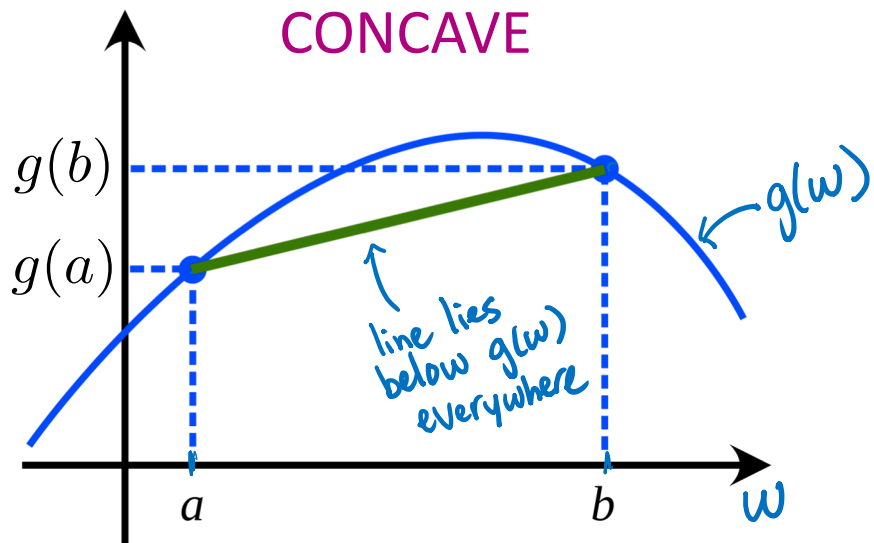


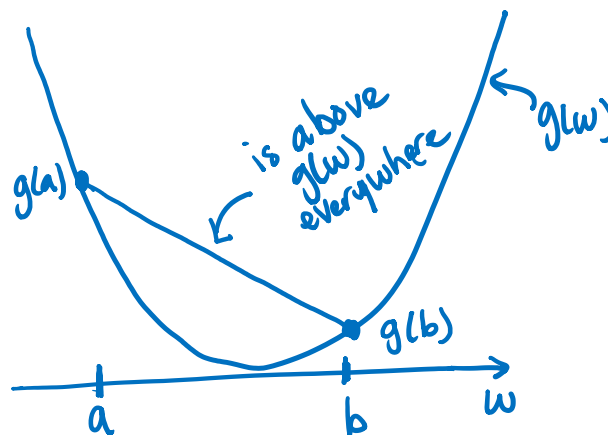
Gradient descent

Convex/concave functions

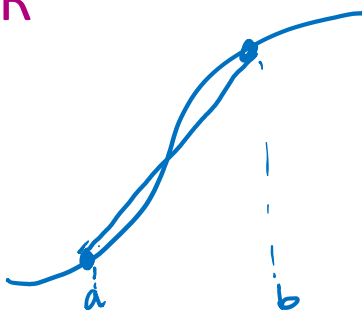
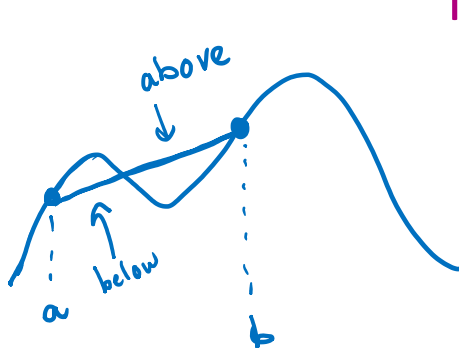
CONCAVE



CONVEX

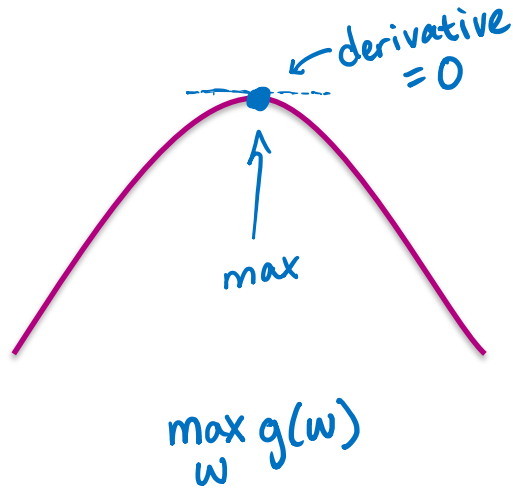


NEITHER

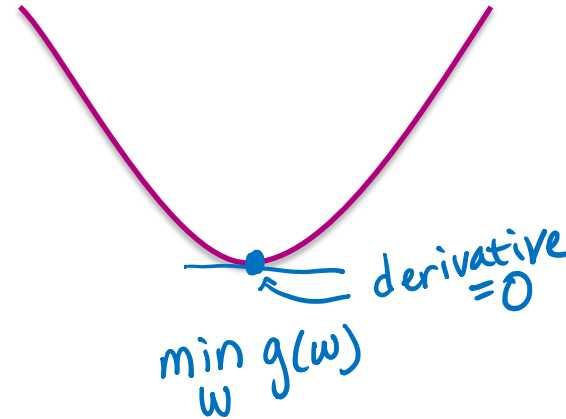


Finding the max or min analytically

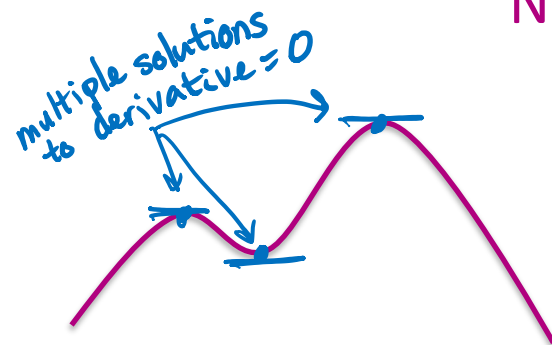
CONCAVE



CONVEX



NEITHER



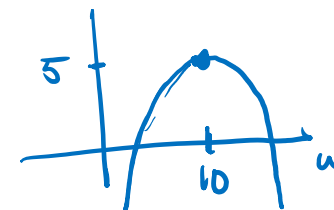
Example:

$$g(w) = 5 - (w-10)^2$$

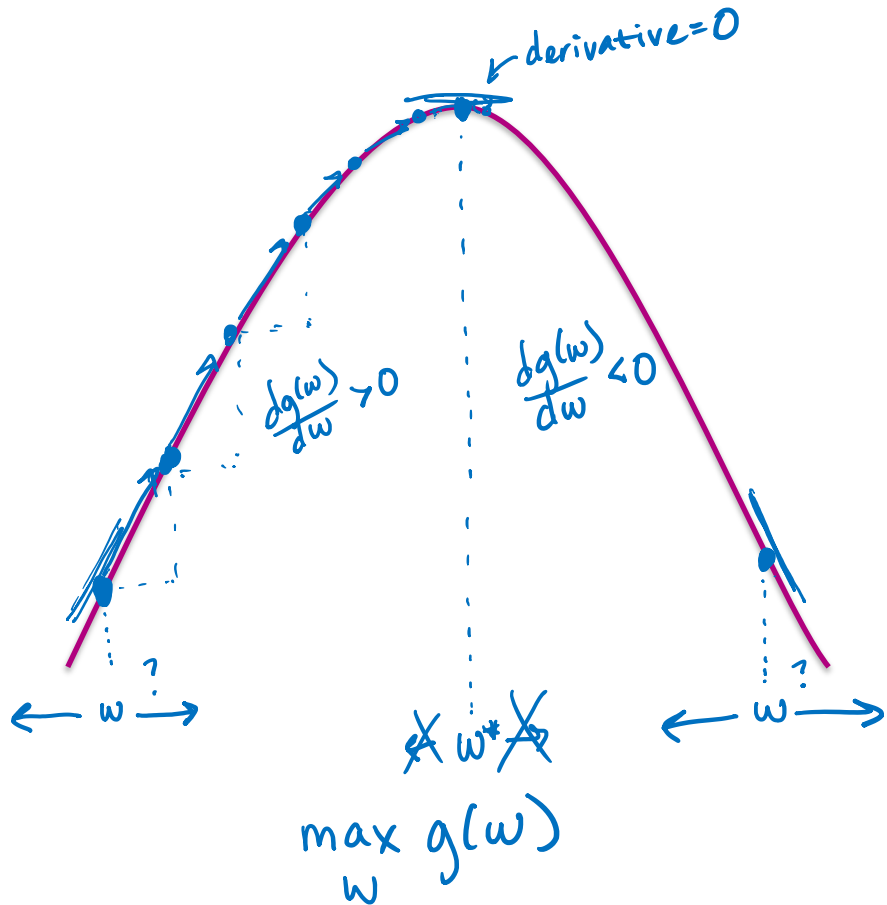
$$\begin{aligned} \frac{dg(w)}{dw} &= 0 - 2(w-10) \cdot 1 \\ &= -2w + 20 \end{aligned}$$

set derivate = 0:

$$\begin{aligned} -2w + 20 &= 0 \\ w &= 10 \end{aligned}$$



Finding the max via hill climbing



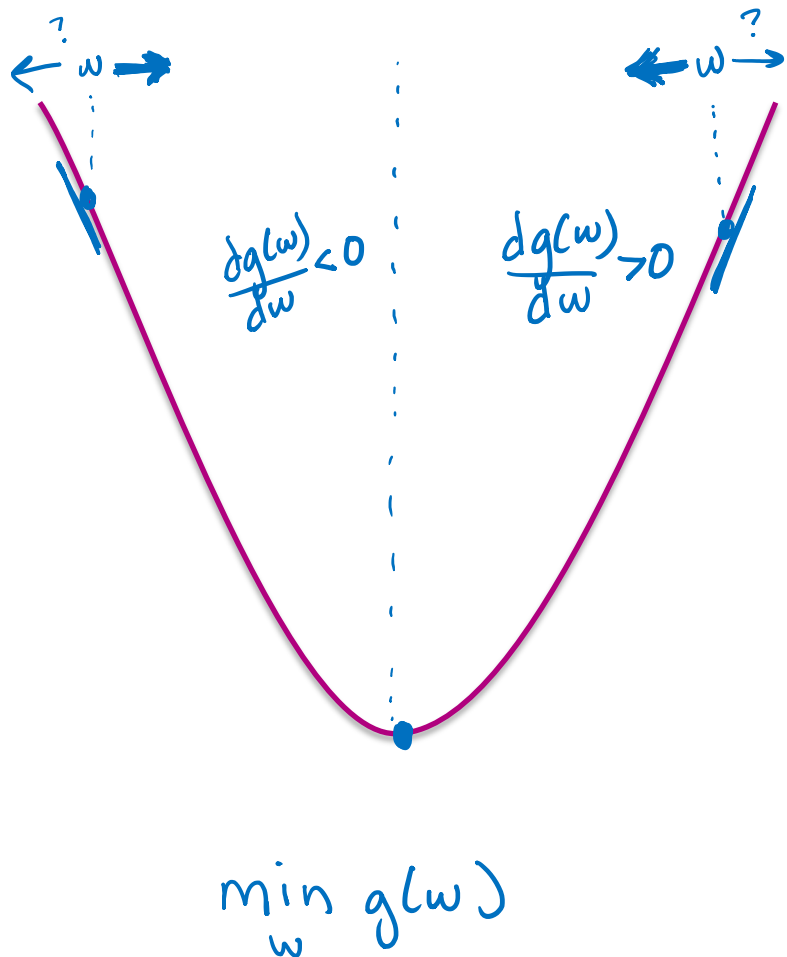
How do we know whether to move w to right or left?
(inc. or dec. the value of w ?)

while not converged

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \frac{dg(w)}{dw}$$

iteration t stepsize

Finding the min via hill descent



When derivative is positive, we want to decrease w and when derivative is negative, we want to increase w

Algorithm:

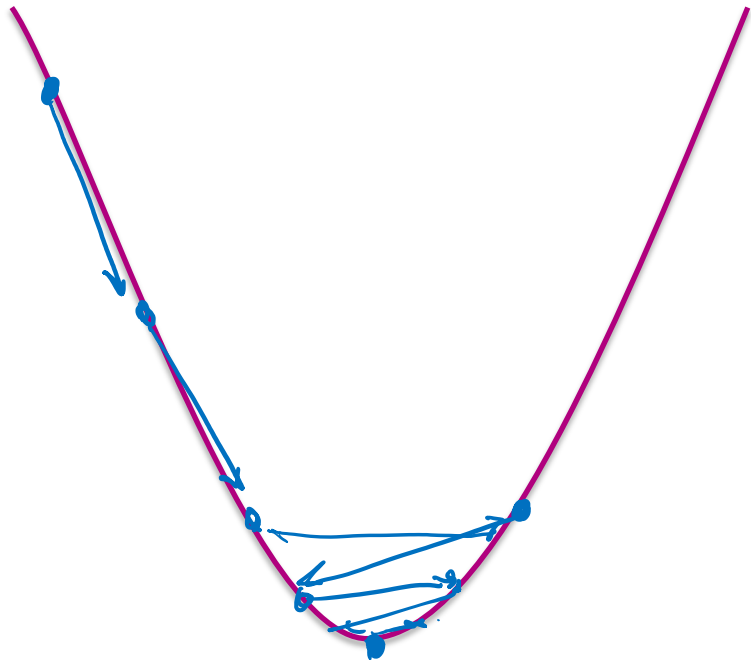
while not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \left. \frac{dg}{dw} \right|_{w^{(t)}}$$

Choosing the stepsize—

Fixed stepsize

η

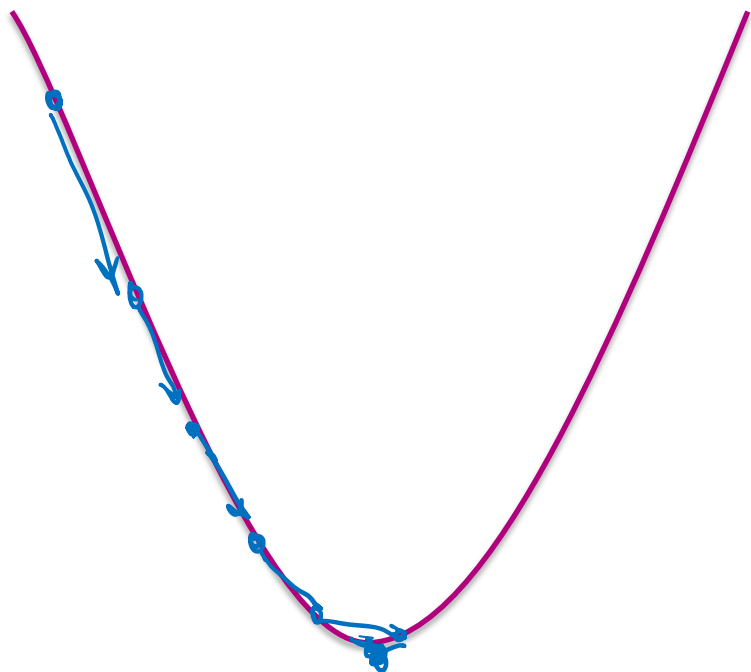


$\eta = 0.1$

Choosing the stepsize—

Decreasing stepsize

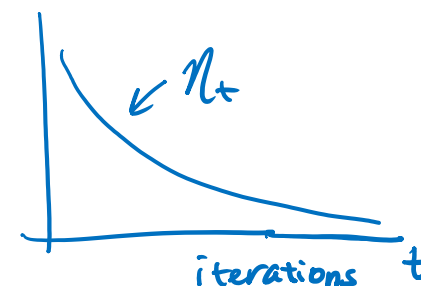
or stepsize schedule



Common choices:

$$\eta_t = \frac{\alpha}{t}$$

$$\eta_t = \frac{\alpha}{\sqrt{t}}$$



Convergence criteria

For convex functions,
optimum occurs when

$$\frac{dg(w)}{dw} = 0$$

In practice, stop when

$$\left| \frac{dg(w)}{dw} \right| < \epsilon$$

threshold to be set

Algorithm:

while not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \left. \frac{dg}{dw} \right|_{w^{(t)}}$$

Moving to higher dimensions

OPTIONAL

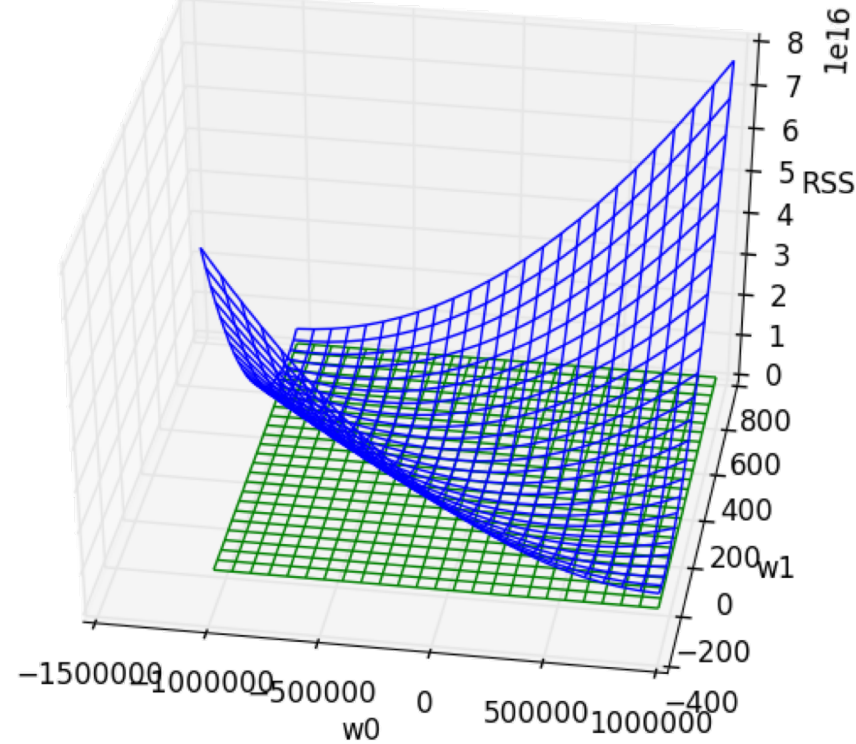
Note: We use the Optional tag to signify that you are not responsible for understanding the following material!

Contour plots

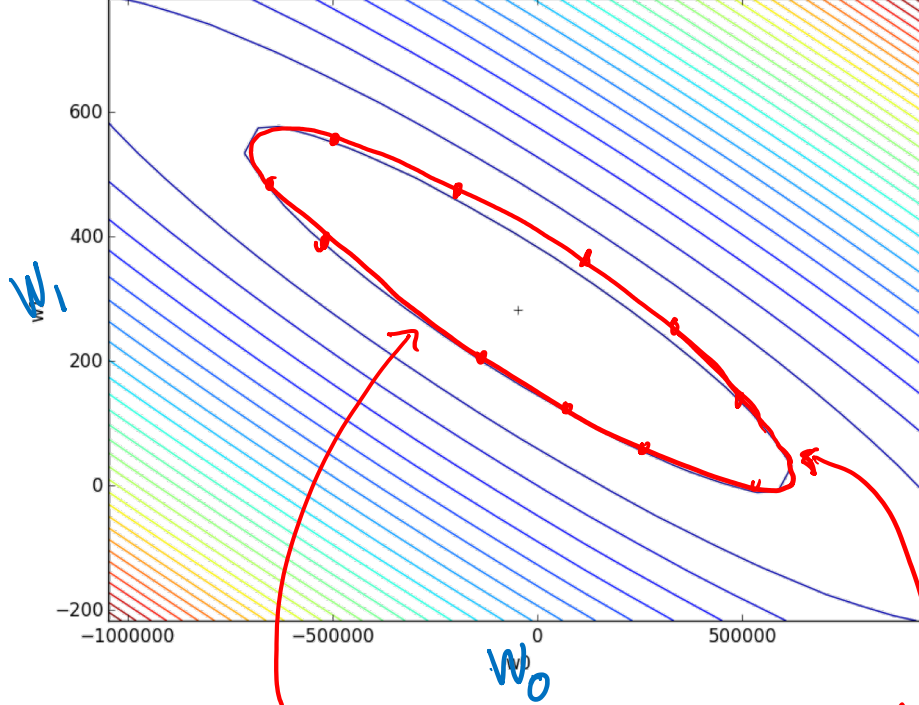
bird's eye view



3D plot of RSS with tangent plane at minimum



Contour plot corresponding to 3D plot of RSS

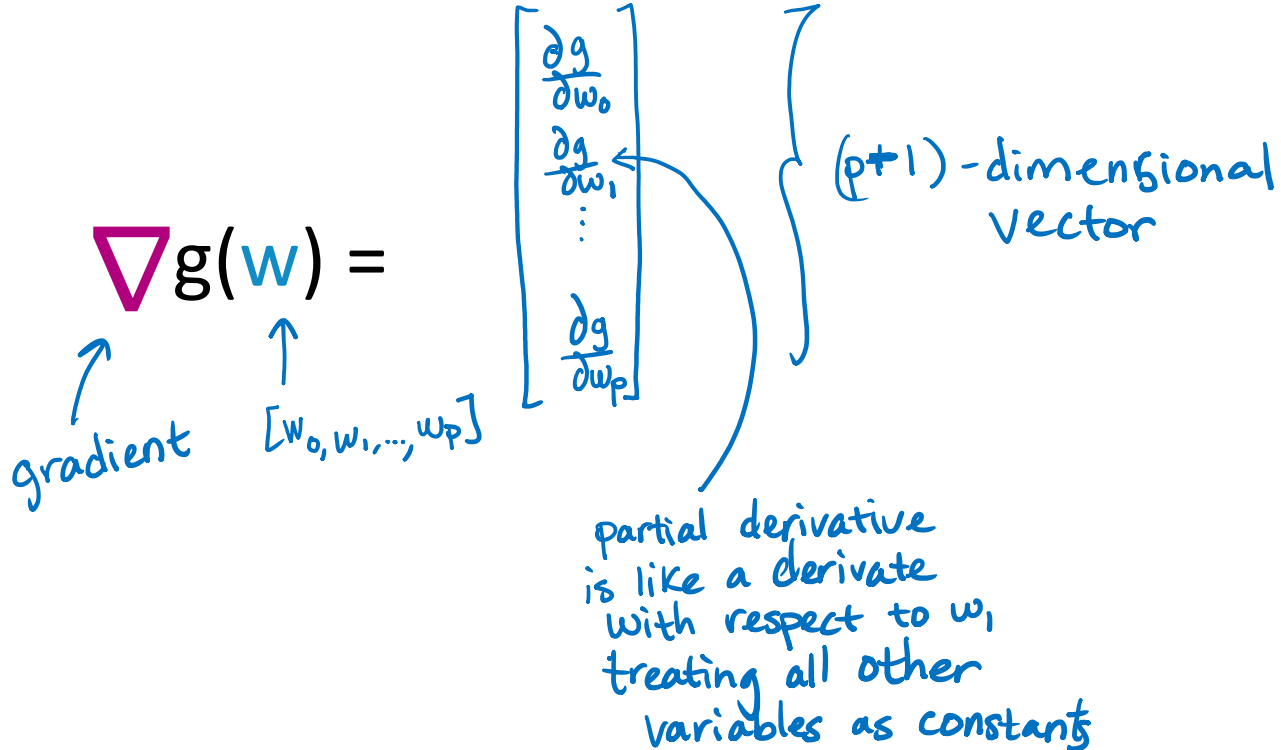
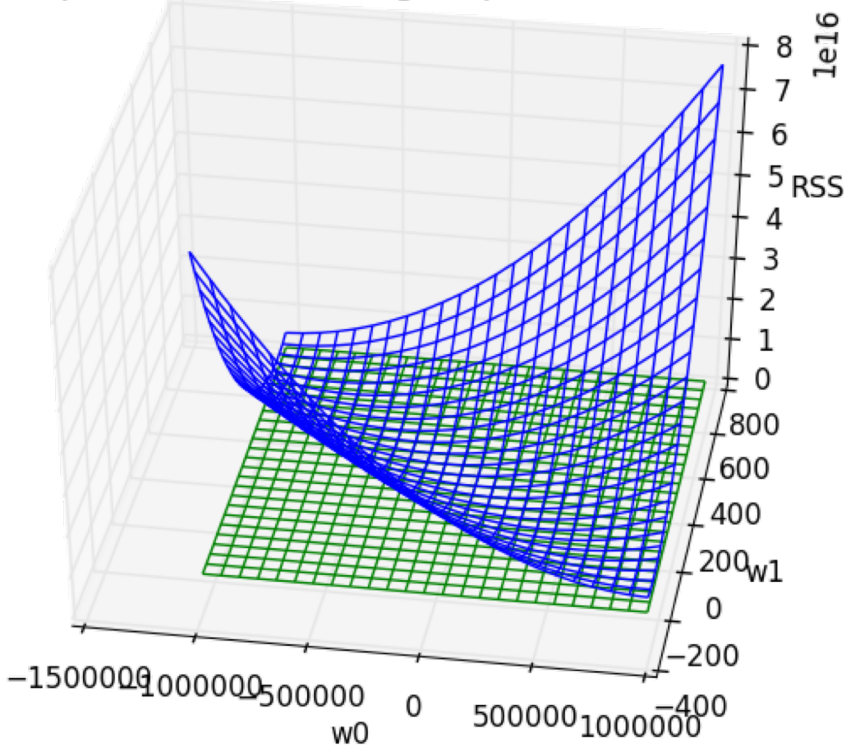


a slice of the 3D surface

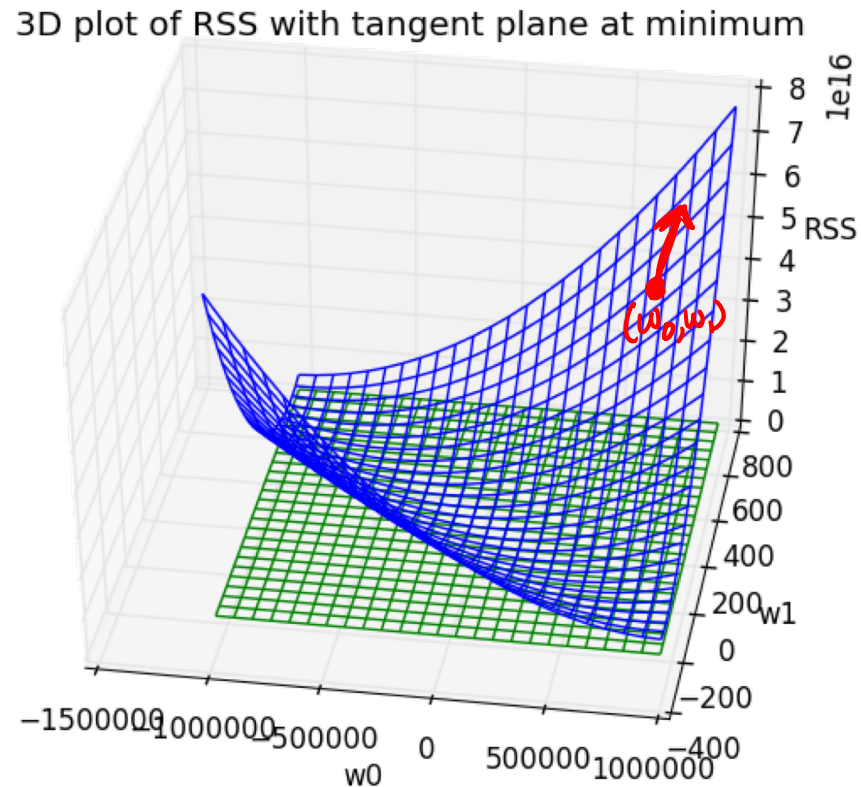
$g(w_0, w_1)$

Moving to multiple dimensions: Gradients

3D plot of RSS with tangent plane at minimum



Gradient example



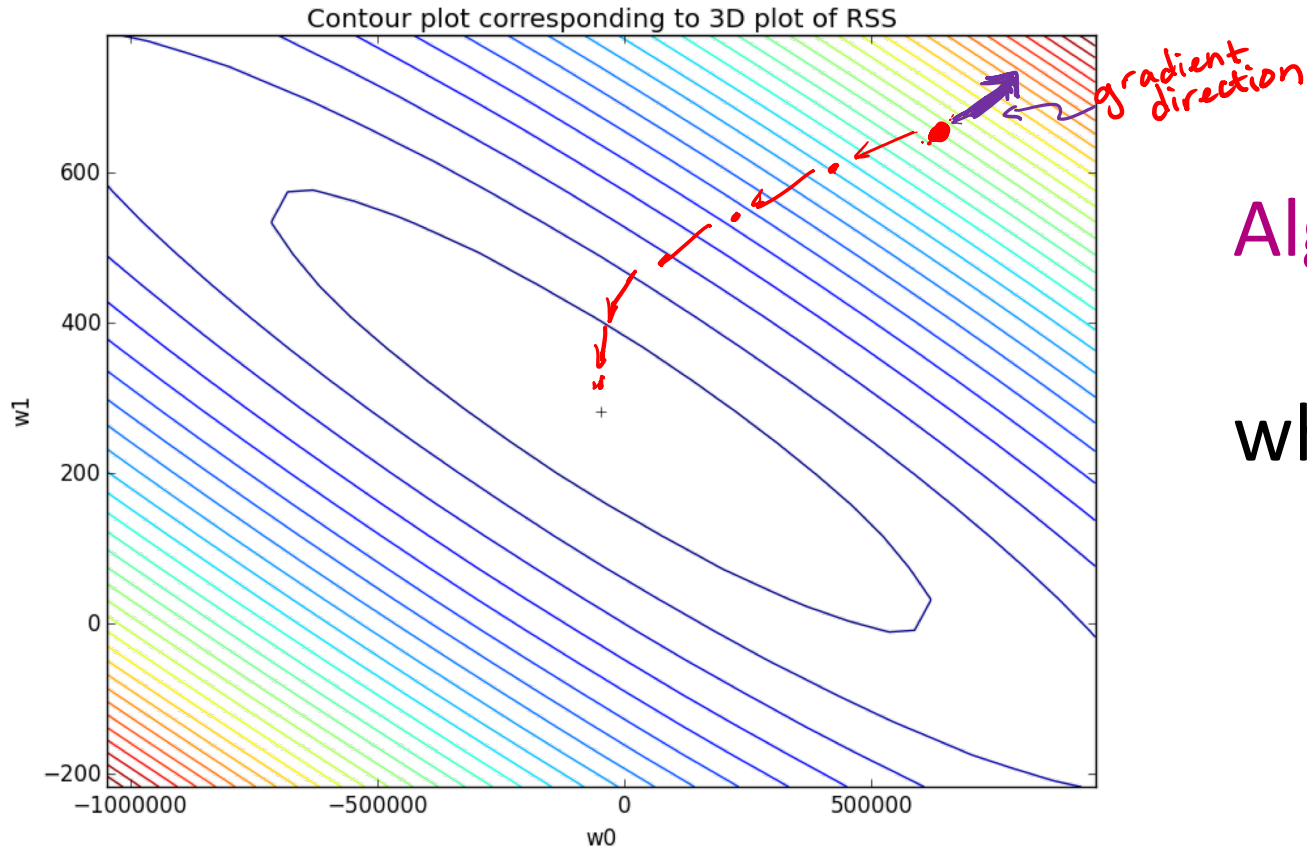
$$g(\mathbf{w}) = 5w_0 + 10w_0w_1 + 2w_1^2$$

$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

$$\frac{\partial g}{\partial w_1} = 10w_0 + 4w_1$$

$$\nabla g(\mathbf{w}) = \begin{bmatrix} 5 + 10w_1 \\ 10w_0 + 4w_1 \end{bmatrix}$$

Gradient descent



Algorithm:

while not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla g(w^{(t)})$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} - \eta \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

Convergence:
 $\|\nabla g(w)\| < \epsilon$

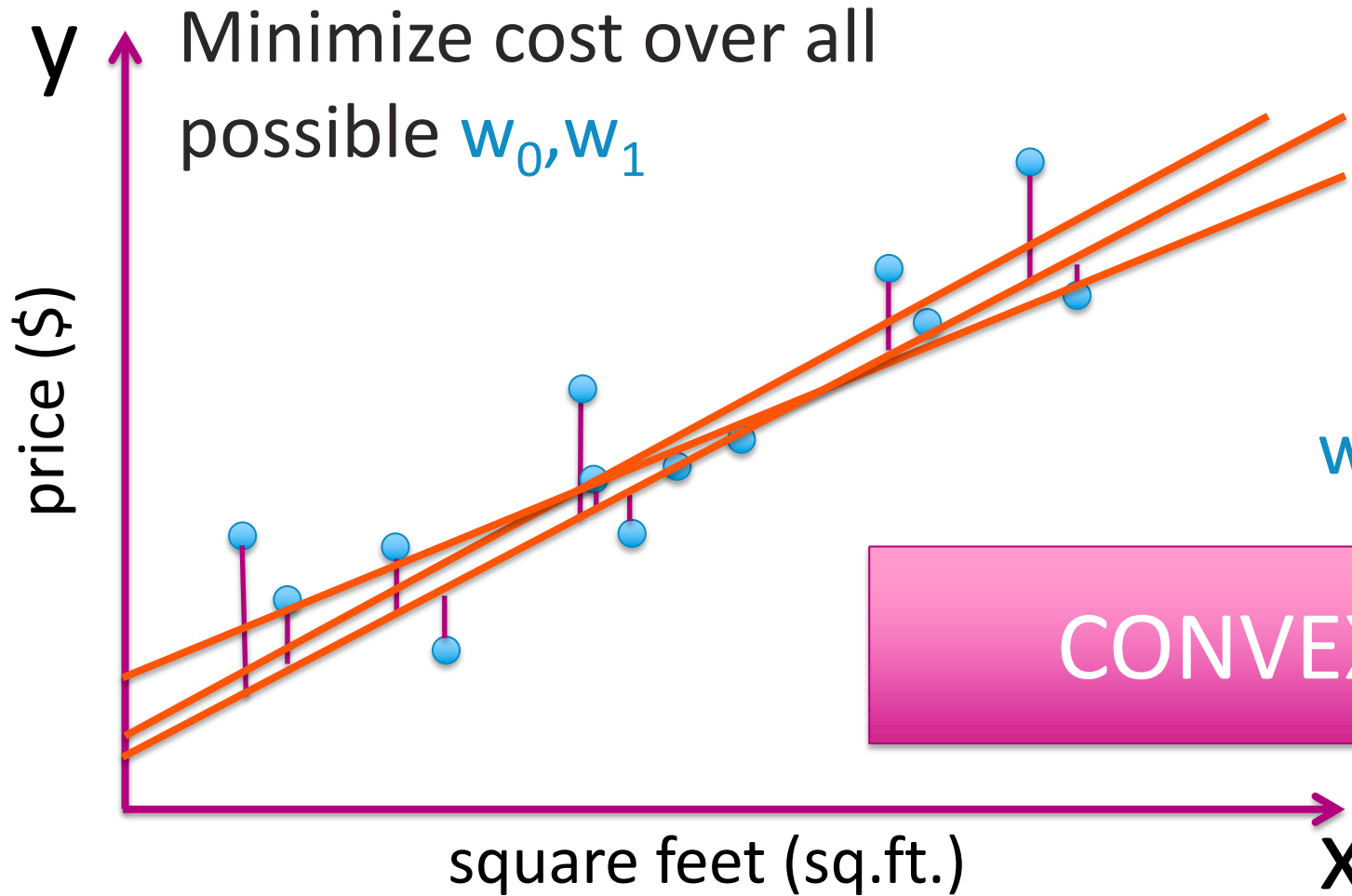
Gradient Descent for Linear Regression

OPTIONAL

Note: We use the Cap to point out that the following section contains advanced topics, passed the level we expect from the class.



Find “best” line



$$\min_{w_0, w_1} \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

⇒ solution is unique
+ gradient descent alg. will converge to minimum

Compute the gradient

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

Aside:

$$\begin{aligned} \frac{d}{dw} \sum_{i=1}^N g_i(w) &= \frac{d}{dw} (g_1(w) + g_2(w) + \dots + g_N(w)) \\ &= \frac{d}{dw} g_1(w) + \frac{d}{dw} g_2(w) + \dots + \frac{d}{dw} g_N(w) \\ &= \sum_{i=1}^N \frac{d}{dw} g_i(w) \end{aligned}$$

In our case

$$g_i(w) = (y_i - [w_0 + w_1 x_i])^2$$

$$\frac{\partial \text{RSS}(w)}{\partial w_0} = \sum_{i=1}^N \frac{\partial}{\partial w_0} (y_i - [w_0 + w_1 x_i])^2$$

same for w_1

Compute the gradient

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

Taking the derivative w.r.t. w_0 _____

$$\sum_{i=1}^N 2 (y_i - [w_0 + w_1 x_i])' \cdot (-1)$$

$$= -2 \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])$$

Compute the gradient

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

Taking the derivative w.r.t. w_1

$$\begin{aligned} & \sum_{i=1}^N 2(y_i - [w_0 + w_1 x_i]) \cdot (-x_i) \\ &= -2 \sum_{i=1}^N (y_i - [w_0 + w_1 x_i]) \underline{x_i} \end{aligned}$$

Compute the gradient

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

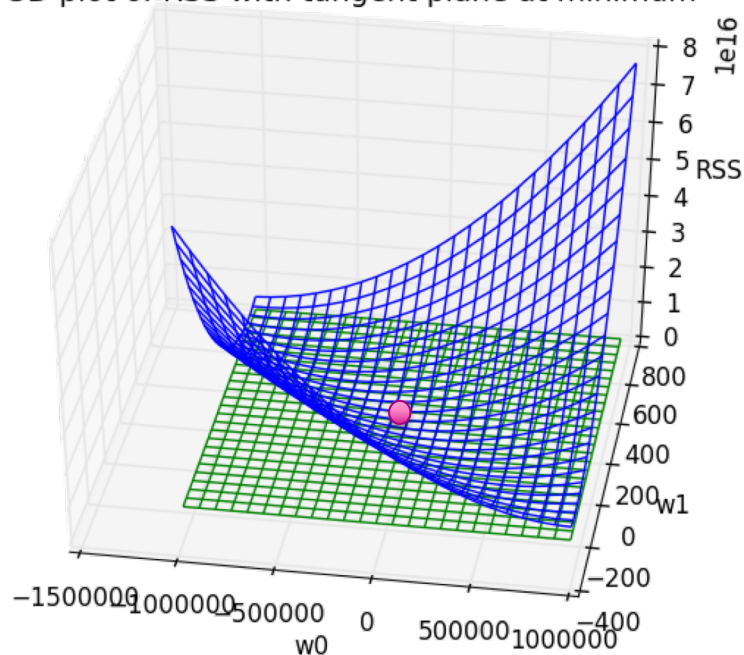
Putting it together:

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

Approach 1: Set gradient = 0

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

3D plot of RSS with tangent plane at minimum



top term: $\hat{w}_0 = \frac{\sum_{i=1}^N y_i}{N} - \hat{w}_1 \frac{\sum_{i=1}^N x_i}{N}$

average sales price $\leftarrow \frac{\sum y_i}{N}$
 estimate of the slope $\leftarrow \hat{w}_1$
 average sq-ft. $\leftarrow \frac{\sum x_i}{N}$

bottom term: $\sum y_i x_i - \hat{w}_0 \sum x_i - \hat{w}_1 \sum x_i^2 = 0$

plug in $\rightarrow \hat{w}_1 = \frac{\sum y_i x_i - \frac{\sum y_i \sum x_i}{N}}{\sum x_i^2 - \frac{\sum x_i \sum x_i}{N}}$

Note:

$$\sum_{i=1}^N y_i$$

$$\sum_{i=1}^N x_i$$

$$\sum_{i=1}^N y_i x_i$$

$$\sum_{i=1}^N x_i^2$$

Approach 2: Gradient descent

Interpreting the gradient:

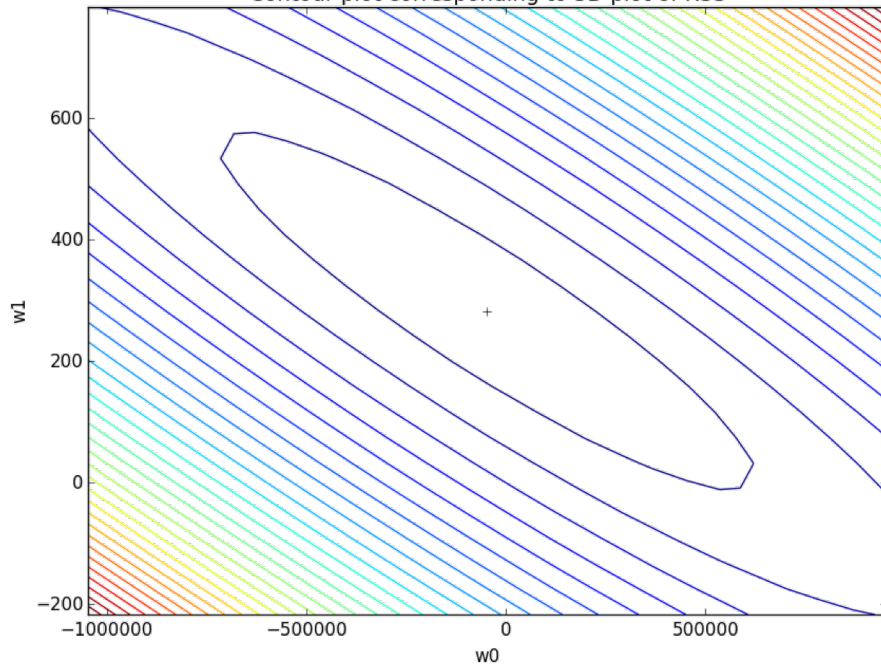
$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] x_i \end{bmatrix}$$

actual house sales observation
predicted value $\hat{y}_i(w_0, w_1)$

Approach 2: Gradient descent

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] x_i \end{bmatrix}$$

Contour plot corresponding to 3D plot of RSS



while not converged $(-2) \cdot (-n)$

$$\begin{bmatrix} w_0^{(t+1)} \\ w_1^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} w_0^{(t)} \\ w_1^{(t)} \end{bmatrix} + 2\eta \begin{bmatrix} \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] \\ \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] x_i \end{bmatrix}$$

If overall, under predicting \hat{y}_i , then $\sum [y_i - \hat{y}_i]$ is positive

→ w_0 is going to increase

similar intuition for w_1 , but multiply by x_i

Comparing the approaches

- For most ML problems, cannot solve $\text{gradient} = 0$
- Even if solving $\text{gradient} = 0$ is feasible, gradient descent can be more efficient
- Gradient descent relies on choosing stepsize and convergence criteria