

Day 3, 9/14/18

Problem: adaptively select learning rate  $\alpha$

Intuition: large  $\alpha$  is good, especially in early stages of optimization, as we advance faster towards the optimum.

But we can overshoot. If we overshoot, we backtrack (i.e., start over, with a reduced  $\alpha$ )

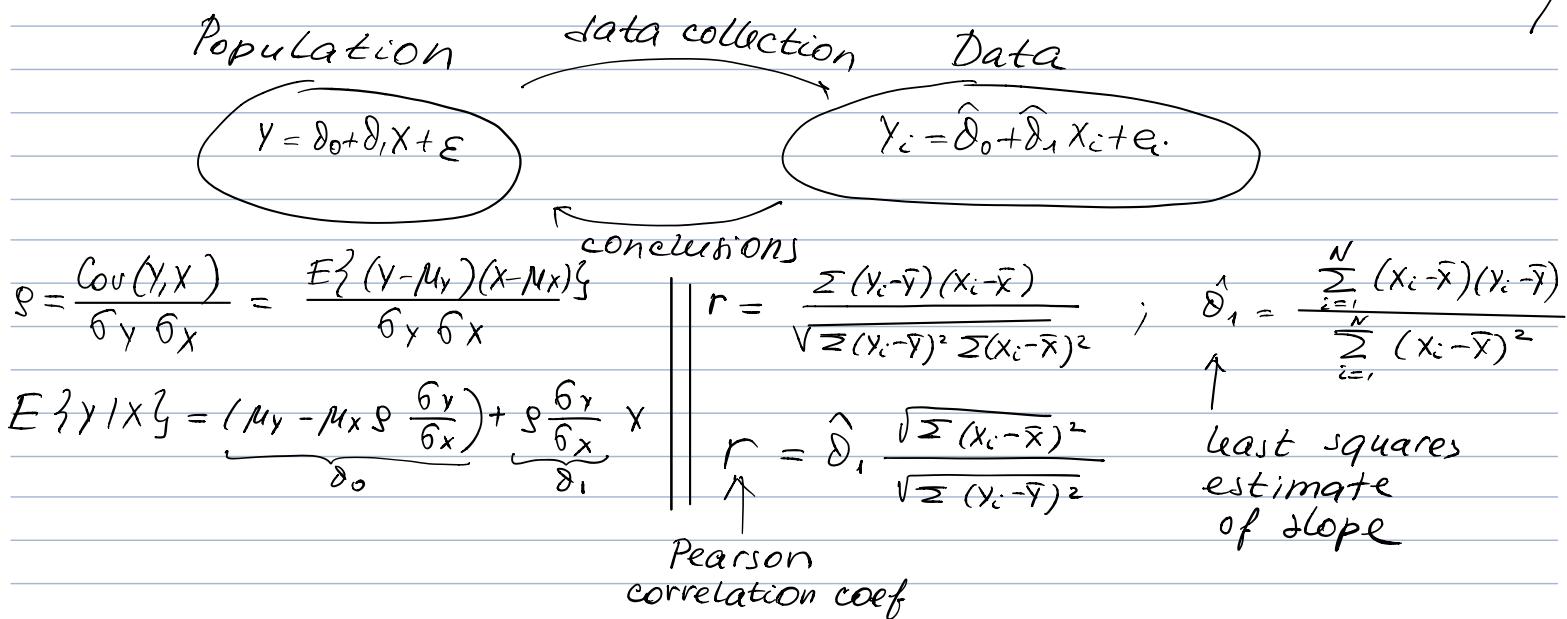
Algo: line search (replaces the highlighted part of batch/stochastic gradient descent)

- Input:  $\delta^{(t)}$ ,  $J(\delta)$ ,  $\nabla J(\delta)$
- Parameters:  $\alpha_{\max}$ ,  $\tau \in [0.5, 0.9]$ , tolerance, maxBacktrack
- $\alpha \leftarrow \alpha_{\max}$
- $\delta \leftarrow \delta^{(t)}$
- $obj \leftarrow J(\delta)$
- Repeat until maxBacktrack {
  - For  $j = 1, \dots, P$  {
    - $\delta_j \leftarrow \delta_j + \alpha \frac{1}{N} \sum_{i=1}^N (y_i - x_i^\top \delta^{(t)})^2$  [batch version]
    - if  $J(\delta_j) < obj - \text{tolerance}$ , then break [improved  $J(\delta)$ ]
    - else  $\alpha \leftarrow \tau \alpha$  [backtrack]
- If maxBacktrack reached, return  $\{\delta^{(t)}, \alpha=0\}$
- else return  $\{\delta, \alpha\}$

## Practical considerations in linear regression

Empirical measures of association when  $(\begin{matrix} X \\ Y \end{matrix}) \sim N \left( \begin{matrix} \mu_X \\ \mu_Y \end{matrix}, \begin{pmatrix} \sigma_x^2 & \text{Cov}(X,Y) \\ \text{Cov}(Y,X) & \sigma_y^2 \end{pmatrix} \right)$

### Pearson correlation



Conclusion: empirical summaries mirror population summaries. However,

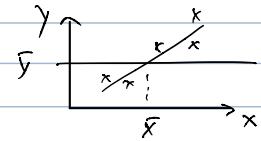
correlation does not extend to multiple predictors

## Coefficient of multiple determination $R^2$

$$\sum_{i=1}^N (y_i - \bar{y})^2 = \sum_{i=1}^N (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

$$\sum_{i=1}^N (y_i - \bar{y})^2 = \underbrace{\sum_{i=1}^N (y_i - \hat{y}_i)^2}_{\text{SSTO}} + \underbrace{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}_{\text{SSE}}$$

[the cross-product = 0]



$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO}, \text{ i.e. the proportion of total variation in } y \text{ that is explained by its linear association with } X$$

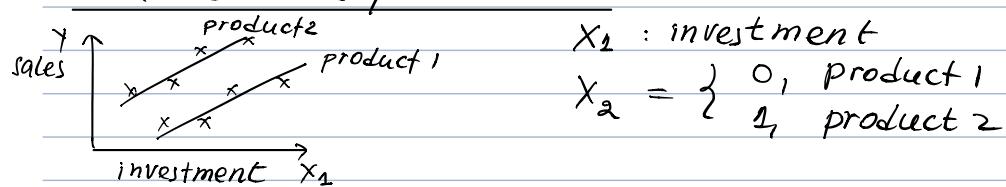
Plugging in the expression for  $\hat{y}_i$ :

$$SSR = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^N (\hat{\delta}_0 + \hat{\delta}_1 x_i - \frac{1}{N} \sum_{j=1}^N \hat{\delta}_0 - \frac{1}{N} \sum_{j=1}^N \hat{\delta}_1 x_j)^2 = \hat{\delta}_1 \sum_{i=1}^N (x_i - \bar{x})^2$$

$$R^2 = \frac{\hat{\delta}_1^2 \sum_{i=1}^N (x_i - \bar{x})^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = r^2 \quad \text{i.e. in linear regression with one predictor, the coefficient of multiple determination is the square of Pearson correlation coefficient}$$

$R^2$  extends to multiple predictors;  $r$  does not. Neither quantify the linearity of the association

## Qualitative predictors



$$y_i = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \epsilon_i$$

$$\begin{aligned} \text{Product 1: } y_i &= \delta_0 + \delta_1 x_1 + \epsilon_i \\ \text{Product 2: } y_i &= (\delta_0 + \delta_2) + \delta_1 x_1 + \epsilon_i \end{aligned} \quad \left. \begin{array}{l} \text{Different intercepts} \\ \text{same slope} \\ \text{same variance} \end{array} \right\}$$

$$y_i = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \delta_3 x_1 x_2 + \epsilon_i$$

$$\begin{aligned} \text{Product 1: } y_i &= \delta_0 + \delta_1 x_1 + \epsilon_i \\ \text{Product 2: } y_i &= (\delta_0 + \delta_2) + (\delta_1 + \delta_3) x_1 + \epsilon_i \end{aligned} \quad \left. \begin{array}{l} \text{Different intercepts} \\ \text{different slopes} \\ \text{same variance} \end{array} \right\}$$

$$3 \text{ products: } x_2 = \begin{cases} 0, \text{ product 1} \\ 1, \text{ product 2} \\ 0, \text{ product 3} \end{cases} \quad x_3 = \begin{cases} 0, \text{ product 2} \\ 0, \text{ product 3} \\ 1, \text{ product 3} \end{cases}$$

Note: coding  $x_2 = \begin{cases} 0, \text{ product 1} \\ 1, \text{ product 2} \\ 2, \text{ product 3} \end{cases}$  is inappropriate, because it assumes that the difference in intercepts between products 3 and 1 is twice as big as between products 2 and 1.

Simpson's paradox: change of association between two variables when condition on another variable

