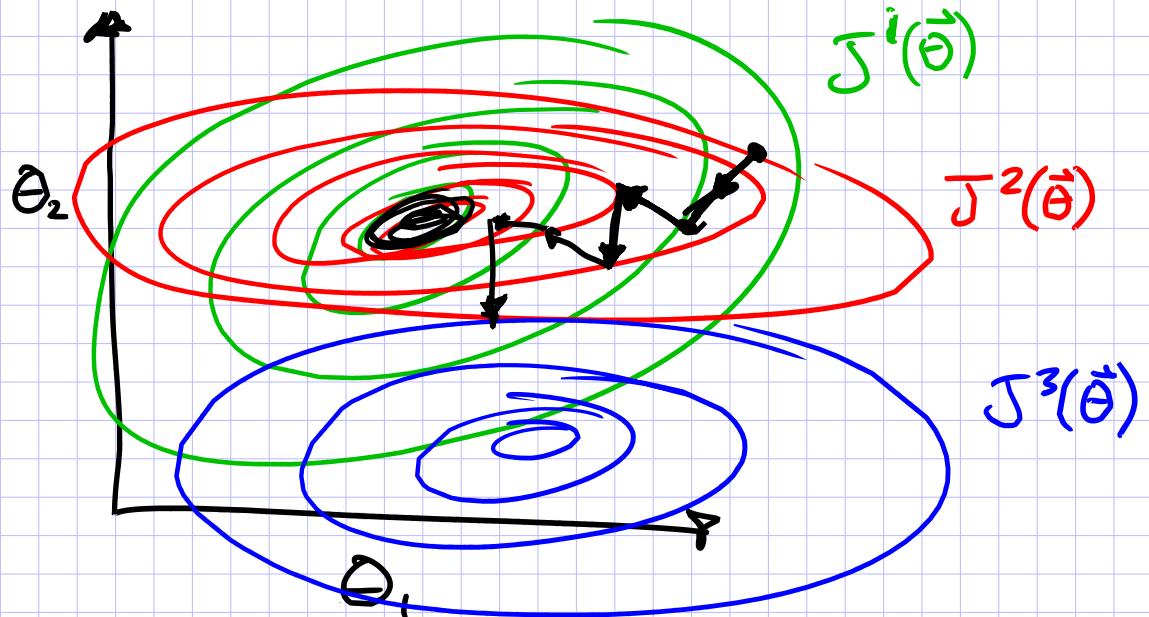
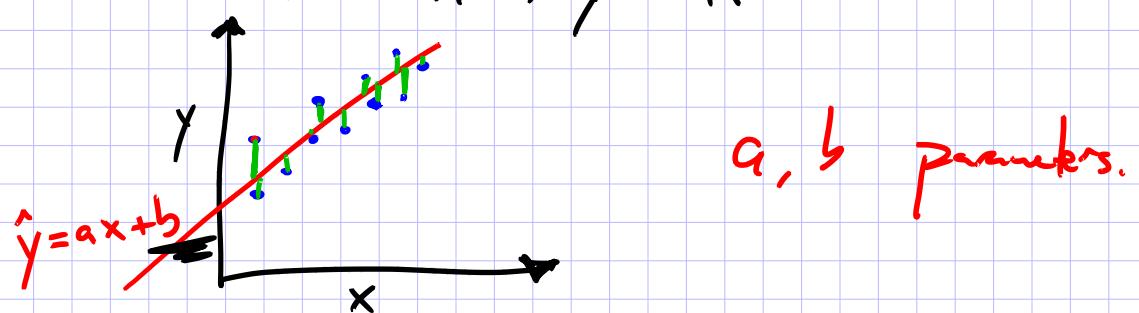


Jelinek Summer School



Linear Reg.

Data $D = \{(\vec{x}^{(i)}, y^{(i)})\}_{i=1}^N = \{(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots\}$

 $\vec{x}^{(i)} \in \mathbb{R}^M, y^{(i)} \in \mathbb{R}$


Linear Fn.

$$y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_M x_M + \theta_0$$

$$y = \underline{\theta_0} + \sum_{n=1}^M \theta_n x_n = \underline{\theta_0} + \vec{\theta}^T \vec{x}$$

$$y = \vec{\theta}^T \vec{x} \text{ if } x_0 = 1$$

$$\frac{\partial h(\vec{x})}{\partial \theta_m} = x_m$$

Lin. Reg. as Fn Approx.

① Assume D generated $\vec{x}^{(i)} \sim p^*(\vec{x})$ unk.
 $\vec{y}^{(i)} = h^*(x^{(i)})$

② Choose Hyp. Space., H

$$H = \{\vec{h}(\vec{x}) = \vec{\theta}^T \vec{x} : \vec{\theta} \in \mathbb{R}^M\}$$

Space of all lin. funs in M-dim

$$\frac{d J^{(i)}(\theta)}{d \theta_m} = \frac{1}{N} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) x_m^{(i)}$$

③ Choose Objective J_m .

$$J(\vec{\theta}) = \sum_{i=1}^N \frac{1}{N} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2$$

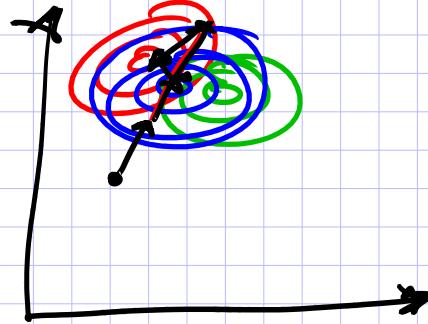
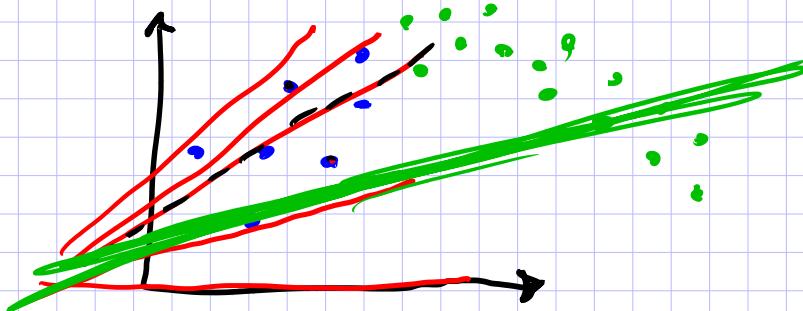
MSE = $\frac{1}{N} \sum_{i=1}^N (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2$

$O(M)$

④ Solve unconstrn opt. prob. (SGD, GD, closed-form)
 $\hat{\vec{\theta}} = \underset{\vec{\theta}}{\operatorname{argmin}} J(\vec{\theta})$

⑤ Predict: $\hat{y} = h(\vec{x}) = \vec{\theta}^T \vec{x}$

SGD:



$O(N^{2.783\dots})$

Closed Form:

$$\vec{\theta} = (\vec{X}^T \vec{X})^{-1} (\vec{X}^T \vec{y})$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$$\vec{X} = \begin{bmatrix} \vec{x}^{(1)} \\ \vdots \\ \vec{x}^{(N)} \end{bmatrix}$$

$N \times M$

Log. Reg.

$$P(D) = \prod_{i=1}^{N/2} p(y^{(i)} | \vec{x}) \prod_{i=N/2+1}^N p(y^{(i)} | \vec{x})$$

$$1 = P(y=0 | \vec{x}^{(i)}) + P(y=1 | \vec{x}^{(i)})$$

$$= \prod_{i=1}^{N/2} (1 - p(y=1 | \vec{x})) \prod_{i=N/2+1}^N p(y=1 | \vec{x})$$

$$\begin{aligned} J(\theta) &= -\log p(D) \\ &= -\log \prod_{i=1}^N p(y^{(i)} | \vec{x}^{(i)}) \\ &= \sum_{i=1}^N -\log p(y^{(i)} | \vec{x}^{(i)}) \\ &= \sum_{i=1}^N \frac{-\log p(y^{(i)} | \vec{x}^{(i)})}{J^{(i)}(\vec{\theta})} \end{aligned}$$

Multinomial Logistic Regression

Data: $D = \{(\vec{x}^{(i)}, y^{(i)})\}$ $y^{(i)} \in \{1, \dots, k\}$

$$\begin{aligned} \text{Model: } p(y | \vec{x}; \vec{\theta}) &= \frac{\exp(\vec{\theta}_y^\top \vec{x})}{Z(\vec{x})} \\ &= \frac{\exp(\vec{\theta}_y^\top \vec{x})}{\sum_{y'=1}^k \exp(\vec{\theta}_{y'}^\top \vec{x})} \end{aligned}$$

Objective: $J(\vec{\theta}) = \sum_{i=1}^N J^{(i)}(\vec{\theta})$ where $J^{(i)}(\vec{\theta}) = -\log p(y^{(i)} | \vec{x}^{(i)})$

Gradient:

$$\nabla J^{(i)}(\vec{\theta}) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_m} \end{bmatrix}$$

$$\frac{\partial J^{(i)}(\vec{\theta})}{\partial \theta_m} = \frac{1}{\sum_{y'=1}^k \exp(\vec{\theta}_{y'}^\top \vec{x}^{(i)})} \left(-\log p(y^{(i)} | \vec{x}^{(i)}) \right)$$

$$\frac{\partial}{\partial \theta_m} \left[-\log p(y | \vec{x}; \vec{\theta}) \right] = ?$$

$$K=3$$

$$\vec{\theta} = \begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vec{\theta}_3 \end{bmatrix}$$

$$\vec{\theta} = \begin{bmatrix} \vec{\theta}_1 \\ \vdots \\ \vec{\theta}_m \end{bmatrix}$$

$$\begin{aligned}
\frac{d}{d\theta_{2m}} \left[-\log p(y|\vec{x}; \vec{\theta}) \right] &= \frac{d}{d\theta_{2m}} \left[-\log \left(\frac{\exp(\theta_y^T \vec{x})}{Z(\vec{x}; \vec{\theta})} \right) \right] \\
&= \frac{d}{d\theta_{2m}} \left(-[\theta_y^T \vec{x} - \log Z(\vec{x}; \vec{\theta})] \right) \\
&= \frac{d}{d\theta_{2m}} \left(-\theta_y^T \vec{x} + \log \sum_{y'} \exp(\theta_{y'}^T \vec{x}) \right) \\
&= \frac{d}{d\theta_{2m}} \left(-\sum_{n=1}^M \theta_{y_n} x_n + \log Z(\vec{x}) \right) \\
&\stackrel{?}{=} -\mathbb{I}(y=z) x_m + \frac{d}{d\theta_{2m}} (\log Z(\vec{x})) \\
\frac{d}{d\theta_{2m}} \log Z(\vec{x}; \vec{\theta}) &= \frac{1}{Z(\vec{x}; \vec{\theta})} \cdot \frac{d}{d\theta_{2m}} \sum_{y'} \exp(\theta_{y'}^T \vec{x}) \\
&= \frac{1}{Z(\vec{x}; \vec{\theta})} \cdot \frac{d}{d\theta_{2m}} \exp(\theta_z^T \vec{x}) \\
&= \frac{\exp(\theta_z^T \vec{x})}{Z(\vec{x}; \vec{\theta})} \cdot \frac{d}{d\theta_{2m}} (\theta_z^T \vec{x}) \\
&= \frac{\exp(\theta_z^T \vec{x})}{Z(\vec{x}; \vec{\theta})} \cdot x_m \\
&= p(z|\vec{x}; \vec{\theta}) x_m \\
&= -\mathbb{I}(y=z) x_m + p(z|\vec{x}; \vec{\theta}) x_m
\end{aligned}$$

L2 Regularization

$$\begin{aligned}
J(\vec{\theta}) &= \left[\frac{1}{N} \sum_{i=1}^N J^{(i)}(\vec{\theta}) \right] + \lambda \sum_{m=1}^M \theta_m^2 \\
&= \left[\frac{1}{N} \sum_{i=1}^N -\log p(y^{(i)}|x^{(i)}; \vec{\theta}) \right] + \lambda \|\vec{\theta}\|_2^2
\end{aligned}$$