

①

$$J(\theta) = \mathbb{E}_{(x,y) \sim P_{data}} [L(f(x, \theta), y)]$$

$$\approx \frac{1}{n} \sum_{i=1}^n L(f(x^{(i)}, \theta), y^{(i)})$$

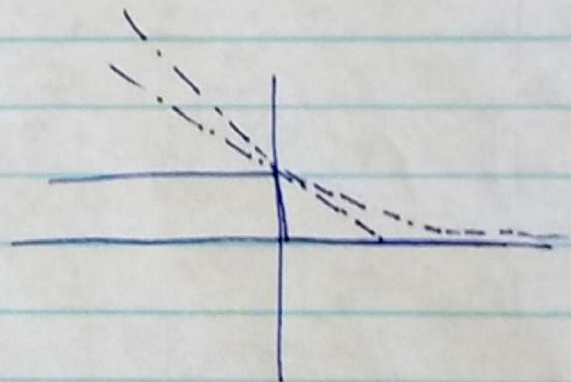
Examples:

$$L(y, \hat{y}) = -\text{Sign}(y \hat{y})$$

Surrogate:

$$L(y, \hat{y}) = (1 - y \hat{y})_+$$

$$L(y, \hat{y}) = (1 - y \hat{y})_+^2$$



$$(x)_+ = \begin{cases} x & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

the surrogate is differentiable and more robust.

$$\nabla_{\theta} J = \mathbb{E}_{(x,y) \sim P_{data}} \left[\nabla_{\theta} L(f(x, \theta), y) \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(f(x, \theta), y_i)$$

Why we care about "Stochastic Gradient Descent"?

— Time consuming to compute $\nabla_{\theta} J$ exactly (full pass)

2

— variance of estimate $\frac{\sigma}{\sqrt{n}}$ but

Cost increases linearly

— How to choose batch size

- learning rate for small batch
- large batch may not fit
- Small batch can function as regularizer

Ill-conditioning

$$f(x) \approx f(x^0) + (x-x^0)^T g + \frac{1}{2} (x-x^0)^T H (x-x^0)$$

$$f(x^0 - \epsilon g) \approx \underbrace{f(x^0)}_{\text{Current value}} - \epsilon g^T g + \frac{1}{2} \epsilon^2 g^T H g$$

if $\frac{1}{2} \epsilon^2 g^T H g$ dominates $\epsilon g^T g$ it is ill-cond

~~It~~ In practice, it might be fine. (Fig 8.1)

In Convex opt, ~~but~~ we can handle this

with Newton method but not in

deep learning

→ b/c Newton method finds critical point which can be saddle point

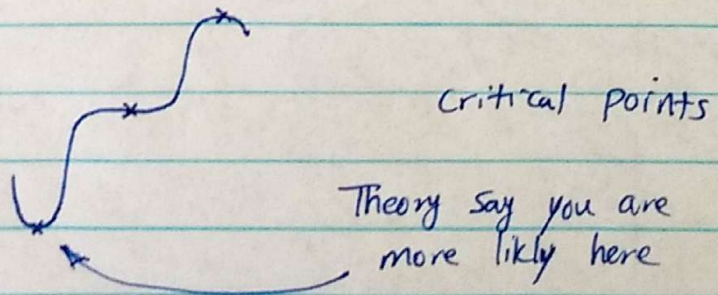
(3)

Local Min.

- Because of symmetry not all local min are bad
- They are bad if they are very different than global -- but we cannot tell
- Take home message: Don't blame everything on local min

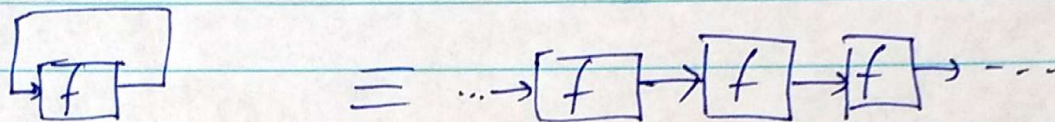
→ Plot $\|g\|_2^2$

- In high dim they are rare anyway b/c we have saddle points more often



(4)

Deep recurrent network:



$$h(\theta) = f(f(f(\dots f(x) \dots))$$

$$\nabla h = \nabla_{\theta} f(x) \nabla_{\theta} f(f(x)) \dots$$

Let's assume at the beginning $f(x) \approx x$

$$\nabla f(x) = \nabla f(f(x)) = \dots = W$$

$$\nabla h = W^t$$

$$W = V \text{diag}(\lambda) V^{-1}$$

$$W^t = V \text{diag}(\lambda^t) V^{-1}$$

$$\text{if } |\lambda_i| < 1$$

$$\rightarrow |\lambda_i|^t \rightarrow 0$$

vanishing gradient

$$\lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

$$\text{if } |\lambda_i| > 1$$

$$\rightarrow |\lambda_i|^t \rightarrow \infty$$

exploding gradient

5

SGD:

$$\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \epsilon_k g$$

\uparrow Step length \nwarrow gradient

Sufficient Cond. for Conv.

$$\sum_{k=1}^{\infty} \epsilon_k = \infty \rightarrow O(k^{-1})$$

$$\sum_{k=1}^{\infty} \epsilon_k^2 < \infty \rightarrow O(k^{-\frac{1}{2}})$$

$$O(k^{-\frac{1}{2}}) < \epsilon_k < O(k^{-1})$$

$$\epsilon_k = (1 - \frac{\tau}{k}) \epsilon_0 + \frac{\tau}{k} \epsilon_\tau$$

choose $\tau, \epsilon_0, \epsilon_\tau$

Good luck! 😊

How should we think about SGD

Regret $J(\theta^k) - \min_{\theta} J(\theta) \iff$ How fast it decays

$O(\frac{1}{k})$ for strong convex

$O(k^{-\frac{1}{2}})$ Convex

non-convex ??