

(1)

7.1 Continuous Optimization

- 7.1 Opt using gradient Desc $227 - 233^6$
- 7.2 Constrained opt & lagrange $233 - 236^3$

As ML is implemented by computer —
the math formulations are expressed
as numerical optimization methods

Training ML = find good parameters

"good" is defined by the objective func

2 main branches of cont opt —
unconstrained & constrained

Assumption that obj func is differentiable

Hence we have a gradient @ each
location in the space to help us
find the optimum value

(2)

by convention, most b) funcs
in ML are intended to be
minimized

global minima

Gradients can be used to help
find the minimum - indicating
whether we should step left or
right

We can solve for all stationary
points of a function by calc'ing
derivatives & setting to 0

Stationary points are real ~~pos~~
roots of the derivatives & points that
have zero gradient

(3)

Example:

$$l(x) = x^4 + 7x^3 + 5x^2 - 17x + 3$$

$$\frac{dl(x)}{dx} = 4x^3 + 21x^2 + 10x - 17$$

this is a cubic equation, it has in general 3 solutions when set to 0

in example 2 = min, 1 = max

to check whether stationary point is a min or max, we take 2nd Derivative & check whether it eq is pos or neg at stationary point

$$\frac{d^2l(x)}{dx^2} = 12x^2 + 42x + 10$$

by substituting est values $x = -4.5, -1.4, 0.7$ we see that middle point < 0
so it is a max
and there's one minimum

(4)

Note - low-order polynomials are solvable analytically high-orders either aren't or are too costly

As a ~~re~~ find the result by starting @ some value $x_0 = -6$ & follow the negative gradient

Gradient tells us the direction to go but not how far (step size)

Furthermore if there are ≥ 1 minimum then our start point could bring us down to the wrong stationary minima

(5)

7.1 Optimization Using Grad Desc

Consider problem of solving real-valued function: $\min F(x)$

$F: \mathbb{R}^d \rightarrow \mathbb{R}$ = ML Problem

F is Differentiable but no Analytical solution

Gradient Descent is a first order opt alg

- This means it only uses the first order derivative
- this is cheap, simple & applicable

- Gradient gives direction of steepest Ascent
- the opposite gives steepest Descent

Gradient points in a direction that is orthogonal to the contour lines lines we wish to optimize

$$w_t = w_0 - \gamma (\nabla F(w_0))^+$$

Start w/ init w_0 & iterate

(6)

for a suitable step size the sequence will converge to a minima

Note, near to the minima
GD is slow & the rate is
inferior to many other methods.
LizZag many times around the
min point

(2)

7.2 Constrained Opt & Lagrange Multi

Previous 7.1 solve min of function

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{where } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Consider additional constraints

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Subject to $g_i(\mathbf{x}) \leq 0$ for all $i=1, \dots, m$

One way of turning constrained prob into unconstrained is to use indicator

$$O(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m I(g_i(\mathbf{x}))$$

$I(z)$ is an infinite step func

$$I(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

But this route is also difficult to solve for

(8)

the solution is to introduce Lagrange multipliers

Idea of Lagrange is to replace step function w/ a linear one

Start $\min f(x)$ Subj to $g_i(x) \leq 0$

Lagrangian:

intro λ to each inequal const

$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x)$$

$$\text{Vectorised} = F(\alpha) + \lambda^T g(x)$$

Introduction of Lagrangian Duals

In opt, Dual = convert prob in one set of vars x (Primal) into another problem w/ new set of vars (λ) (Dual)

(Primal)

Start: $\min_{\mathbf{x}} f(\mathbf{x})$ Subj to $g_i(\mathbf{x}) \leq 0$
for all $i = 1, \dots, m$

Dual: $\max_{\lambda \in \mathbb{R}^m} D(\lambda)$ Subj to $\lambda \geq 0$

Dual = $\lambda \cdot D(\lambda) = \min_{\mathbf{x}} f(\mathbf{x}, \mathbf{\lambda})$
var