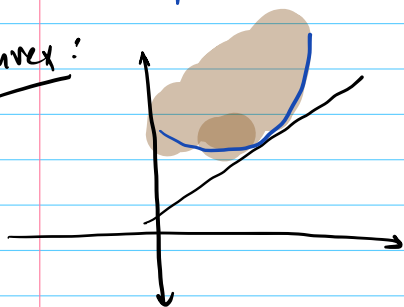Risk $R(\theta) = E[\ell(Y, \theta)]$

Optimization
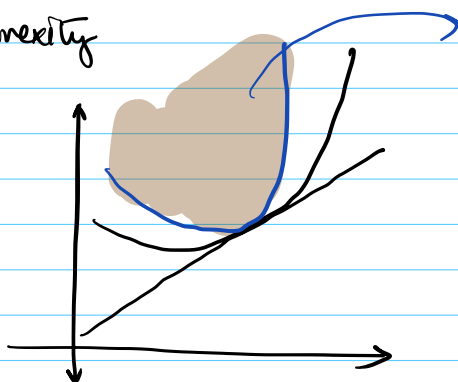$$\nabla R(\theta_0) = 0$$

Assumption $R(\cdot)$ is strongly convex

<span style="color:orange">inner product</span>

Convex:

$$f(x) \geq f(x^*) + \langle \nabla f(x^*), x - x^* \rangle$$
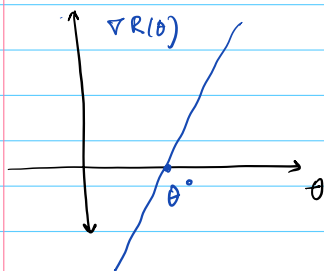
Strong Convexity
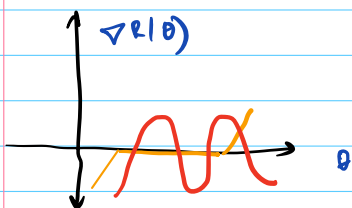
Is also above the parabola.

(rules out flat pts.)

$$f(x) \geq f(x^*) + \langle \nabla f(x^*), x - x^* \rangle + \frac{\alpha}{2} \| x - x^* \|^2$$
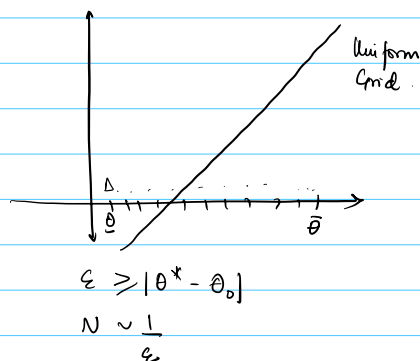
We want $\nabla R.(\theta)$ to look like this.

$\nabla R(\theta)$

$\theta$

$\theta^{\circ}$

Strong Convexity rules out cases like :—

$\nabla R(\theta)$

$\theta$

1.  Grid search. $[\underline{\theta}, \bar{\theta}]$

$$\theta_i = \underline{\theta} + i\Delta$$

(A) $\min_{\theta_i} \| \nabla R(\theta_i) \|$ , N steps

$$\Delta = \frac{\bar{\theta} - \underline{\theta}}{N} \quad (\text{Accuracy})$$

uniform Grid

$\Delta$

$\underline{\theta}$

$\bar{\theta}$

$\varepsilon \geqslant |\theta^* - \theta_0|$

$N \sim \frac{1}{\varepsilon}$

You do not have a closed form.

2. Bisection Algorithm

$a_0 = \underline{\theta}$ , $b_0 = \bar{\theta}$

$c_1 = \dfrac{\bar{\theta} + \underline{\theta}}{2}$

$\text{sign}(\nabla R(a_0)) \neq \text{sign}(\nabla R(b_0))$

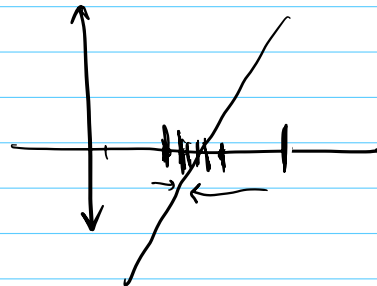If $\text{sign}(\nabla R(b_0)) \neq \text{sign}(\nabla R(c_1))$

$\quad a_1 = c_1, \quad b_1 = b_0$

$a_k, b_k, \quad c_{k+1} = \dfrac{a_k + b_k}{2}$

If $\text{sign}(\nabla R(c_{k+1})) \neq \text{sign}(\nabla R(b_k))$

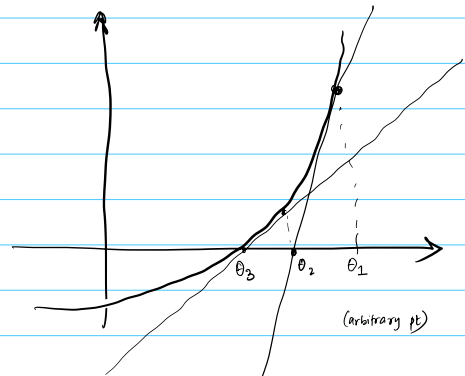$\quad a_{k+1} = c_{k+1}, \quad b_{k+1} = b_k \qquad [a_k, b_k]$

$|b_k - a_k| = \varepsilon = \dfrac{|\bar{\theta} - \underline{\theta}|}{2^k} \qquad N \sim \log \dfrac{1}{\varepsilon}$

3.    Newton Algorithm

$$H_\theta = \left( \frac{\partial^2 R(\theta)}{\partial \theta_i \, \partial \theta_j} \right)_{ij}$$

In one dimension, this would just be double derivative.

⊛ Iteration $R$: $\theta_k$

$$\nabla R(\theta_k) + H(\theta_k)(\theta - \theta_k) = 0$$

$$\theta_{k+1} = \theta_k - H(\theta_k)^{-1} \nabla R(\theta_k)$$

$$\underbrace{\Delta R(\theta_k)} \approx \underbrace{R'(\theta_k) \cdot (\theta_{k+1} - \theta_k)}_{} + \frac{1}{2} \underbrace{R''(\theta_k) \cdot (\theta_{k+1} - \theta_k)}_{}$$

$$\underline{R(\theta_{k+1}) - R(\theta_k)} \approx \underbrace{\nabla R(\theta_k)}_{} \qquad \underbrace{H(\theta_k)}_{}$$

$$0 =$$

$$R(\theta) = E[\ell(Y, \theta)]$$

we can't keep along the R.V.

A sample: $\{Y_i\}_{i=1}^N$ which has a cutoff → 1) data collection is expensive; 2) most of the data is historical; there is only so much data then.

## Analogy Principle

a set of $N$ numbers

We take the sample & create a new distribution "replicating" sampling scheme $\{y_i\}_{i=1}^N$

New r.v. → $P(Z = y_i) = \dfrac{1}{N}$   ($Z$ takes value with equal probability)

empirical risk:

hats are for empirical stuff →

$$\hat{R}(\theta) = E_N[\ell(Z, \theta)]$$
$$= \frac{1}{N} \sum_{i=1}^N \ell(y, \theta)$$