

4/20/20

Regression Quality Metrics

trained model : $\text{model}(\tilde{x}, \tilde{w}^*) = \tilde{x}^T \tilde{w}^*$

$\uparrow\uparrow$
LS LAD

Prediction :

$\tilde{x}^o : \text{model}(\tilde{x}^o, \tilde{w}^*) = y_o$

Judging the quality of trained model

$$\rightarrow \underline{\text{MSE}} = \frac{1}{P} \sum_{p=1}^P (\text{model}(\tilde{x}_p, \tilde{w}^*) - y_p)^2$$

\uparrow
LS, LAD

$$\rightarrow \text{MAD} = \frac{1}{P} \sum_{p=1}^P |\text{model}(\tilde{x}_p, \tilde{w}^*) - y_p|$$

MIN
absolute

deviations

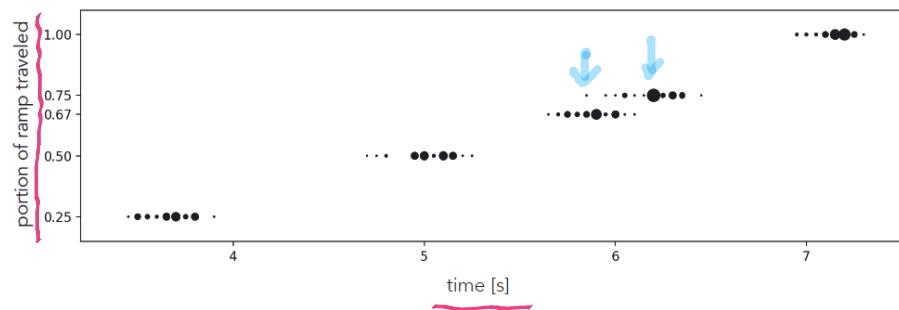


Figure 5.9 Figure associated with Example 5.7. See text for details.

Weighted Regression

emphasize or de-emphasize points

- deal w/ duplicates

due to quantization or binning

assume b_p is the number of identical input pair
 $(\tilde{x}_p, \tilde{y}_p)$

Error

$$\underbrace{(\text{model}(\tilde{x}_1, \tilde{w}) - y_1)^2 + \dots + (\text{model}(\tilde{x}_1, \tilde{w}) - y_1)^2}_{b_1}$$

$$+ \underbrace{(\text{model}(\tilde{x}_p, \tilde{w}) - y_p)^2 + \dots + (\text{model}(\tilde{x}_p, \tilde{w}) - y_p)^2}_{b_p}$$

$$g(\tilde{w}) = \frac{1}{b_1 + b_2 + \dots + b_p} \sum_{p=1}^P b_p (\text{model}(\tilde{x}_p, \tilde{w}) - y_p)^2$$

• Weight points by confidence

confidence in the trustworthiness of
each data point

Fig 5.10

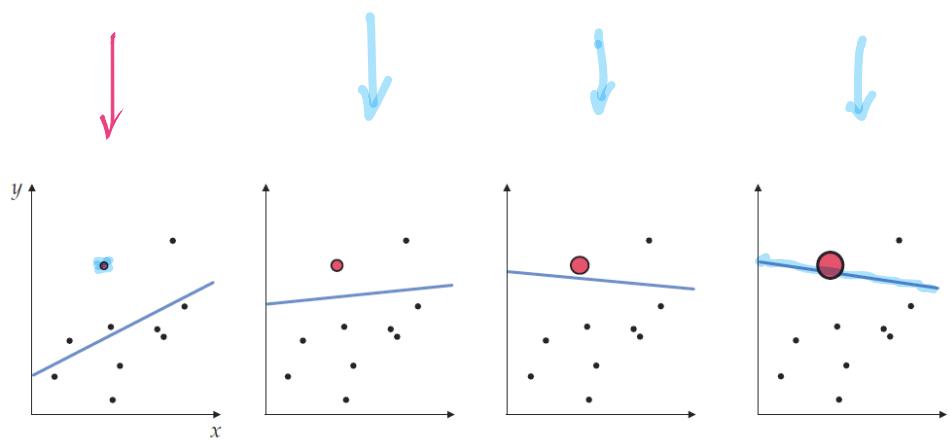


Figure 5.10 Figure associated with Example 5.8. See text for details.

Multi-output Regression

both input & output are vector-valued

P input pairs

$$(\bar{x}_1, \bar{y}_1), \dots, (\bar{x}_P, \bar{y}_P)$$

$$\bar{x}_P = \begin{bmatrix} x_{1,P} \\ x_{2,P} \\ \vdots \\ x_{N,P} \end{bmatrix}_{N \times 1}, \quad \bar{y}_P = \begin{bmatrix} y_{0,P} \\ y_{1,P} \\ \vdots \\ y_{C-1,P} \end{bmatrix}_{1 \times C}$$

$\bar{x}_i \in \mathbb{R}^N$
 $\bar{y}_i \in \mathbb{R}^C$

assume: a linear relationship holds between

\bar{x}_P & just the c-th dimension of \bar{y}_P

$$\rightarrow \boxed{\tilde{x}_P^T \tilde{w}_c \approx y_{c,P}}, \quad P=1, \dots, P$$

scalar

$$\tilde{w}_c = [b_c \quad w_{1,c} \quad \dots \quad w_{N,c}]^T$$

assume: above applies to all C entries of output

$$W = \begin{bmatrix} \frac{\tilde{w}_0}{b_0} & \frac{\tilde{w}_1}{b_1} & \frac{\tilde{w}_c}{b_c} & \frac{\tilde{w}_{c-1}}{b_{c-1}} \\ w_{1,0} & w_{1,1} & w_{1,c} & w_{1,c-1} \\ \vdots & \vdots & \vdots & \vdots \\ w_{N,0} & w_{N,1} & w_{N,c} & w_{N,c-1} \end{bmatrix}_{(N+1) \times C}$$

$$\tilde{x}_p^T W = \left[\tilde{x}_p^T \tilde{w}_0 \quad \tilde{x}_p^T \tilde{w}_1 \quad \dots \quad \tilde{x}_p^T \tilde{w}_{c-1} \right]$$

$$\boxed{\tilde{x}_p^T W \approx \bar{y}_p, \quad p=1, \dots, P}$$

$1 \times (N+1)$ $(N+1) \times C$ $1 \times c$

Error at the p -th point: $\underline{(\tilde{x}_p^T W - \bar{y}_p)}$ $1 \times c$

ℓ_2 : Least squares:

$$g(W) = \frac{1}{P} \sum_{p=1}^P \| \tilde{x}_p^T W - \bar{y}_p \|_2^2 = \frac{1}{P} \sum_{p=1}^P \sum_{c=0}^{C-1} (\tilde{x}_p^T \tilde{w}_c - y_{c,p})^2$$

ℓ_1 : LAD

$$g(W) = \frac{1}{P} \sum_{p=1}^P \| \tilde{x}_p^T W - \bar{y}_p \|_1 = \frac{1}{P} \sum_{p=1}^P \sum_{c=0}^{C-1} |\tilde{x}_p^T \tilde{w}_c - y_{c,p}|$$

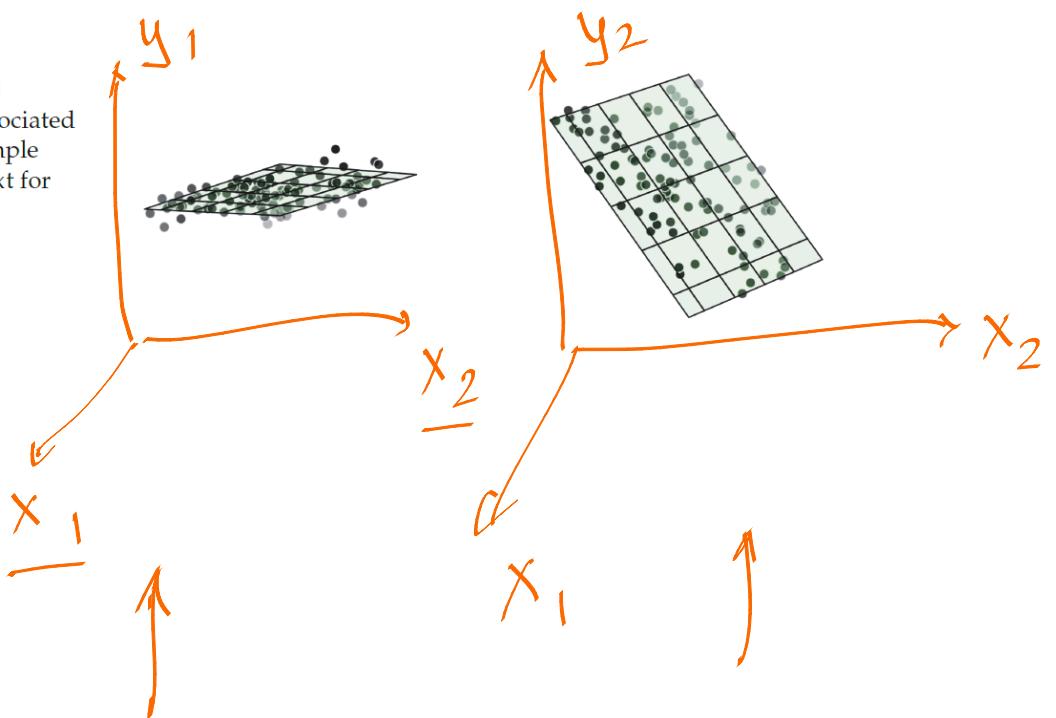
$$\underline{g(W)} = \sum_{c=0}^{C-1} \left(\frac{1}{P} \sum_{p=1}^P |\tilde{x}_p^T \tilde{w}_c - y_{c,p}| \right) = \sum_{c=0}^{C-1} g_c(\tilde{w}_c)$$

$g_c(\tilde{w}_c)$

Independent
optimization
for each \tilde{w}_c

$$N=2, C=2$$

Figure 5.11
Figure associated
with Example
5.9. See text for
details.



Vector Norms

$\bar{x} \in \mathbb{R}^N$

$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

Euclidean norm or ℓ_2 norm

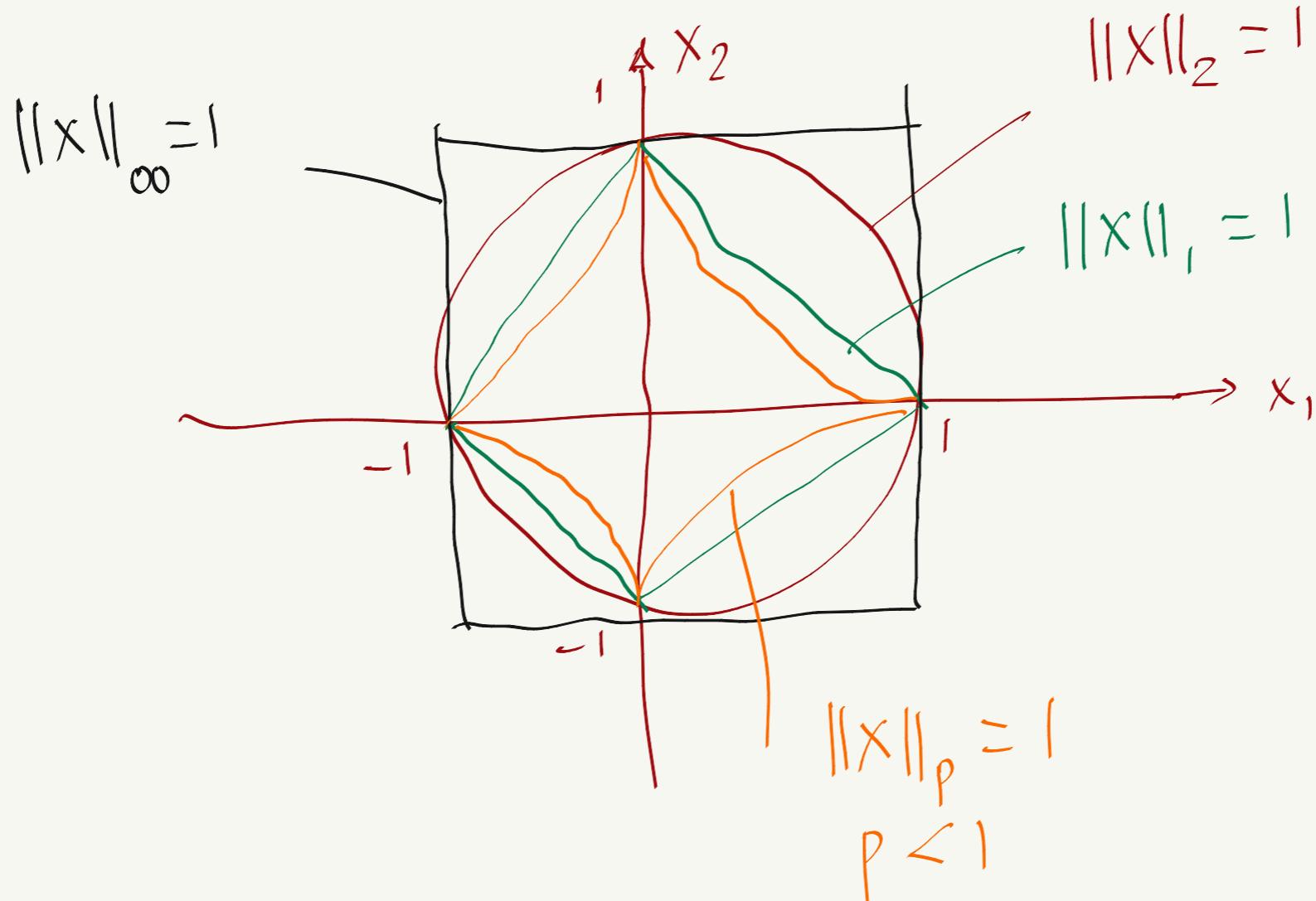
$$\ell_2: \|x\|_2 = \sqrt{\sum_{n=1}^N x_n^2}$$

$$\ell_1: \|x\|_1 = \sum_{n=1}^N |x_n|$$

$$\ell_\infty: \|x\|_\infty = \max_n |x_n|$$

$$\underline{\ell_p}: \|x\|_p = \left(\sum_{n=1}^N |x_n|^p \right)^{1/p}$$

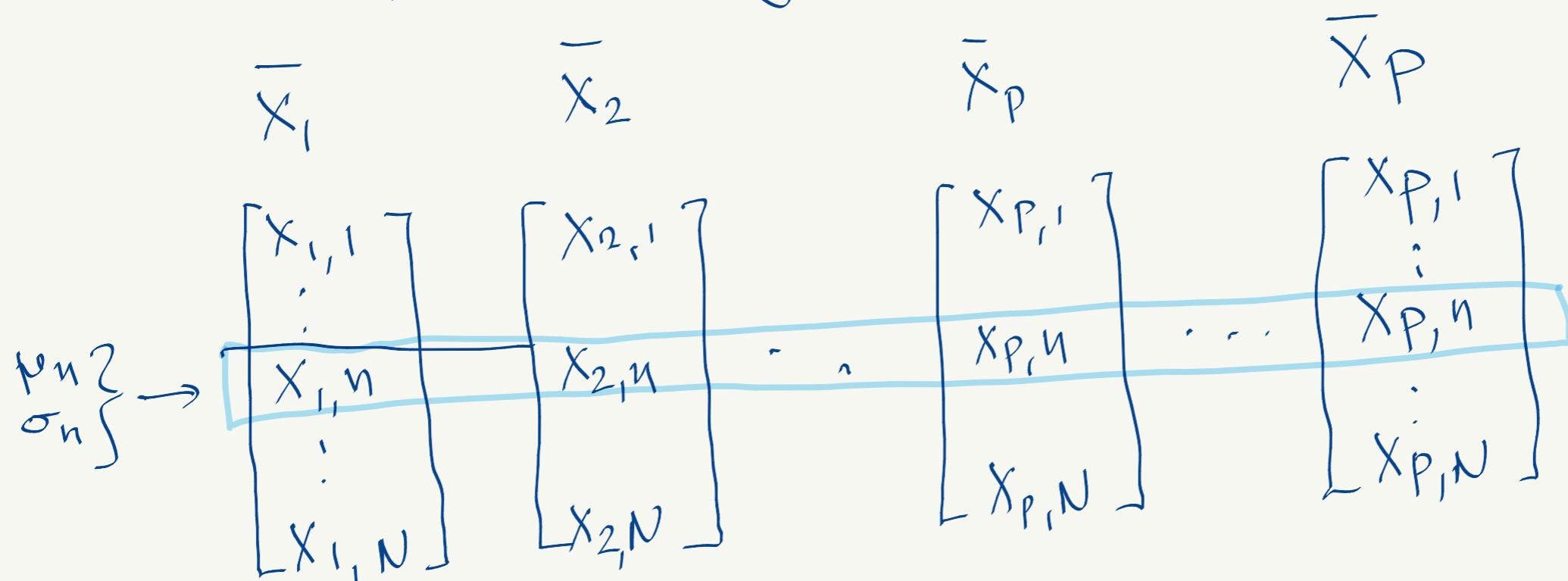
UNIT BALL



In problem 5.9 in the homework
 you are asked to do standard normalization
 of the input-features

Standard Normalization (9.3)

how? why?



n-th input feature

$$\begin{aligned}
 & \text{invertible} \\
 & x_{P,n} \leftarrow \frac{x_{P,n} - \mu_n}{\sigma_n}, \quad n=1, \dots, N \\
 & \mu_n = \frac{1}{P} \sum_{p=1}^P x_{p,n} \\
 & \sigma_n = \left(\frac{1}{P} \sum_{p=1}^P (x_{p,n} - \mu_n)^2 \right)^{1/2}
 \end{aligned}$$

except when $\sigma_n = 0$ (this means $x_{1,n} = \text{constant}$)

the n-th feature is redundant
 \Rightarrow is removed

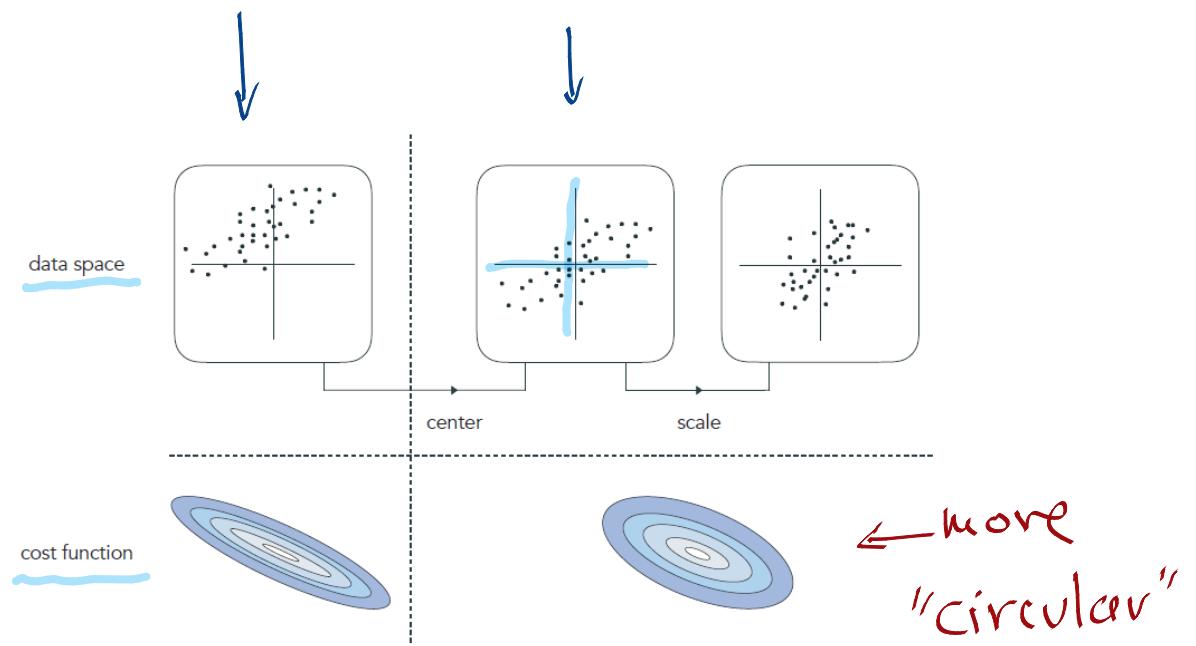


Figure 9.14 (Standard normalization illustrated. The input space of a generic dataset (top-left panel) as well as its mean-centered (top-middle panel) and scaled version (top-right panel). As illustrated in the bottom row, where a prototypical cost function corresponding to this data is shown, standard normalization results in a cost function with less elliptical and more circular contours compared to the original cost function.

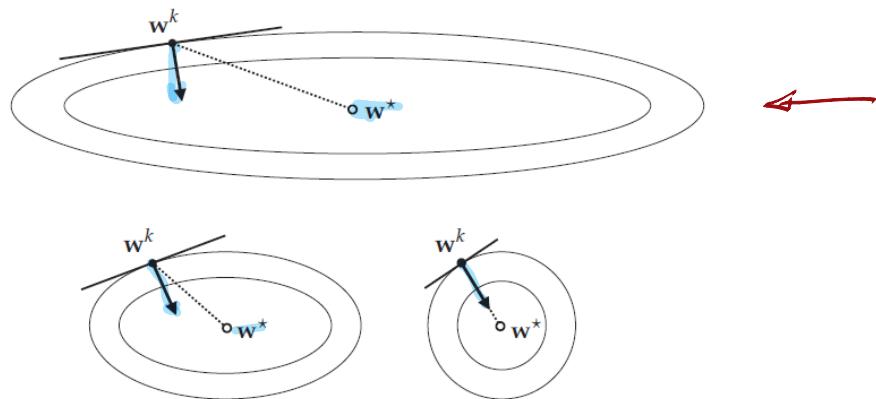
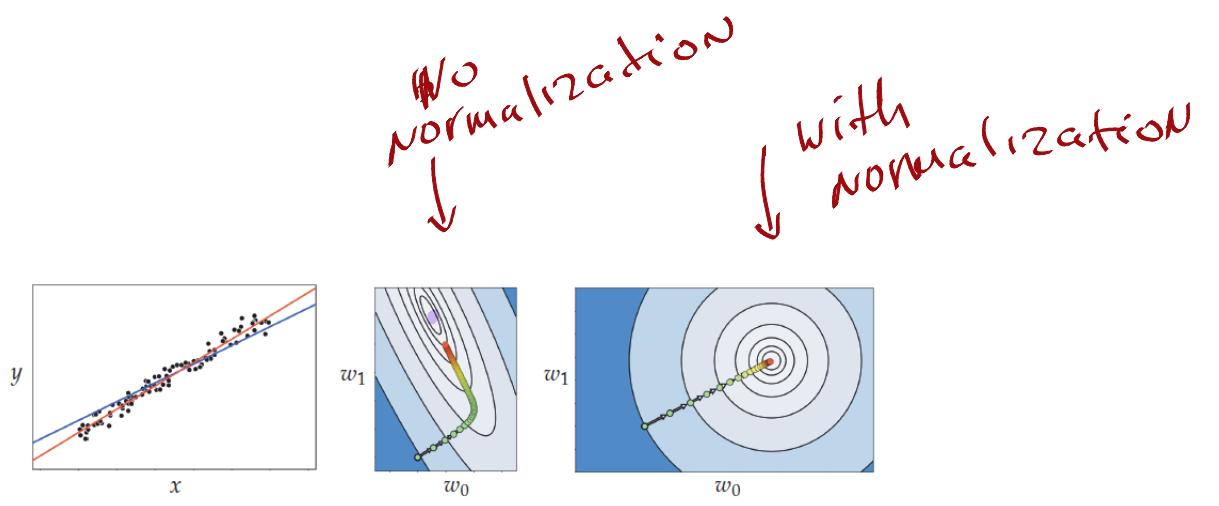


Figure 9.15 In standard normalizing input data we temper its associated cost function's often-elliptical contours, like those shown in the top panel, into more circular ones as shown in the bottom-left and bottom-right panels. This means that the gradient descent direction, which points away from the minimizer of a cost function when its contours are elliptical (leading to the common zig-zagging problem with gradient descent), points more towards the function's minimizer as its contours become more circular. This makes each gradient descent step much more effective, typically allowing the use of much larger steplength parameter values α , meaning that measurably fewer steps are required to adequately minimize the cost function.

Speed up convergence
(when GD is used)



$$\begin{bmatrix} \bar{w}^0 = [0 \ 0]^T \\ \alpha = 10^{-1} \end{bmatrix} \xrightarrow{\text{100 it}} \begin{array}{l} \text{blue line} \\ \downarrow \end{array} \quad \xrightarrow{\text{20 it}} \begin{array}{l} \text{red line} \\ \downarrow \end{array}$$