

Homework 2

Yuanteng Chen

1. Why Learning Rates Cannot be Too Big.

(a) For what values of learning rate $\eta > 0$ is the recurrence (3) stable?

Solution:

$$W_{t+1} = (1 - 2\eta\sigma^2)W_t + 2\eta\sigma y$$

$$2\eta\sigma y = 2\eta\sigma^2 \cdot \frac{y}{\sigma}$$

$$\therefore W_{t+1} = (1 - 2\eta\sigma^2) \left(W_t - \frac{y}{\sigma} \right) + \frac{y}{\sigma}$$

$$W_1 = (1 - 2\eta\sigma^2) \left(W_0 - \frac{y}{\sigma} \right) + \frac{y}{\sigma}$$

$$W_2 = (1 - 2\eta\sigma^2) \left(W_1 - \frac{y}{\sigma} \right) + \frac{y}{\sigma}$$

$$= (1 - 2\eta\sigma^2) \left[(1 - 2\eta\sigma^2) \left(W_0 - \frac{y}{\sigma} \right) + \frac{y}{\sigma} - \frac{y}{\sigma} \right] + \frac{y}{\sigma}$$

$$= (1 - 2\eta\sigma^2)^2 \left(W_0 - \frac{y}{\sigma} \right) + \frac{y}{\sigma}$$

\vdots

$$W_{t+1} = (1 - 2\eta\sigma^2)^{t+1} \left(W_0 - \frac{y}{\sigma} \right) + \frac{y}{\sigma}$$

$$|1 - 2\eta\sigma^2| < 1 \Leftrightarrow \text{recurrence is stable}$$

$$\therefore 0 < \eta < \frac{1}{\sigma^2}$$

(b). get within a factor $(1 - \varepsilon)$ of w^*

||

$$|W_t - w^*| < \varepsilon |w^*|$$

$$\therefore W_{t+1} = (1 - 2\eta\sigma^2)^{t+1} (W_0 - \frac{y}{\sigma}) + \frac{y}{\sigma}$$

$$W^* = \frac{y}{\sigma}$$

$$|W_t - \frac{y}{\sigma}| < \epsilon \left| \frac{y}{\sigma} \right|$$

$$W_t - \frac{y}{\sigma} = (1 - 2\eta\sigma^2)^t (W_0 - \frac{y}{\sigma})$$

$$\therefore \text{initial condition } W_0 = 0$$

$$\therefore |(1 - 2\eta\sigma^2)^t (-\frac{y}{\sigma})| < \epsilon \left| \frac{y}{\sigma} \right|$$

$$\therefore (1 - 2\eta\sigma^2)^t < \epsilon$$

$$\log (1 - 2\eta\sigma^2)^t < \log \epsilon$$

$$\therefore \log (1 - 2\eta\sigma^2) < 0$$

$$\therefore t > \frac{\log \epsilon}{\log (1 - 2\eta\sigma^2)}$$

$$(c) \begin{bmatrix} \sigma_L & 0 \\ 0 & \sigma_S \end{bmatrix} \begin{bmatrix} W[1] \\ W[2] \end{bmatrix} = \begin{bmatrix} y[1] \\ y[2] \end{bmatrix}$$

$$\sigma_L > \sigma_S$$

initial condition $w=0$, a single learning rate η

$$\text{Solution: assume } \Sigma = \begin{bmatrix} \sigma_L^2 & 0 \\ 0 & \sigma_S^2 \end{bmatrix}$$

same as (1):

$$W_{t+1} = (E - 2\eta\Sigma^2)W_t + 2\eta\Sigma y$$

$$\Sigma^2 = \begin{bmatrix} \sigma_L^2 & 0 \\ 0 & \sigma_S^2 \end{bmatrix} \quad E - 2\eta\Sigma^2 = \begin{bmatrix} 1 - 2\eta\sigma_L^2 & 0 \\ 0 & 1 - 2\eta\sigma_S^2 \end{bmatrix}$$

$$\therefore \begin{cases} |1-2\eta\sigma_1^2| < 1 & \therefore \sigma_1 \gg \sigma_5 \\ |1-2\eta\sigma_5^2| < 1 & \therefore \eta < \frac{1}{\sigma_1^2} < \frac{1}{\sigma_5^2} \end{cases}$$

(d) depending on η, σ_1, σ_5 , which of the two dimensions is converging faster and which one is converging slower?

Solution: in c) $W_{t+1} = (1-2\eta\sigma^2)^{t+1} (W_0 - \frac{y}{\sigma}) + \frac{y}{\sigma}$

$$\therefore \begin{cases} W_{[1]t+1} = (1-2\eta\sigma_1^2)^{t+1} (W_0 - \frac{y_{[1]}}{\sigma_1}) + \frac{y_{[1]}}{\sigma_1} \\ W_{[2]t+1} = (1-2\eta\sigma_5^2)^{t+1} (W_0 - \frac{y_{[2]}}{\sigma_5}) + \frac{y_{[2]}}{\sigma_5} \end{cases}$$

$$\text{if } |1-2\eta\sigma_1^2| < |1-2\eta\sigma_5^2| < 1$$

then $W_{[1]}$ is converging faster,

else $W_{[2]}$ is converging faster.

(e) when $|1-2\eta\sigma_1^2| = |1-2\eta\sigma_5^2|$, we get the fastest overall convergence to the solution.

$$\therefore \sigma_1 \gg \sigma_5 \therefore 2 = 2\eta(\sigma_1^2 + \sigma_5^2)$$

$$\eta = \frac{1}{\sigma_1^2 + \sigma_5^2}$$

(f): the speed of converging of σ_i ($i \neq 1, 5$)

is between ϵ_1 and ϵ_5 ,

so they will not influence the choice of possible learning rates.

(9) I have no idea

2. Accelerating Gradient Descent with Momentum.

$$L(w) = \|y - Xw\|_2^2$$

$$w_{t+1} = w_t - \eta z_{t+1}$$

$$z_{t+1} = (1-\beta)z_t + \beta g_t$$

the gradient descent update:

$$w_{t+1} = (I - 2\eta(X^T X))w_t + 2\eta X^T y$$

$$w^* = (X^T X)^{-1} X^T y$$

(18) $w_{t+1} = w_t - \eta z_{t+1}$

$$z_{t+1} = (1-\beta)z_t + \beta(2X^T X w_t - 2X^T y)$$

$$x_t = V^T(w_t - w^*)$$

$$a_t = V^T z_t$$

(a) ① $w_{t+1} = w_t - \eta z_{t+1}$

$$V^T w_{t+1} = V^T w_t - \eta V^T z_{t+1}$$

$$V^T w_{t+1}[i] = V^T w_t[i] - \eta V^T z_{t+1}[i]$$

$$\therefore X_t = V^T (W_t - W^*)$$

$$\therefore V^T W_{t+1} = X_{t+1} + V^T \cdot W^*$$

$$V^T W_t = X_t + V^T \cdot W^*$$

$$\therefore X_{t+1} + V^T \cdot W^* = X_t + V^T \cdot W^* - \eta V^T Z_{t+1}$$

$$X_{t+1} = X_t - \eta V^T Z_{t+1}$$

$$\therefore a_t = V^T Z_t$$

$$\therefore X_{t+1} = X_t - \eta a_{t+1}$$

$$X_{t+1}[i] = X_t[i] - \eta a_{t+1}[i]$$

$$\textcircled{2} Z_{t+1} = (1-\beta)Z_t + \beta(2X_t^T X W_t - 2X_t^T y)$$

$$V^T Z_{t+1} = (1-\beta)V^T Z_t + \beta V^T (2X_t^T X W_t - 2X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta V^T (2V \Sigma U^T U \Sigma V^T W_t - 2X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta V^T (2V \Sigma^2 V^T W_t - 2X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta (2V^T V \Sigma^2 V^T W_t - 2V^T X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta (2\Sigma^2 V^T W_t - 2V^T X_t^T y)$$

$$\therefore X_t = V^T (W_t - W^*)$$

$$V X_t = W_t - W^* \therefore W_t = V X_t + W^*$$

$$\therefore a_{t+1} = (1-\beta)a_t + \beta (2\Sigma^2 V^T (V X_t + W^*) - 2V^T X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta (2\Sigma^2 X_t + 2\Sigma^2 V^T W^* - 2V^T X_t^T y)$$

$$\begin{aligned}\therefore W^* &= (X^T X)^{-1} X^T y \\ &= (V \Sigma^2 V^T)^{-1} X^T y \\ &= V \Sigma^{-2} V^T X^T y\end{aligned}$$

plug W^* into :

$$a_{t+1} = (1-\beta)a_t + \beta(2\Sigma^2 X_t + 2\Sigma^{-2} V^T V \Sigma^{-2} V^T X_t^T y - 2V^T X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + \beta(2\Sigma^2 X_t + 2V^T X_t^T y - 2V^T X_t^T y)$$

$$a_{t+1} = (1-\beta)a_t + 2\beta \Sigma^2 X_t$$

\Downarrow

$$a_{t+1}[i] = (1-\beta)a_t[i] + 2\beta \sigma_i^2 X_t[i]$$

$$\therefore \begin{cases} W_{t+1} = W_t - \eta Z_{t+1} \\ Z_{t+1} = (1-\beta)Z_t + \beta(2X_t^T X W_t - 2X_t^T y) \end{cases}$$

\Downarrow

$$\begin{cases} a_{t+1}[i] = (1-\beta)a_t[i] + 2\beta \sigma_i^2 X_t[i] \\ X_{t+1}[i] = X_t[i] - \eta a_{t+1}[i] \end{cases}$$

cb) $\begin{bmatrix} a_{t+1}[i] \\ X_{t+1}[i] \end{bmatrix} = R_i \begin{bmatrix} a_t[i] \\ X_t[i] \end{bmatrix}$ Derive R_i

$$\therefore \begin{cases} a_{t+1}[i] = (1-\beta)a_t[i] + 2\beta \sigma_i^2 X_t[i] \\ X_{t+1}[i] = X_t[i] - \eta a_{t+1}[i] \end{cases}$$

$$X_{t+1}[i] = X_t[i] - \eta a_{t+1}[i]$$

$$= X_t[i] - \eta (1-\beta) a_t[i] - 2\eta \beta \sigma_i^2 X_t[i]$$

$$= (1-2\eta \beta \sigma_i^2) X_t[i] - \eta (1-\beta) a_t[i]$$

$$\therefore \begin{cases} a_{t+1}[i] = (1-\beta) a_t[i] + 2\beta \sigma_i^2 X_t[i] \\ X_{t+1}[i] = (1-2\eta \beta \sigma_i^2) X_t[i] - \eta (1-\beta) a_t[i] \end{cases}$$

$$\therefore R_i = \begin{bmatrix} 1-\beta & , & 2\beta \sigma_i^2 \\ \eta(1-\beta) & , & 1-2\eta \beta \sigma_i^2 \end{bmatrix}$$

$$c) \quad R_i \vec{x} = \lambda \vec{x}$$

$$(R_i - \lambda E) \vec{x} = 0$$

$$|R_i - \lambda E| = 0$$

$$\begin{vmatrix} 1-\beta-\lambda & 2\beta \sigma_i^2 \\ \eta(1-\beta) & 1-2\eta \beta \sigma_i^2-\lambda \end{vmatrix} = 0$$

$$1-\beta-\lambda-2\eta \beta \sigma_i^2 + 2\eta \beta^2 \sigma_i^2 + 2\eta \beta \sigma_i^2 \lambda - \lambda + \beta \lambda + \lambda^2 \\ + 2\eta \sigma_i^2 \beta (1-\beta) = 0$$

$$\lambda^2 - (2-\beta-2\eta \beta \sigma_i^2) \lambda + (1-\beta) = 0$$

$$\Delta = (2-\beta-2\eta \beta \sigma_i^2)^2 - 4(1-\beta)$$

$$\begin{cases} \Delta \geq 0 \Leftrightarrow \text{real eigenvalues} \\ \Delta < 0 \Leftrightarrow \text{complex eigenvalues.} \end{cases}$$

(d) when λ is repeated: $\Delta = 0$

$$(2 - \beta - 2\eta\beta\sigma i^2)^2 = 4(1 - \beta)$$

$$2 - \beta - 2\eta\beta\sigma i^2 = \pm 2\sqrt{1 - \beta}$$

$$\eta = \frac{2 - \beta \pm 2\sqrt{1 - \beta}}{2\beta\sigma i^2}$$

$$\therefore \text{highest } \eta = \frac{2 - \beta + 2\sqrt{1 - \beta}}{2\beta\sigma i^2}$$

(e) when λ is real

$$\eta > \frac{2 - \beta + 2\sqrt{1 - \beta}}{2\beta\sigma i^2} \quad \text{or} \quad \eta < \frac{2 - \beta - 2\sqrt{1 - \beta}}{2\beta\sigma i^2}$$

(f) when λ is complex:

$$\frac{2 - \beta + 2\sqrt{1 - \beta}}{2\beta\sigma i^2} < \eta < \frac{2 - \beta - 2\sqrt{1 - \beta}}{2\beta\sigma i^2}$$

$$(g) \text{ optimal rate} = \frac{(\sigma_{\max}/\sigma_{\min})^2 - 1}{(\sigma_{\max}/\sigma_{\min})^2 + 1}$$

$$\sigma_{\max}^2 = 5 \quad \sigma_{\min}^2 = 0.05$$

$$\therefore \text{rate} = \frac{100-1}{100+1} = \frac{99}{101}$$

using ordinary gradient descent:

$$\left(\frac{99}{101}\right)^{T_1} \leq 99.5\%$$

using this learning rate with momentum:

in (c):

$$\text{we got } \lambda_1, \lambda_2 = \sqrt{1-\beta} = \sqrt{0.9} < 1$$

\therefore the higher one of λ_1, λ_2 is $\geq \sqrt{0.9}$

\therefore the convergence rate $r \geq \sqrt{0.9}$

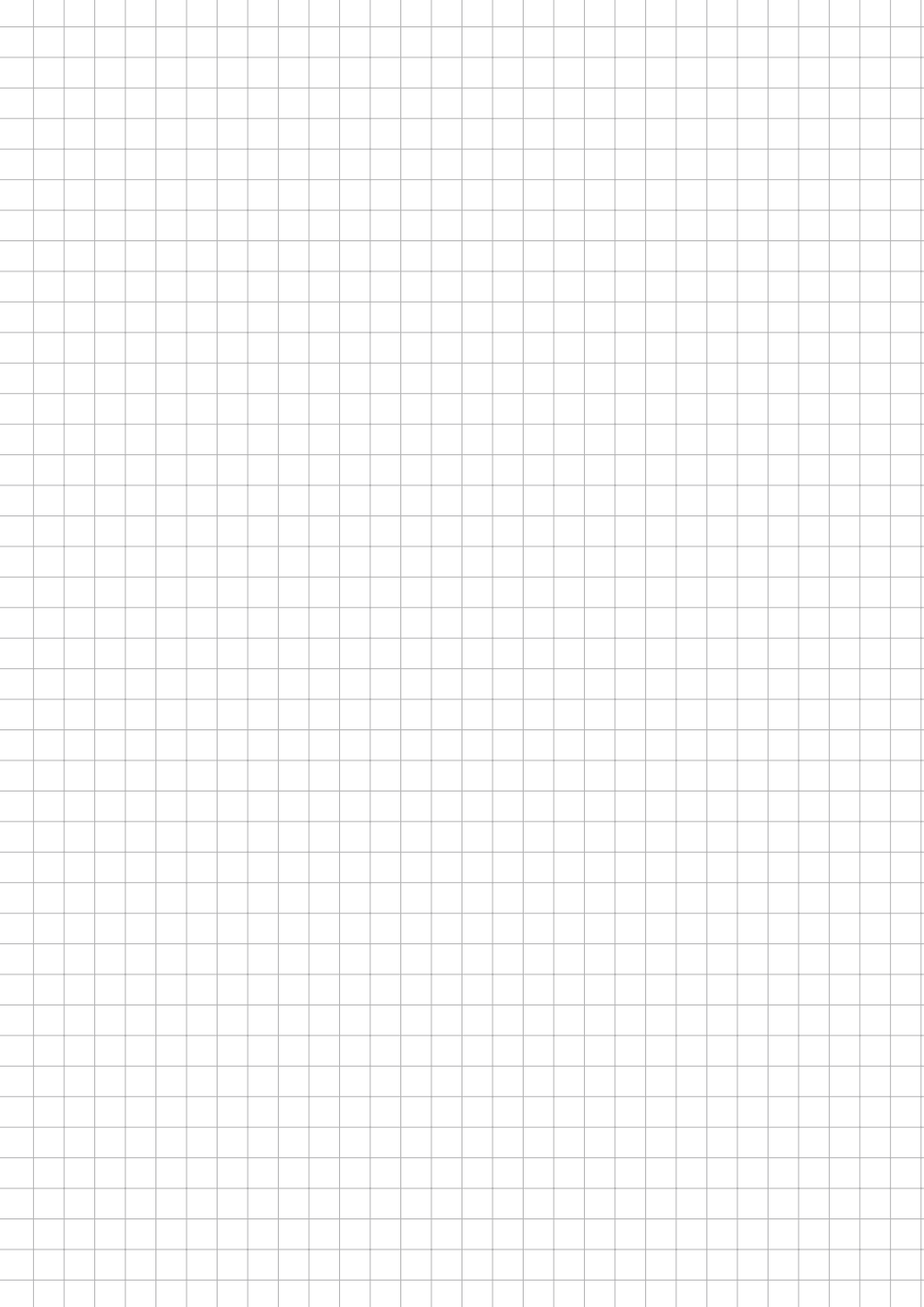
$$(c) \quad T_2 \leq 99.5\%$$

$$\log(r)^{T_2} \leq \log 99.5\%$$

$$T_2 \log r \leq \log 99.5\%$$

$$T_2 \geq \frac{\log 99.5\%}{\log r}$$

$$T_1 > T_2 \quad \left(\sqrt{0.9} < \frac{99}{101} \right)$$



3- Regularization and Instance Noise

$$\tilde{X}_i = X_i + N_i \quad N_i \sim \mathcal{N}(0, \sigma^2 I_n)$$

$$\tilde{X} = \begin{bmatrix} \tilde{X}_1^T \\ \tilde{X}_2^T \\ \vdots \\ \tilde{X}_m^T \end{bmatrix} \quad \bar{X}_i \in \mathbb{R}^n \text{ and } y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

$$\arg \min_w E[||\tilde{X}w - y||^2]$$

$$(a) E[||\tilde{X}w - y||^2]$$

$$= E\left[\sum_{i=1}^m (\tilde{X}_i^T w - y_i)^2\right]$$

$$= \sum_{i=1}^m E[(X_i + N_i)^T w - y_i]^2$$

$$= \sum_{i=1}^m E[(X_i^T w + N_i^T w - y_i)^2]$$

$$= \sum_{i=1}^m E[(X_i^T w - y_i + N_i^T w)^2]$$

$$= \sum_{i=1}^m E[(X_i^T w - y_i)^2 + 2(N_i^T w)(X_i^T w - y_i) + (N_i^T w)^2]$$

$$= \sum_{i=1}^m E[(X_i^T w - y_i)^2] - 2E[(N_i^T w)(X_i^T w - y_i)] + E[(N_i^T w)^2]$$

$$= \sum_{i=1}^m [(X_i^T w - y_i)^2 - 2E[(N_i^T w)(X_i^T w - y_i)] + E[(N_i^T w)^2]]$$

$$\because N_i \in \mathcal{N}(0, \sigma^2 I_n)$$

$$\therefore 2E[(N_i^T w)(X_i^T w - y_i)] = 0$$

$$E[N_i N_i^T] = \sigma^2 I_n$$

$$\therefore = \sum_{i=1}^m (X_i^T w - y_i)^2 + w^T \sigma^2 I_n w$$

$$= (Xw - y)^2 + \sigma^2 \|w\|^2 \cdot m$$

⇓

is equivalent to a regularized least squares problem:

$$\arg \min_w \frac{1}{m} \|Xw - y\|^2 + \lambda \|w\|^2$$

(b) $\tilde{x}_i = x + Nt, Nt \in (0, \sigma^2)$

$$L(w) = \frac{1}{2} (\tilde{x}w - y)^2 \quad w_0 = 0$$

$$\frac{\partial L}{\partial w} = (\tilde{x}w - y) \tilde{x}$$

$$= (w(x + Nt) - y)(x + Nt)$$

$$= w(x + Nt)^2 - y(x + Nt)$$

$$= w(x^2 + 2xNt + Nt^2) - y(x + Nt)$$

$$\therefore w_{t+1} = w_t - \eta \cdot \frac{\partial L}{\partial w}$$

$$= w_t - \eta [w_t(x^2 + 2xNt + Nt^2) - y(x + Nt)]$$

$$= w_t (1 - \eta(x^2 + 2xNt + Nt^2)) - \eta y(x + Nt)$$

as Nt is i.i.d. $E(Nt) = 0$ $E(Nt^2) = \sigma^2$

$$E(w_{t+1}) = E(w_t) \cdot E(1 - \eta(x^2 + 2xNt + Nt^2)) - E(\eta y(x + Nt))$$

$$E(w_{t+1}) = E(w_t) \cdot E(1 - \eta(x^2 + \sigma^2)) - E(\eta y x)$$

$$= E(w_t) (1 - \eta(x^2 + \sigma^2)) - \eta y x$$

(c) . For what values of learning rate η do we expect

the expectation of the learned weight to converge using gradient descent?

Solution:

gradient descent to converge

$$\Leftrightarrow -1 < 1 - \eta(x^2 + \sigma^2) < 1$$

$$0 < \eta < \frac{2}{x^2 + \sigma^2}$$

(d) what would we expect $E(w_t)$ to converge as $t \rightarrow \infty$? How does this differ from the situation without noise?

Solution:

when $E(w_t)$ to converge:

$$\frac{\partial L}{\partial w} = w(x^2 + 2xN_t + N_t^2) - y(x - N_t)$$

$$E\left(\frac{\partial L}{\partial w}\right) = w(x^2 + \sigma^2) - yx = 0$$

$$w = \frac{yx}{x^2 + \sigma^2}$$

$$w = \frac{y}{x} \frac{1}{1 + \frac{\sigma^2}{x^2}}$$

without noise: $w = \frac{y}{x}$

there is a scalar value $\frac{1}{1 + \frac{\sigma^2}{x^2}}$

when noise is added to x .

7. (a) CSDN, ChatGPT

(b)

(c) 16 hours.