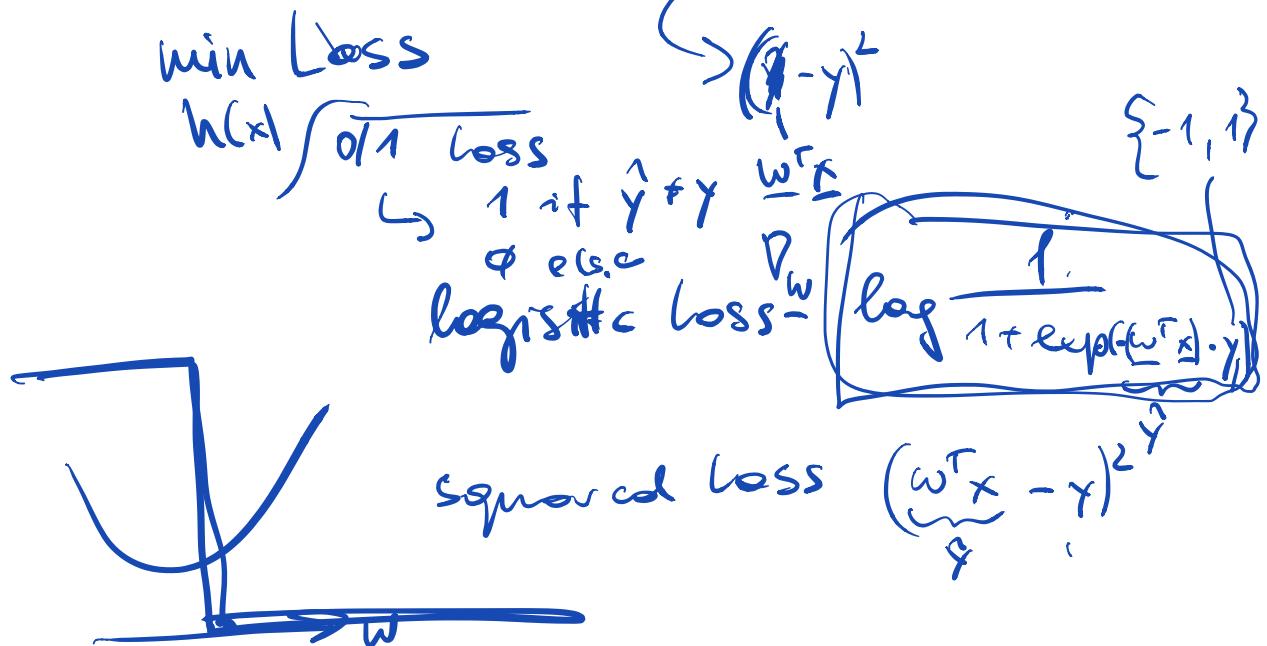


$$\begin{array}{c}
 \text{Data} \\
 \begin{matrix}
 \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_n \\
 y_1 & y_2 & \dots & y_n
 \end{matrix}
 \end{array}
 \quad
 \begin{array}{l}
 \text{hypothesis} \\
 \text{space} \\
 h(\underline{x}) = \underline{w}^T \underline{x}
 \end{array}
 \quad
 \begin{array}{l}
 \text{loss-fct.} \\
 \text{squared loss} \\
 (y - \hat{y})^2
 \end{array}$$

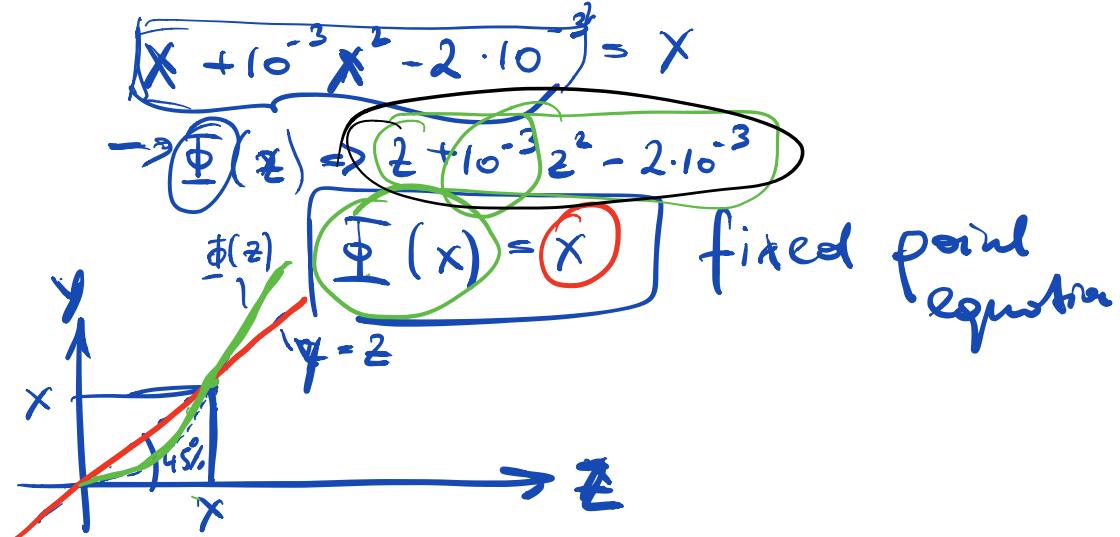


Data	Hypothesis Space	Loss-Fct.
$(\underline{x}_1^{(1)}, y^{(1)}) \dots (\underline{x}_m^{(m)}, y^{(m)})$	$h(\underline{x}) = \underline{w}^T \underline{x}$	$(\hat{y} - y)^2$
Linear Regression		
$\min_{\underline{w} \in \mathbb{R}^n}$	$\frac{1}{m} \sum (\underline{w}^T \underline{x}_i^{(i)} - y^{(i)})^2$	
		training error $\sum (\hat{y}_i - y_i)^2$

find number x s.t. $x^2 = 2$
 allow odds/mults.

$$x^2 = 2$$

$$+x/-2 \cdot 10^{-3} \cancel{+10^{-3}} \quad 10^{-3} \cdot x^2 = 2 \cdot 10^{-3}$$



fixed point iteration

initial guess $x^{(0)}$ $\rightarrow x^{(1)} = \underline{\Phi}(x^{(0)}) \rightarrow x^{(2)} = \underline{\Phi}(x^{(1)})$

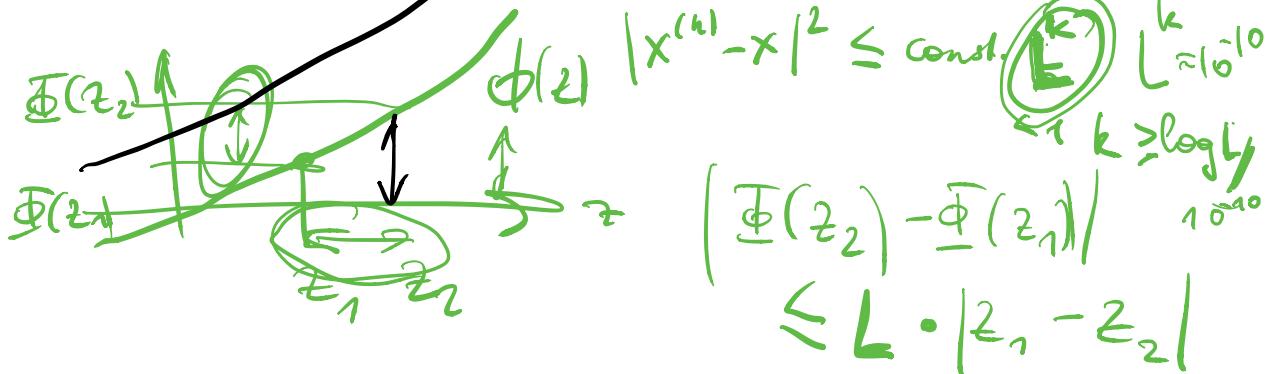
$$\boxed{x^{(k+1)} = \underline{\Phi}(x^{(k)})}$$

if $\underline{\Phi}$ is Lipschitz ~~constant~~ continuous with

constant $L < 1 \Rightarrow$

$$x^{(k)} \rightarrow x = \sqrt{2}$$

$$|x^{(k)} - x|^2 \leq \text{const. } L^k \quad L \approx 10^{-10}$$



$$\min_{\underline{w}} \mathcal{E}(\underline{w}) = \frac{1}{m} \| \underline{y} - \underline{X} \underline{w} \|_2^2$$

$\begin{pmatrix} \underline{y}^{(1)} \\ \vdots \\ \underline{y}^{(m)} \end{pmatrix}$ $(\underline{X}^{(n)})^\top$

$$-\alpha \cdot \nabla \mathcal{E}(\underline{w}^{(opt)}) = \phi$$

$$\underline{w}^{(opt)} - \alpha \nabla \mathcal{E}(\underline{w}^{(opt)}) = \underline{w}^{(opt)}$$

$\Phi(\underline{w}) \rightarrow \underline{w} - \alpha \nabla \mathcal{E}(\underline{w}) \rightarrow \underline{w}^{(k+1)} = \Phi(\underline{w}^{(k)})$

step-size / learning rate

$$\nabla \mathcal{E}(\underline{w}) = -\frac{2}{m} \underline{X}^\top (\underline{y} - \underline{X} \underline{w})$$

$$\Phi(\underline{w}) = \underline{w} + \left(\frac{2\alpha}{m} \right) [\underline{X}^\top \underline{y} - \underline{X}^\top \underline{X} \underline{w}]$$

$$\underline{I} \cdot \underline{w}$$

$$= \left[\underline{I} - \frac{2\alpha}{m} \underline{X}^\top \underline{X} \right] \underline{w} + \left(\frac{2\alpha}{m} \right) \underline{X}^\top \underline{y}$$

doesn't depend
on \underline{w}

$$\underline{I}(\underline{w}^{(1)}) - \underline{\Phi}(\underline{w}^{(2)}) =$$

$$\left\| \left(\underline{I} - \frac{2\alpha}{m} \underline{X}^\top \underline{X} \right) \left(\underline{w}^{(1)} - \underline{w}^{(2)} \right) \right\| \leq L \left\| \underline{w}^{(1)} - \underline{w}^{(2)} \right\|$$

matrix \underline{C} vector \underline{w}

$$\|\underline{C}\underline{w}\| \leq \|\underline{C}\| \cdot \|\underline{w}\|$$

Gilbert Strang

$$\|\underline{C} = \frac{2\alpha}{m} \underline{X}^T \underline{X}\|_2$$

spectral norm

$$\max_{\text{eigenvectors}} |\lambda_L|$$

$$\lambda_L = 1 - \frac{2\alpha}{m}$$

$$\underline{X} \in \left(\begin{array}{c} (\underline{x}^{(1)})^T \\ \vdots \\ (\underline{x}^{(m)})^T \end{array} \right)$$

positive semi-definite eigenvector matrix

$$\begin{array}{ccccccccc} & & & & & & & & \\ \text{---} & \text{---} \\ | & x & x & x & x & x & x & x & x \end{array}$$

$$\left(\frac{2\alpha}{m} \underline{X}^T \underline{X} \right)$$

$$\left(\frac{1}{m} \underline{X}^T \underline{X} \right) \leq 2$$

max eigenvalue

$$\underline{C} = \begin{pmatrix} \frac{1}{m} \sum_{i=1}^m x_1^{(i)} x_1^{(i)} \\ \vdots \\ \frac{1}{m} \sum_{i=1}^m x_n^{(i)} x_n^{(i)} \end{pmatrix}$$

empirical cross correlation matrix of features

Momentum

$$\underline{\underline{w}}^{(k+1)} = \underline{\underline{w}}^{(k)} - \alpha \nabla \mathcal{E}(\underline{\underline{w}}^{(k)}) + \beta (\underline{\underline{w}}^{(k)} - \underline{\underline{w}}^{(k-1)})$$

new guess
for optimal weights gradient descent momentum

Deep learning
Book

Stochastic Gradient Descent

$$\underline{\underline{w}}^{(k+1)} = \underline{\underline{w}}^{(k)} - \alpha \nabla \mathcal{E}(\underline{\underline{w}}^{(k)}) + \beta (\underline{\underline{w}}^{(k)} - \underline{\underline{w}}^{(k-1)})$$

$$\mathcal{E}(\underline{\underline{w}}) \stackrel{!}{=} \frac{1}{m} \| \mathbf{y} - \mathbf{X} \underline{\underline{w}} \|^2$$

$$\nabla \mathcal{E}(\underline{\underline{w}}) = \left(-\frac{2}{m} \mathbf{X}^T (\mathbf{y} - \mathbf{X} \underline{\underline{w}}) \right)$$

$$-\frac{2}{m} \sum_{i=1}^m \mathbf{x}^{(i)} (\mathbf{y}^{(i)} - \mathbf{x}^T \underline{\underline{w}}^{(i)}) \quad \text{Squared Loss} \quad \sum_{i=1}^m \mathbf{x}^T \mathbf{x}^{(i)}$$

$$\approx \text{randomly select subset } \mathcal{B} \subseteq \{1, \dots, m\}$$

$$\left(\sum_{i \in \mathcal{B}} \mathbf{x}^{(i)} (\mathbf{y}^{(i)} - \mathbf{x}^T \underline{\underline{w}}^{(i)}) \right) \quad n/m$$

$$\nabla \mathcal{E}(\omega) = -\frac{2}{m} \left(\underline{\underline{X}}^T \underline{\underline{Y}} - \underline{\underline{X}}^T \underline{\underline{X}} \underline{\omega} \right)$$

pre-compute $\underline{\underline{X}}^T \underline{\underline{X}} \approx \text{cov}(\underline{\underline{X}})$

Covariance matrix