

MATH FOR ML : UNIT 2 (CALC + OPT) REVIEW

- I DERIVATIVE NOTIONS (total, gradient, Jacobian, etc.)
- II TAYLOR SERIES + IMPLICATIONS.
- III OPTIMIZATION w or w/o CONSTRAINTS.
- IV GRADIENT DESCENT.

DERIVATIVE NOTIONS

- Single-variable calculus : only two directions (pos. or neg.)
- Multi-variable calculus : inf. many directions

$$F: \mathbb{R}^d \rightarrow \mathbb{R}^n$$



second derivative
Hessian
 ↑ For $F: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla F(x_0)^T d$$

DIRECTIONAL DERIVATIVE

when $d=e_i$:

PARTIAL DERIVATIVE

(Scalar in \mathbb{R} that says "how fast F changes in d " direction.)

(change in F in a basis direction.)

In general:

$$\begin{bmatrix} -\nabla f_1(x)^T \\ \vdots \\ -\nabla f_n(x)^T \end{bmatrix}$$

In Jacobian case:

$$\begin{bmatrix} -\nabla F \\ \vdots \\ -\nabla F_n \end{bmatrix} \begin{bmatrix} | \\ d \\ | \end{bmatrix}$$

= $\begin{bmatrix} \text{how much } F_1 \text{ changes} \\ \vdots \\ \text{how much } F_n \text{ changes} \\ \text{in direction } d. \end{bmatrix}$

$$\nabla F(x) = \left(\frac{\partial F}{\partial x_1}, \dots, \frac{\partial F}{\partial x_d} \right)$$

| written down

when $n=1$

Gradient
 $\nabla f(x)^T \in \mathbb{R}^{1 \times d}$

when n is arbitrary.

Jacobian
 $\nabla f(x) \in \mathbb{R}^{n \times d}$

only make sense AT a fixed $\hat{x}_0 \in \mathbb{R}^d$.

TOTAL DERIVATIVE
 (Lin. Transformation)
 $Df_{\hat{x}_0}: \mathbb{R}^d \rightarrow \mathbb{R}^n$

is

is

gives...

TAYLOR SERIES.

$$\text{At } x_0: F(x) = F(x_0) + F'(x_0)(x-x_0) + \frac{F''(x_0)}{2!}(x-x_0)^2 + \dots$$

Lin. Approximation.

second order approx.

① **DESCENT LEMMA**
 $f(x_t) \leq f(x_{t-1}) - \frac{\mu}{2} \|\nabla f(x_{t-1})\|^2$

+ Remainder: Bounded by smoothness on Quad. Form.

③ **Proof of First order Convexity**

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

$$F(x) = F(x_0) + \nabla F(x_0)^T (x-x_0) + \frac{1}{2} (x-x_0)^T \nabla^2 F(x_0) (x-x_0) + \dots$$

② Opt. Necessary cond.

Opt. SUFF. cond.

$$\nabla F(x_0) = 0, \quad \nabla^2 F(x_0) \text{ PSD.}$$

$$\nabla F(x_0) = 0, \quad \nabla^2 F(x_0) \text{ PD.}$$

(otherwise: Candidate local min > some close point)

(PD / PSD Notation: $A \geq 0$ PSD, $A > 0$ PD)

OPTIMIZATION PROBLEMS

Necessary: This has to be true for P to be true.

Sufficient: If this is true then P is true.

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

x^* = global minimizer.

$$\begin{array}{l} \text{Minimize } F(\vec{x}) \\ \text{Subject to } \vec{x} \in \mathcal{C}. \end{array}$$

③ CONVEX $\mathcal{C} = \text{CONVEX}$
 $F = \text{CONVEX}$

All local Minima
are global minima.

GD is guaranteed to minimize!

local minima.

① UNCONSTRAINED ($\mathcal{C} = \mathbb{R}^d$)

(i) Necessary: $\nabla F(x^*) = 0$,
 $\nabla^2 F(x^*)$ PSD.

(ii) Sufficient: $\nabla F(x^*) = 0$,
 $\nabla^2 F(x^*)$ PD.

② CONSTRAINED $\mathcal{C} \subset \mathbb{R}^d$

A) Equality Constraints

$$\begin{array}{l} \text{minimize } F(\vec{x}) \\ \text{subject to } h_1(x) = 0 \\ \vdots \\ h_m(x) = 0 \end{array}$$

Solve using Lagrangian:

$$L(\vec{x}, \vec{\lambda}) = F(x) + \sum_{i=1}^m \lambda_i h_i(\vec{x})$$

(Mind regular) Reduce
Points: $L(x, \lambda)$

B) Inequality Constraints

Subject to $g_1(x) \leq 0$ solve using Lagrangian.

$$\dots g_m(x) \leq 0 \quad L(x, \lambda, \mu) = F(x) + \sum \lambda_i h_i(x) + \sum \mu_j g_j(x)$$

complementary slackness

$$\mu^* \geq 0 \quad \& \quad \mu_j^* g_j(x^*) = 0 \quad \forall j \in [r]$$

GRADIENT DESCENT

$$f(x_t) \leftarrow f(x_{t-1}) - \eta \nabla f(x_{t-1})$$

- Taylor's Thrm
- Smoothness of f
(Bounded Hessian eigenvals).

IF CONVEX F .

...

DESCENT LEMMA

$$f(x_t) \leq f(x_{t-1}) - \frac{\eta}{2} \|\nabla f(x_{t-1})\|^2$$

("Updates make the function smaller if $\nabla f(x_{t-1}) \neq 0$ ".)

GD ON CONVEX

$$F(x_t) \leq \frac{B}{T} \|x_0 - x^*\|^2 + F(x^*)$$

("GD eventually reaches the global minimum".)

- Use Descent Lemma.
- First-order Def. of Convexity.
- Potential Argument.

QUADRATIC FORM.

Quadratic Form $x^T A x \leftarrow \text{In } d=1: x a x = a x^2$

Quadratic Function $x^T A x + b^T x + c$. (A is matrix, b vector, c scalar)

Taylor Series: $f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2}(x-x_0)^2 + \dots$

$f^d \downarrow$

$f(x) = f(x_0) + \nabla f(x_0)^T (x-x_0) + (x-x_0)^T \nabla^2 f(x_0) (x-x_0)$

① MVC Quad. Function: $x^T A x + b^T x + c$

$$a x^2 + b x + \frac{b^2}{4a}$$

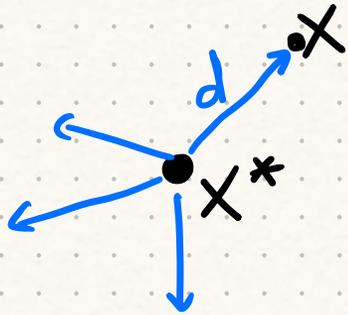
② Single Variable Quad: $a x^2 + b x + c \rightarrow \left(\sqrt{a} x + \frac{b}{2\sqrt{a}} \right)^2 - \frac{b^2}{4a} + c$

$$\left(\sqrt{a} x + \frac{b}{2\sqrt{a}} \right)^2 - \left(\frac{b^2}{4a} - c \right)$$

HS: $(x-\alpha)^2 + \beta \Rightarrow$ so b and c are only responsible for shifting the parabola.

INTUITION FOR NECESSARY / SUFFICIENT.

From top down on \mathbb{R}^d



$$F(x) = F(x^*) + \nabla F(x^*)^T (x - x^*) + (x - x^*)^T \nabla^2 F(x^*) (x - x^*)$$



$$F(x^* + d) = F(x^*) + \nabla F(x^*)^T d + d^T \nabla^2 F(x^*) d.$$

SUFFICIENT: $\nabla F(x^*) = 0$ \rightarrow Minimum.
 $\nabla^2 F(x^*)$ PD

\Rightarrow

$$F(x^* + d) = F(x^*) + \underbrace{\nabla F(x^*)^T d}_0 + \underbrace{d^T \nabla^2 F(x^*) d}_{> 0}$$

$$F(x^* + d) = F(x^*) + \text{Positive \#}$$

NECESSARY: $\nabla F(x^*) = 0$ \leftarrow Minimum.
 $\nabla^2 F(x^*)$ PSD

$$F(x^* + d) = F(x^*) + 0 + \text{Negative \#}$$

\rightarrow for some d .