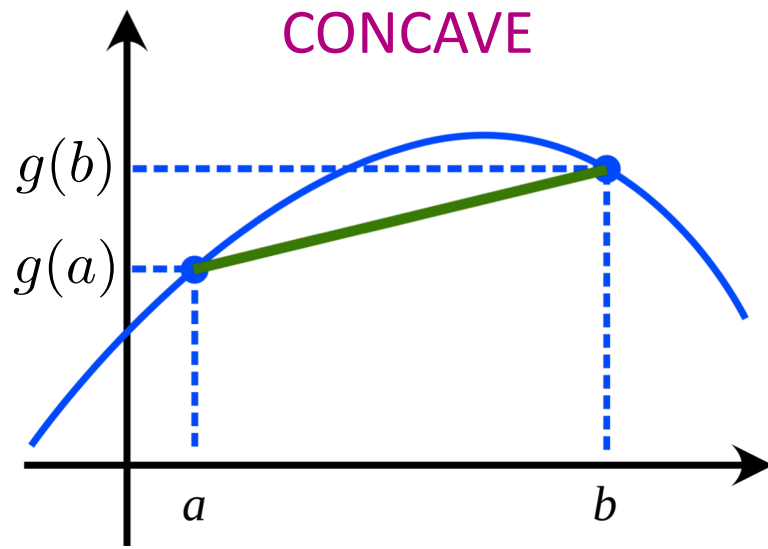


# Gradient descent

# Convex/concave functions



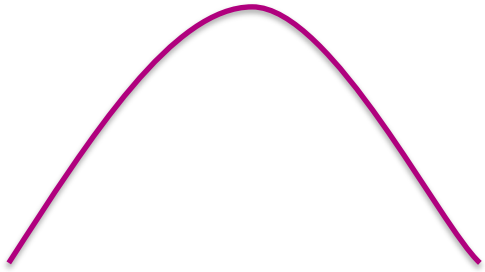
CONCAVE

CONVEX

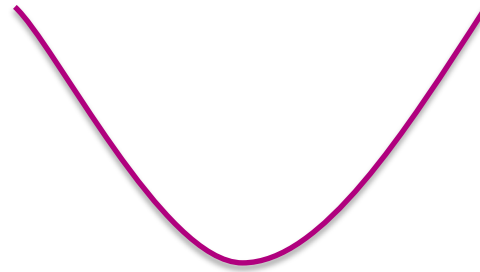
NEITHER

# Finding the max or min analytically

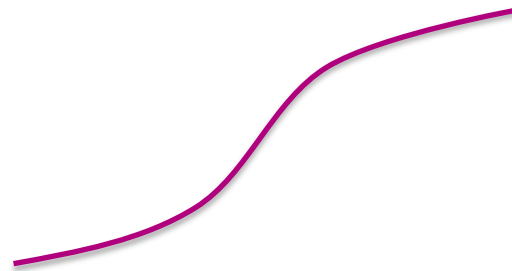
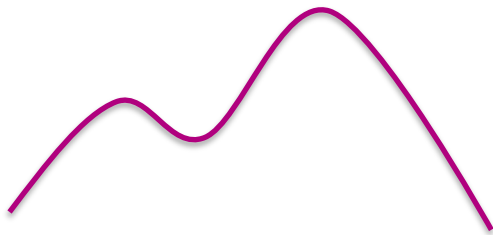
CONCAVE



CONVEX



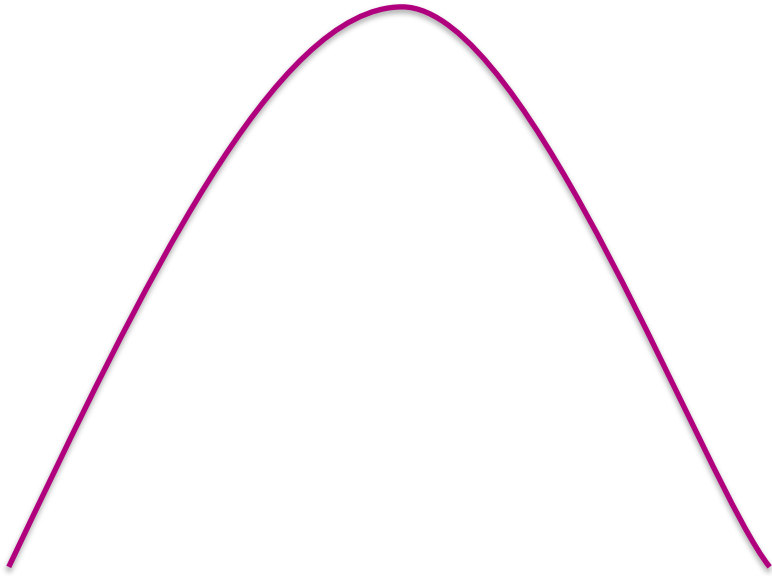
NEITHER



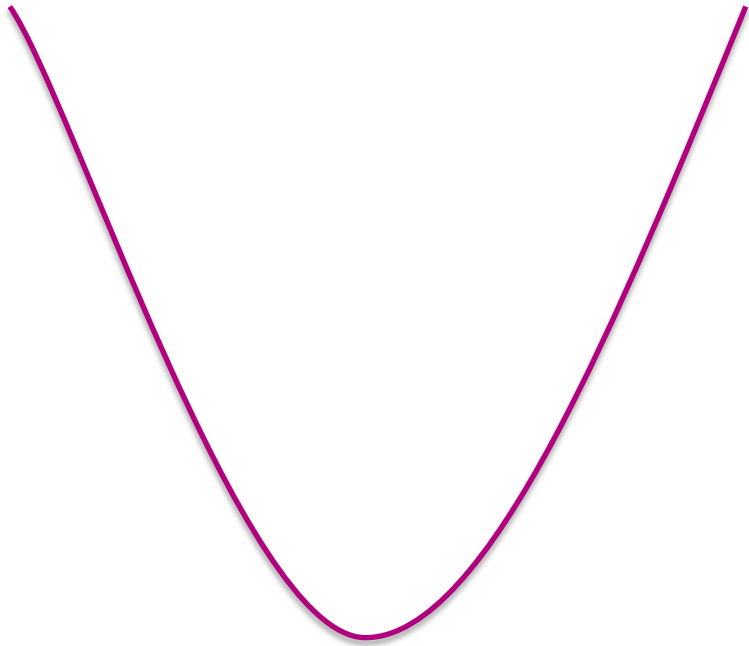
Example:

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

# Finding the max via hill climbing



# Finding the min via hill descent

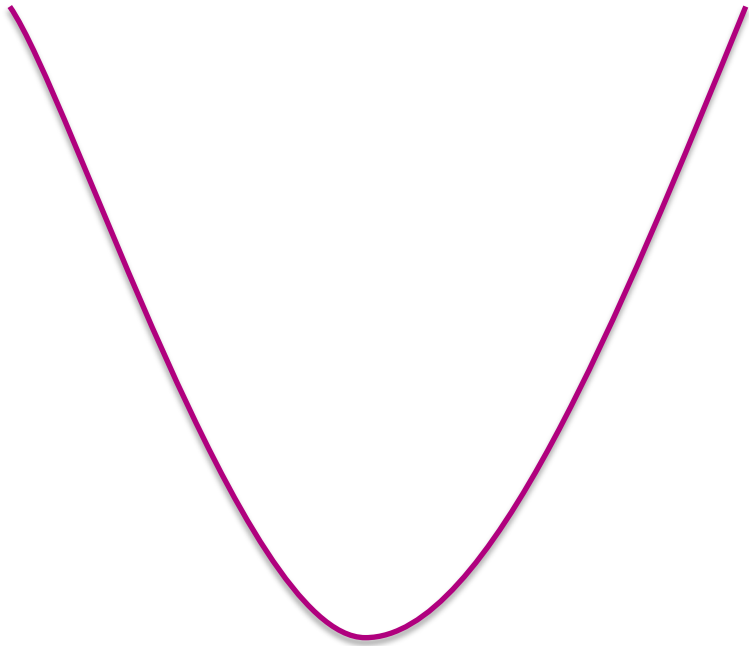


Algorithm:

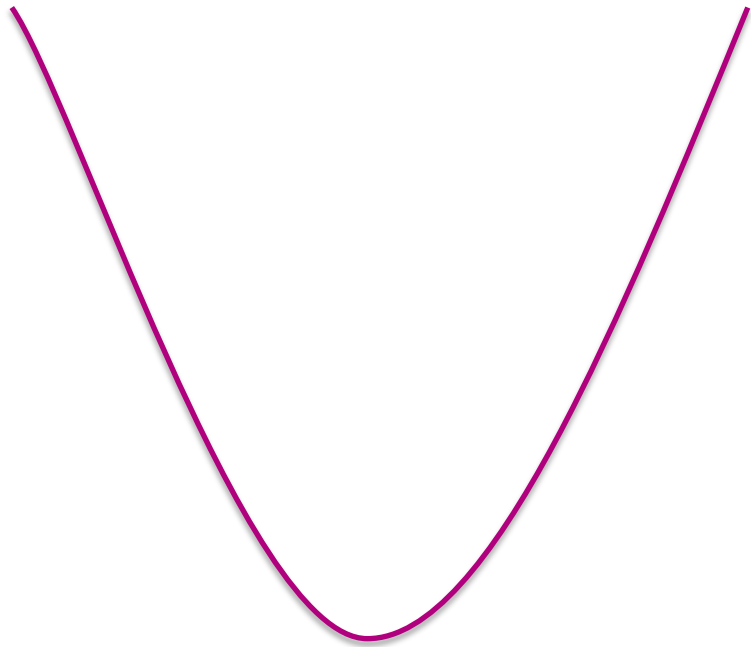
while not converged

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

# Choosing the stepsize— Fixed stepsize



# Choosing the stepsize— Decreasing stepsize



Common choices:

# Convergence criteria

For convex functions,  
optimum occurs when

In practice, stop when

Algorithm:

while not converged

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



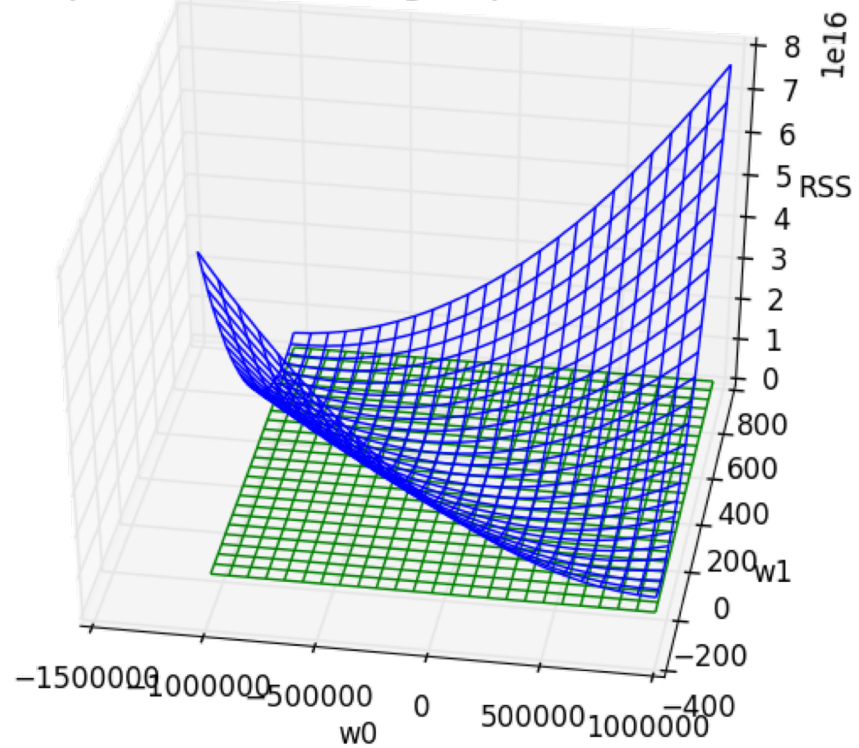
# Moving to higher dimensions

**OPTIONAL**

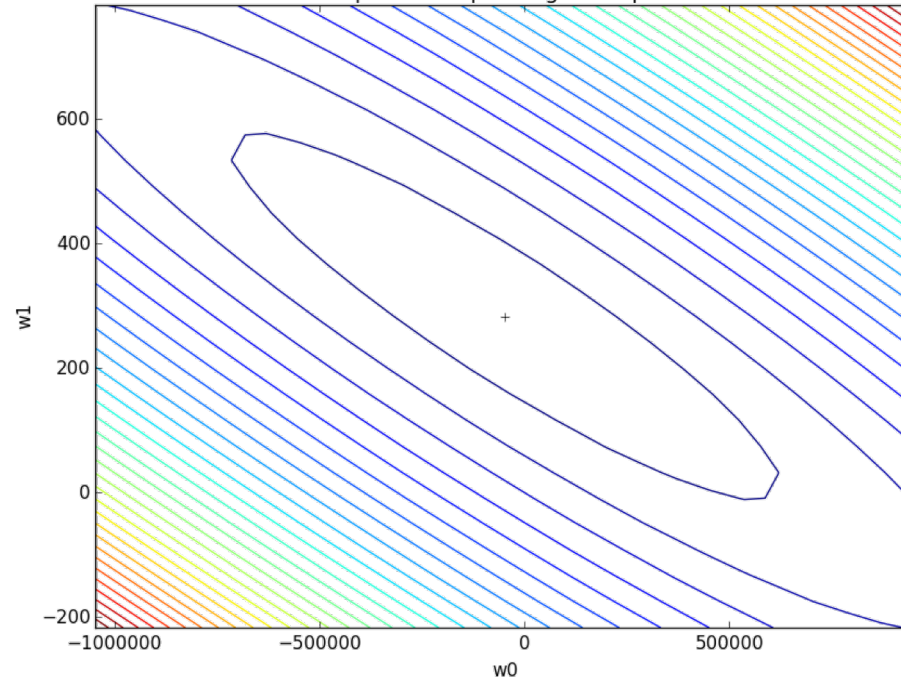
Note: We use the Optional tag to signify that you are not responsible for understanding the following material!

# Contour plots

3D plot of RSS with tangent plane at minimum

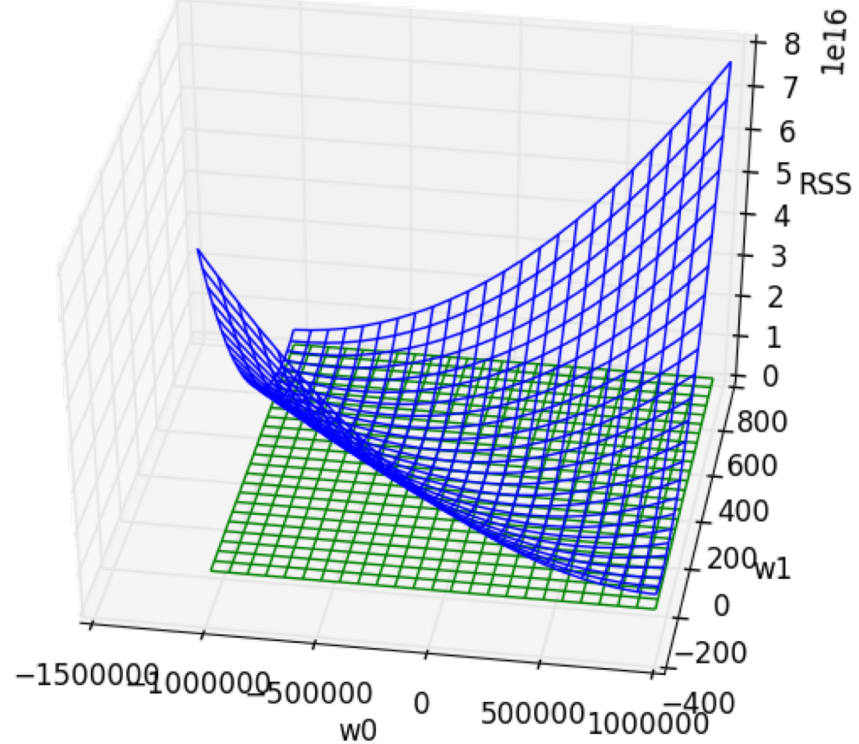


Contour plot corresponding to 3D plot of RSS



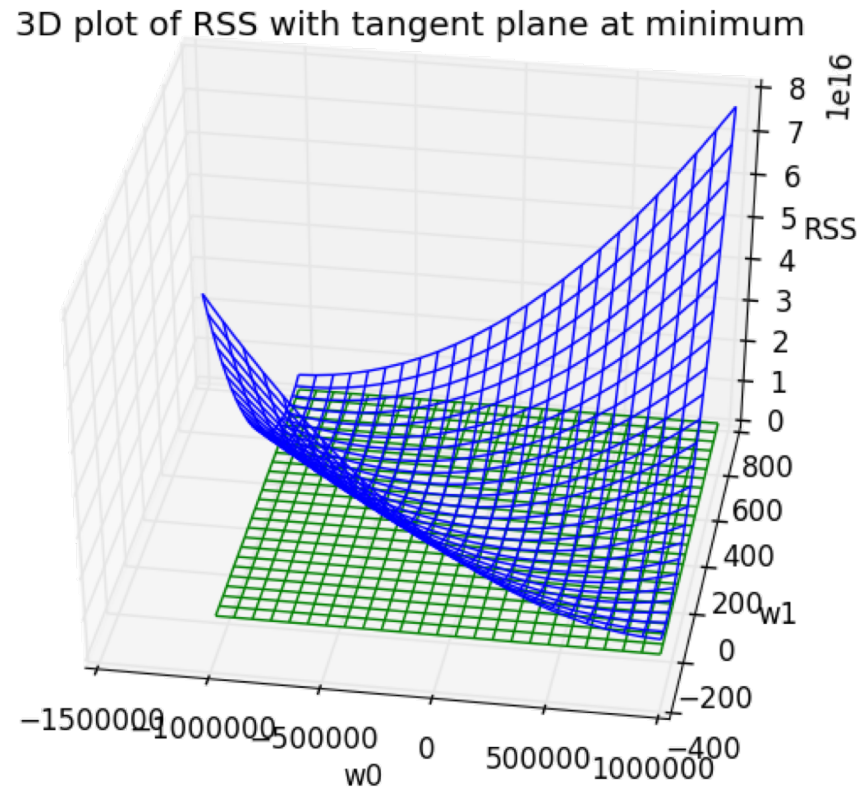
# Moving to multiple dimensions: Gradients

3D plot of RSS with tangent plane at minimum



$$\nabla g(\mathbf{w}) =$$

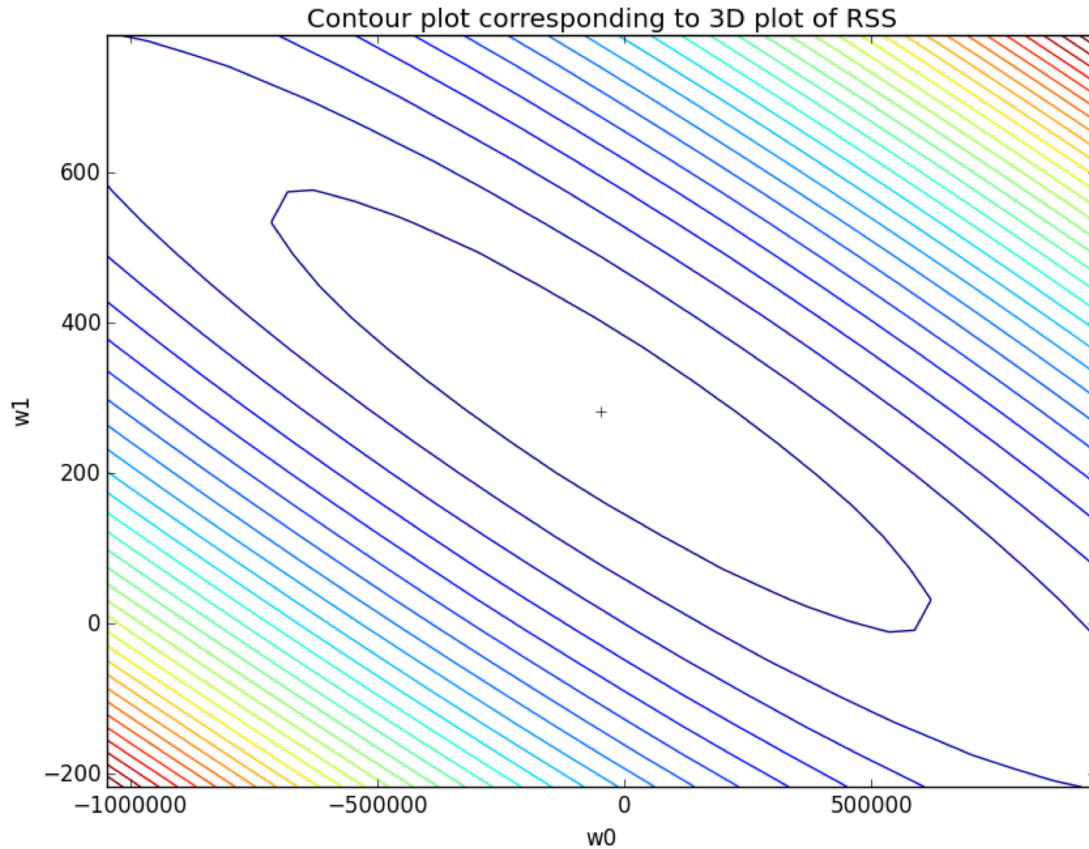
# Gradient example



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

$$\nabla g(\mathbf{w}) =$$

# Gradient descent



Algorithm:

while not converged

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

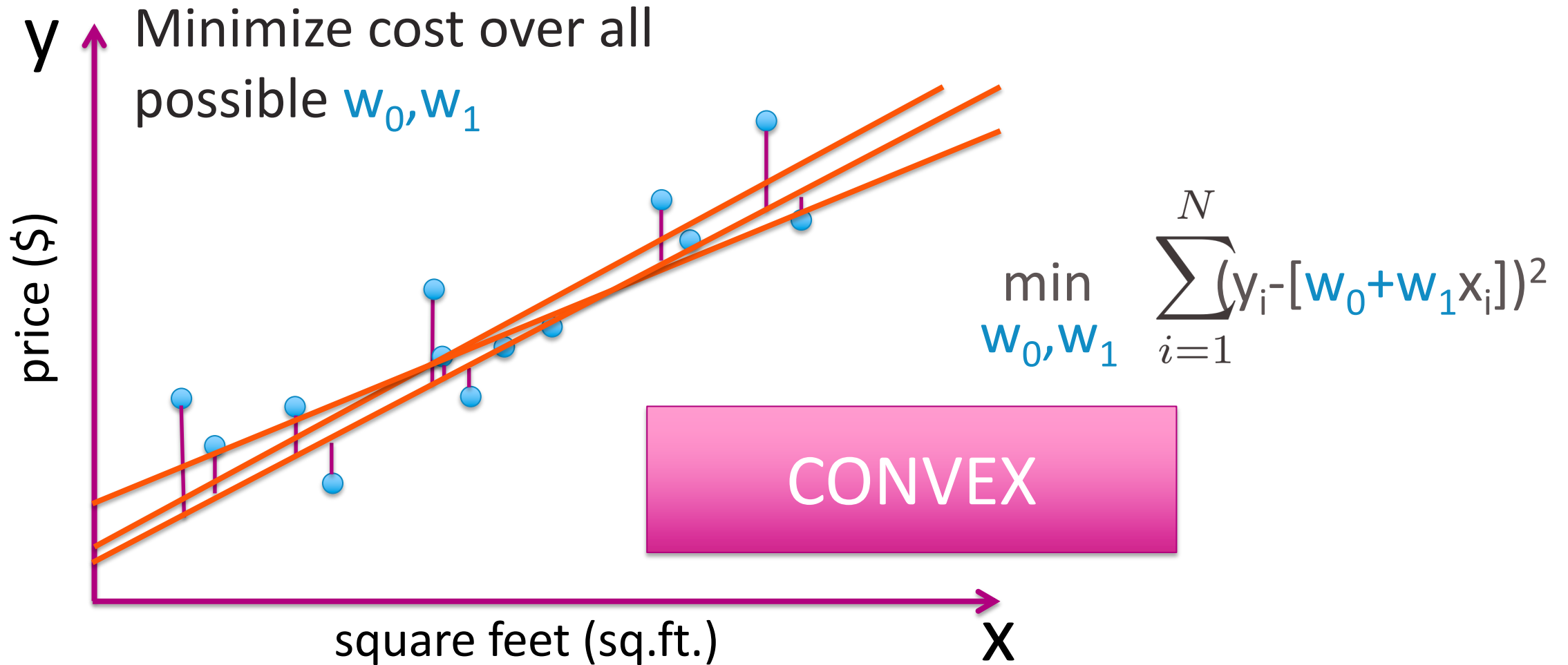
# Gradient Descent for Linear Regression

**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



# Compute the gradient

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

Aside:

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



# 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.  $w_0$

# 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.  $w_1$

# Compute the gradient

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

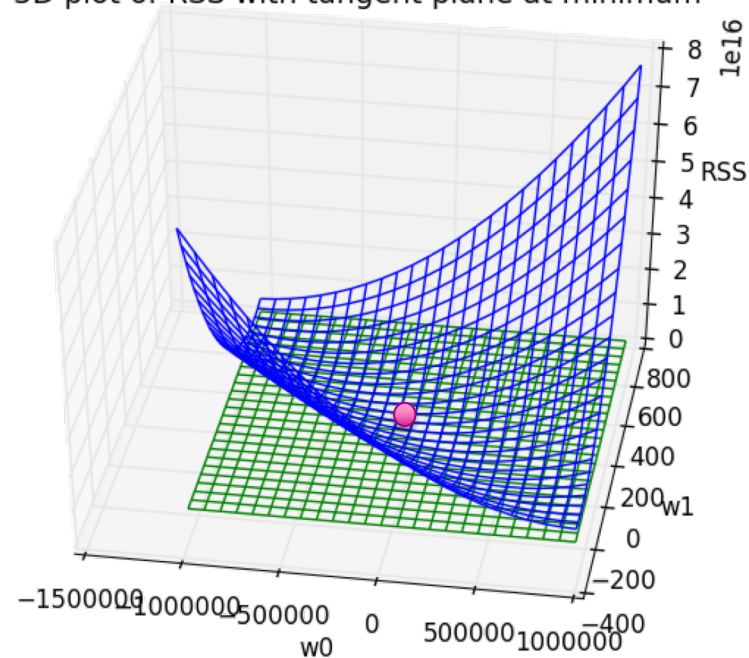
Putting it together:

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

# 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



# Approach 2: Gradient descent

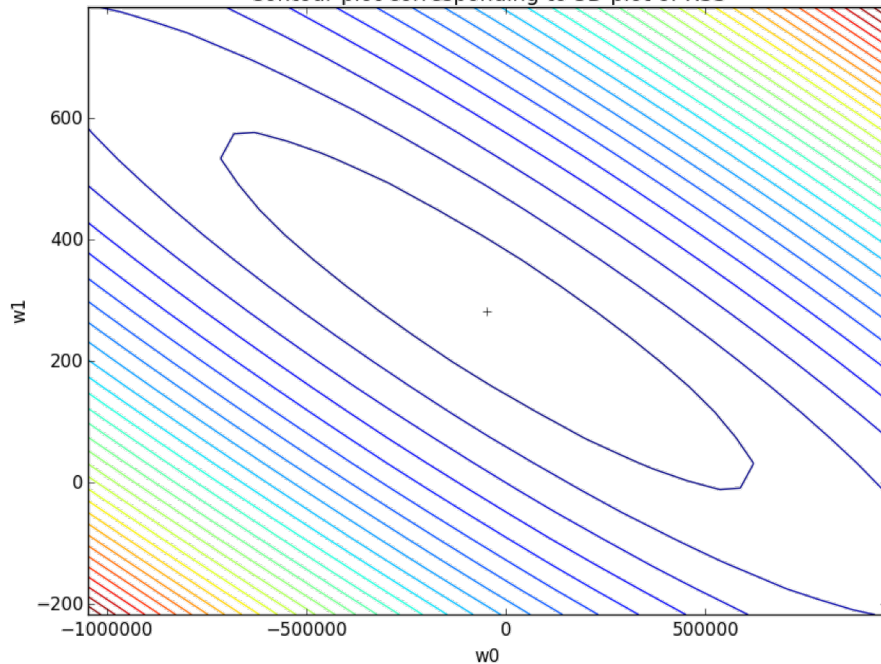
Interpreting the gradient:

$$\nabla_{\mathbf{w}_0, \mathbf{w}_1} \text{RSS}(\mathbf{w}_0, \mathbf{w}_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (\mathbf{w}_0 + \mathbf{w}_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (\mathbf{w}_0 + \mathbf{w}_1 x_i)] x_i \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] x_i \end{bmatrix}$$

# 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



# Comparing the approaches

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