} \\ &= \prod_{i=1}^n \exp(\sigma(\langle x_i, \theta \rangle)). \\ &\quad \times \exp((1-y_i)(1-\sigma(\langle x_i, \theta \rangle))) \end{aligned}$$

$$= \exp(y_i \sigma(\langle x_i, \theta \rangle)) + (1-y_i) (1-\sigma(\langle x_i, \theta \rangle)).$$

the negative log-likelihood is

$$\left\{ \begin{array}{l} l(\theta) = -\log(L(\theta, y_1, \dots, y_n)) = -\sum_{i=1}^n [y_i \sigma(\langle x_i, \theta \rangle) + (1-y_i) (1-\sigma(\langle x_i, \theta \rangle))] \end{array} \right.$$

this function is known as the binary cross-entropy

c'est une bonne idée de travailler sur $l(\theta + h) - l(\theta)$

vaut mieux calculer les dérivées partielles

$$\frac{\partial g(\theta)}{\partial \theta} = \partial \frac{\sigma(\langle x, \theta \rangle)}{\partial \theta}$$

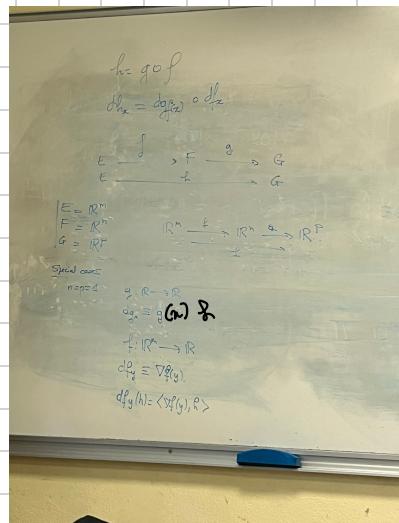
$$\frac{1}{1+e^{-}}$$

deux fonctions différentes

$$g = R \circ f$$

$$dg_x = df_x \circ dg_x$$

$$g'_x = f'(f(x)) f'(x).$$



$$dR_x = dg_{f(x)} \circ df_x$$