

Introduction to Machine Learning



Time For Some Cool Demonstrations!

Contents

- ❖ What is Machine Learning
- ❖ Types of Machine Learning
- ❖ Linear Regression - with 1 feature
- ❖ Gradient Descent - Learning Algorithm
- ❖ Learning Rate - How to decide ?
- ❖ Effect of High-leverage points
- ❖ Effect of High-end outliers
- ❖ Linear Regression - with multiple features

What is Machine Learning

"the field of study that gives computers the ability to learn without being explicitly programmed."

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E."

Types of Machine Learning

Supervised Learning

- Classification
- Regression
- Ranking

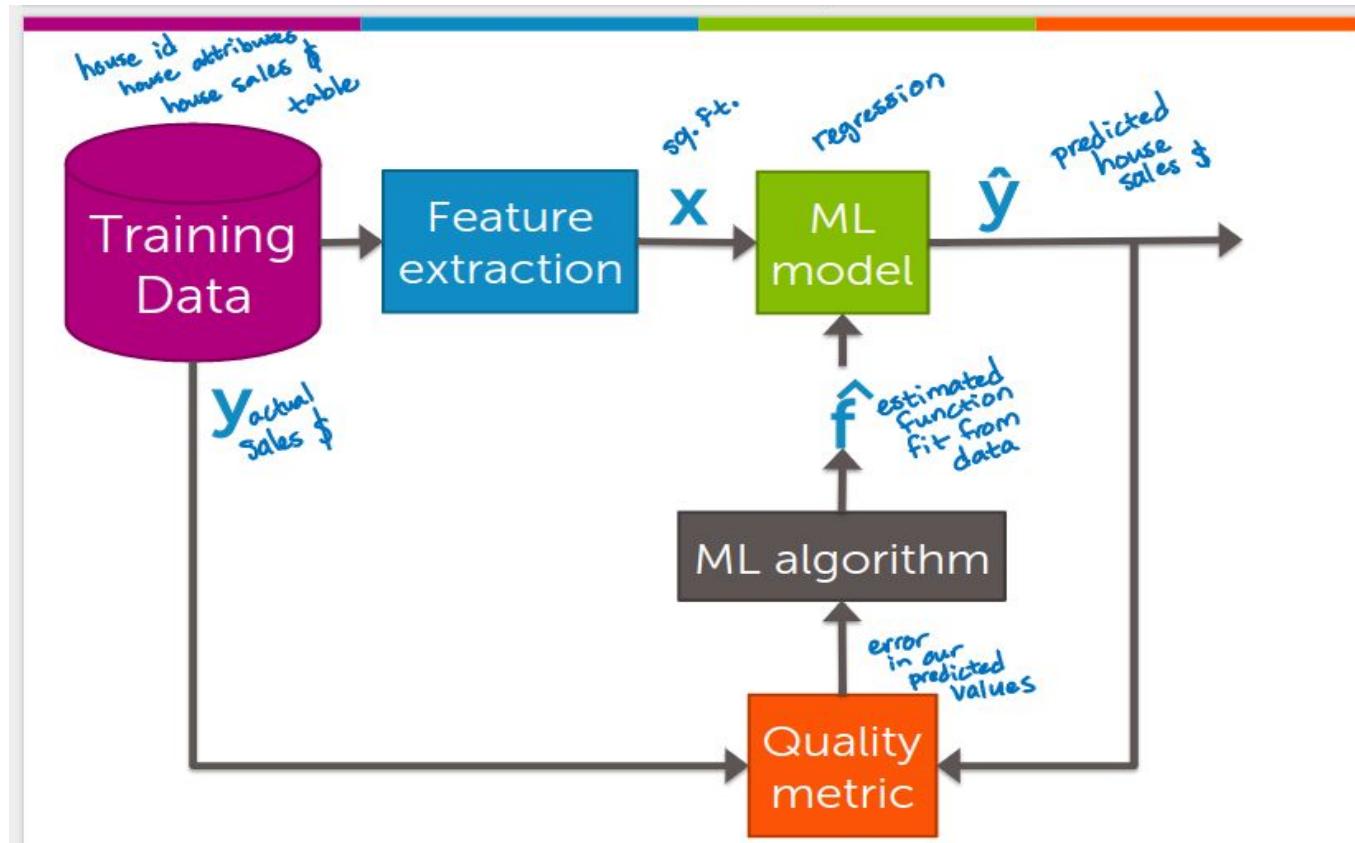
Unsupervised Learning

- Clustering
- Association Mining
- Segmentation
- Dimension Reduction

Reinforcement Learning

- Decision Process
- Reward System
- Recommendation Systems

How does supervised learning work ?



Regression Model

How much is my house worth?



Look at recent sales in my neighborhood

- How much did they sell for?



Data



input *output*
 $(x_1 = \text{sq.ft.}, y_1 = \$)$



$(x_2 = \text{sq.ft.}, y_2 = \$)$



$(x_3 = \text{sq.ft.}, y_3 = \$)$



$(x_4 = \text{sq.ft.}, y_4 = \$)$



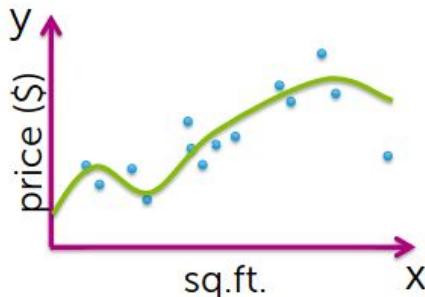
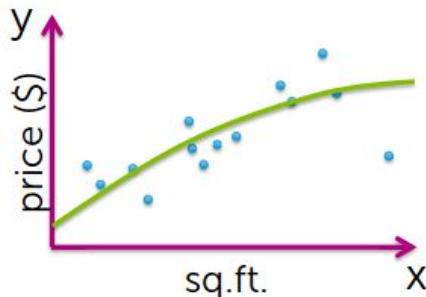
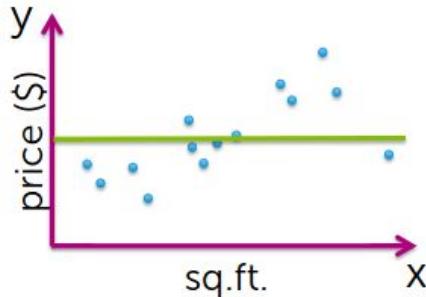
$(x_5 = \text{sq.ft.}, y_5 = \$)$

⋮

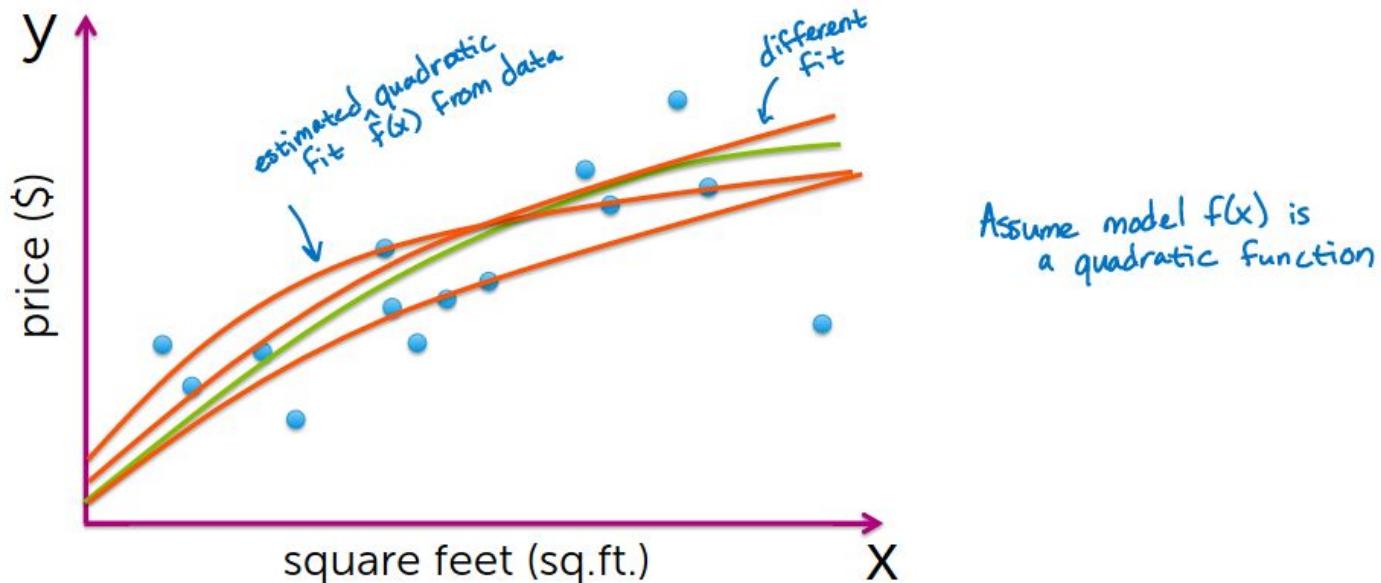
Input vs. Output:

- y is the quantity of interest
- assume y can be predicted from x

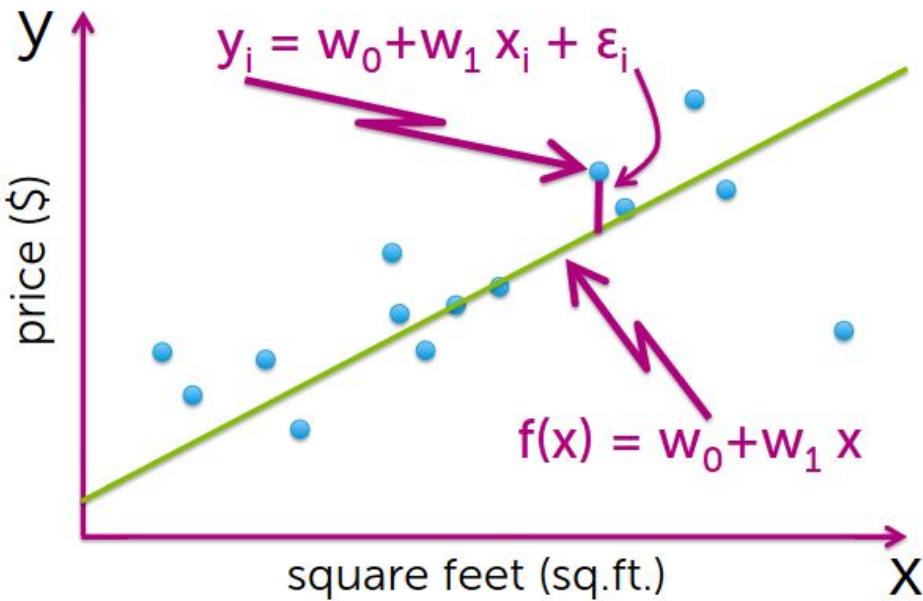
Task 1– Which model $f(x)$?



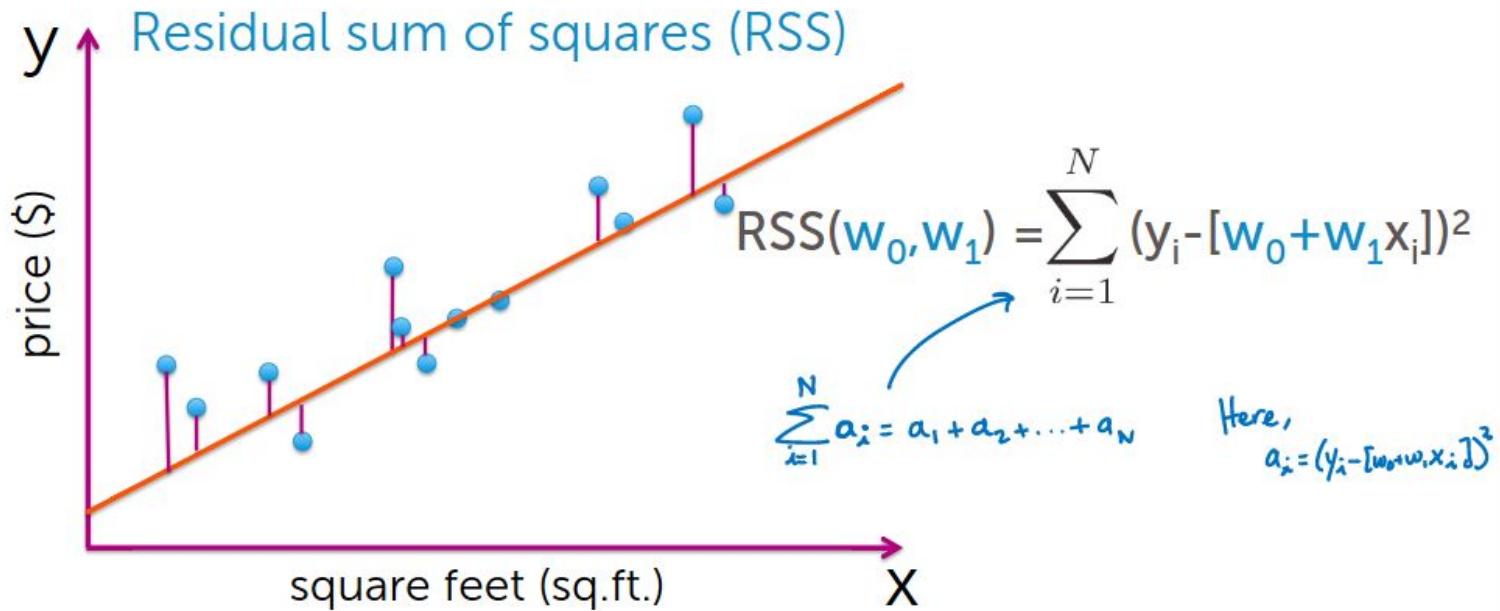
Task 2 – For a given model $f(x)$, estimate function $\hat{f}(x)$ from data



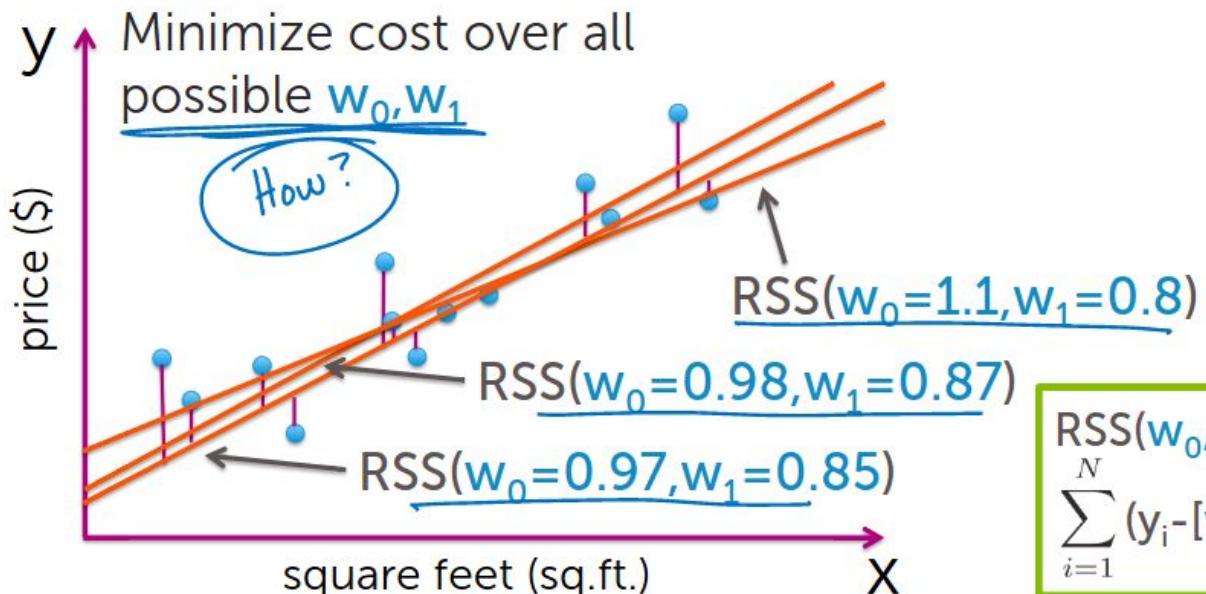
Simple linear regression model



"Cost" of using a given line



Find “best” line

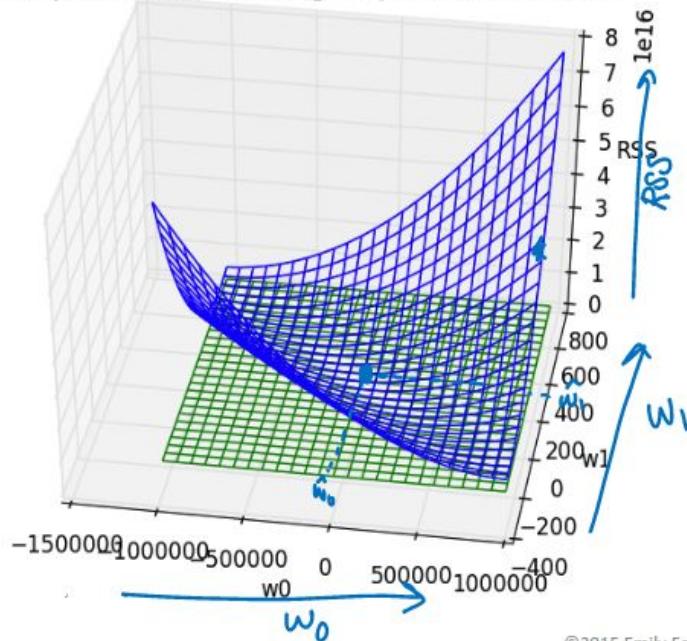


Recall:

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

Minimizing the cost

3D plot of RSS with tangent plane at minimum

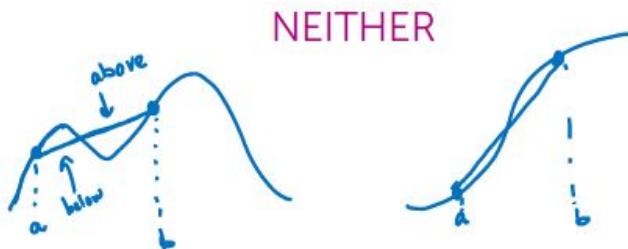
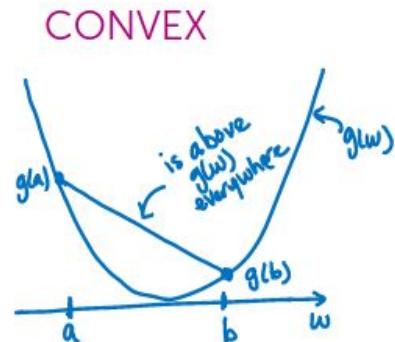
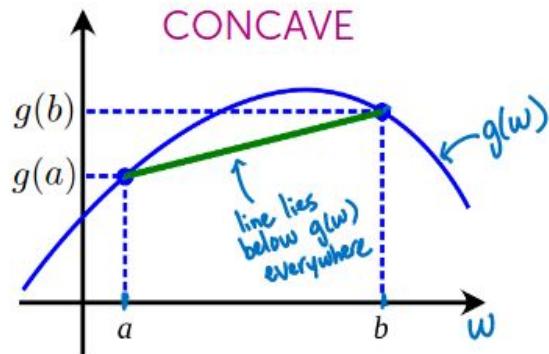


Minimize function
over all possible w_0, w_1

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

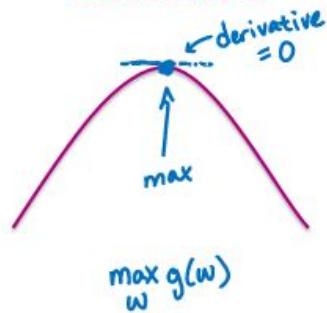
RSS(w_0, w_1) is a function
of 2 variables = $g(w_0, w_1)$

Convex/concave functions



Finding the max or min analytically

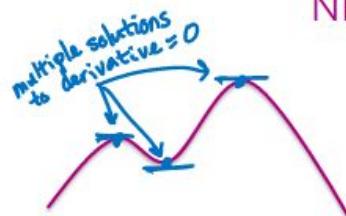
CONCAVE



CONVEX



NEITHER

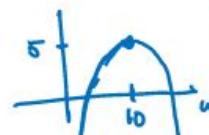


Example:

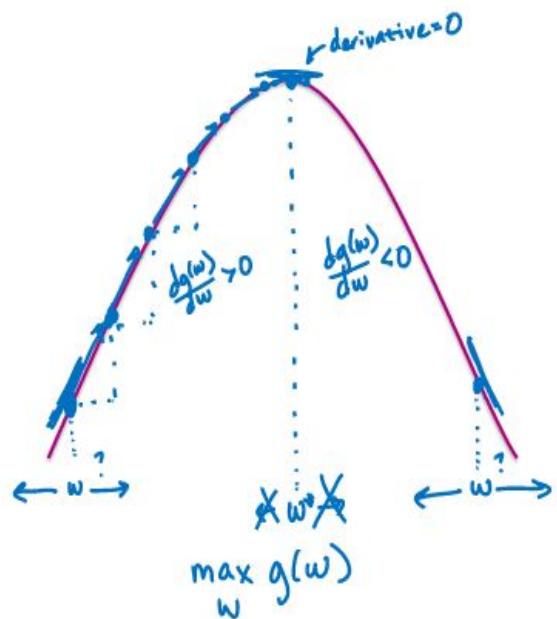
$$g(w) = 5 - (w-10)^2$$

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

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



Finding the max via hill climbing



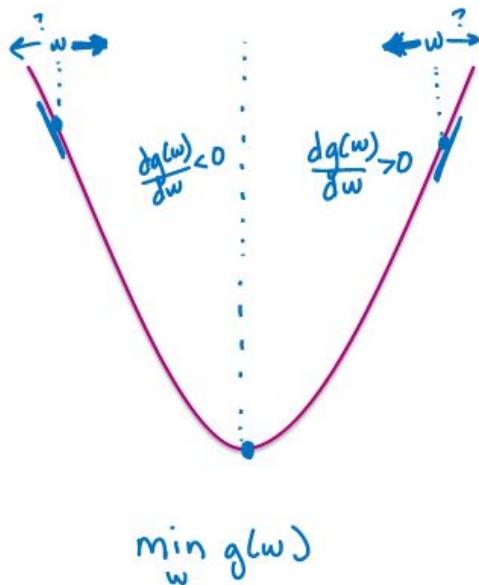
How do we know whether to move w to right or left?
(inc. or dec. the value of w ?)

while not converged

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

iteration 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

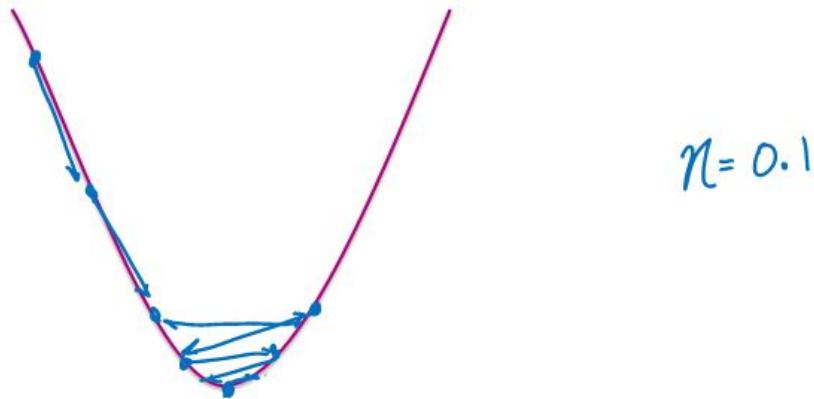
Algorithm:

while not converged

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

Choosing the stepsize— Fixed stepsize

η



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

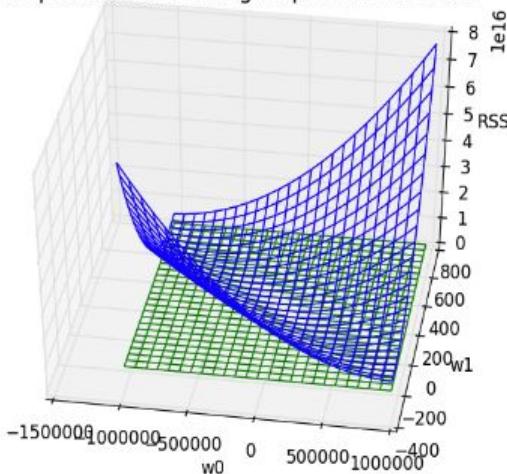
Algorithm:

while not converged

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

Moving to multiple dimensions: Gradients

3D plot of RSS with tangent plane at minimum



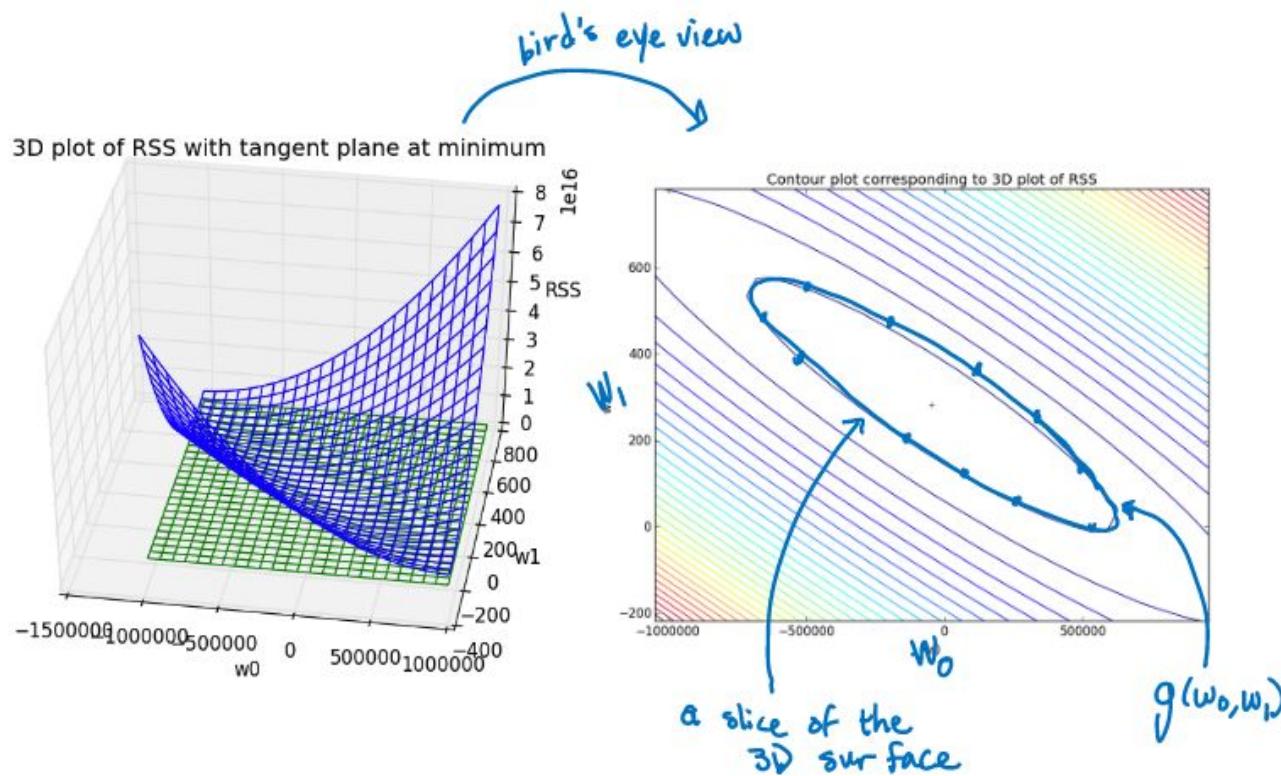
$$\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{bmatrix}$$

gradient \uparrow $[w_0, w_1, \dots, w_p]$

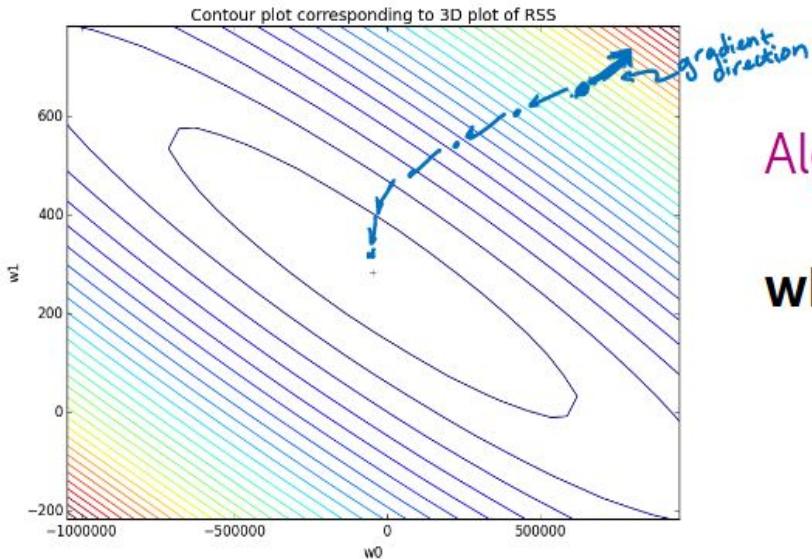
$(p+1)$ -dimensional vector

partial derivative is like a derivate with respect to w_i , treating all other variables as constants

Contour plots



Gradient descent

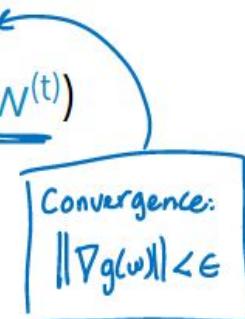


Algorithm:

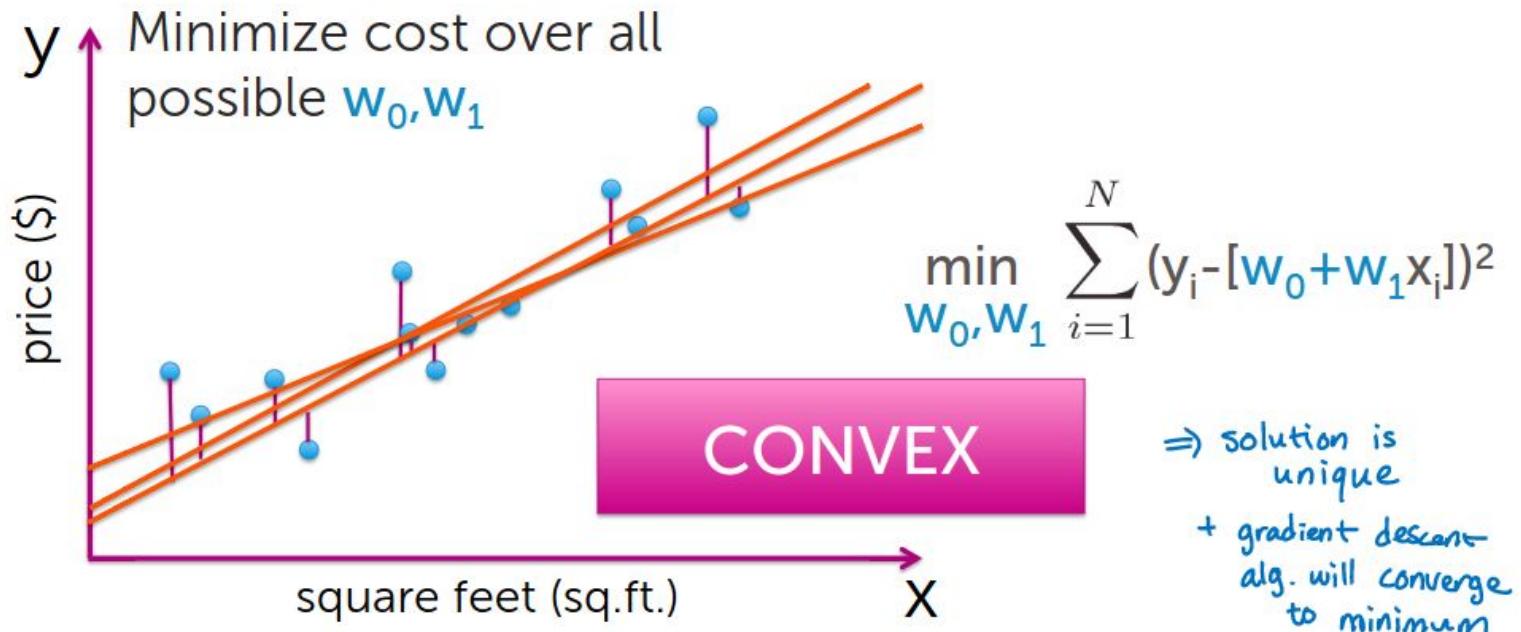
while not converged

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

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



Find “best” line



Compute the gradient

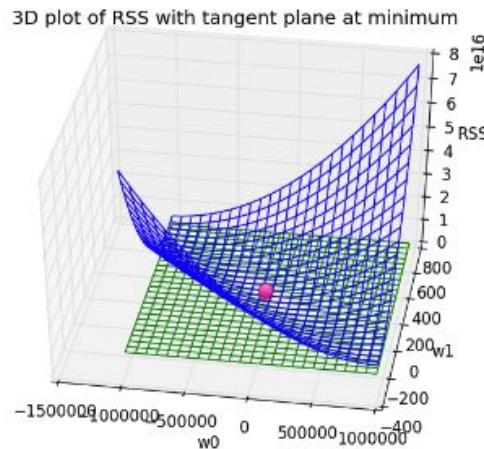
$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

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

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



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 house price
estimate of the slope
average sq.ft.

bottom term:

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

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

Note:

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

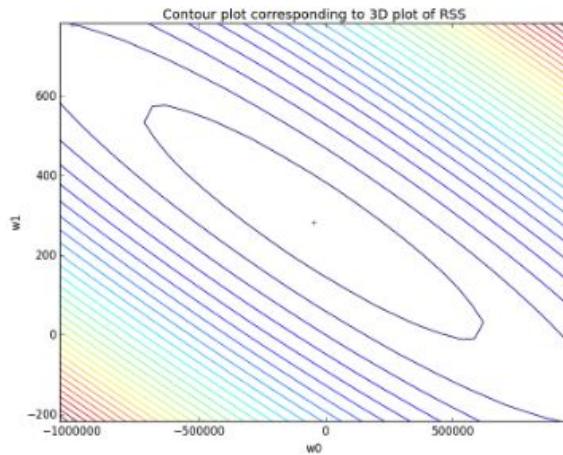
$$\sum_{i=1}^N x_i$$

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

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

Approach 2: Gradient descent

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



while not converged

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

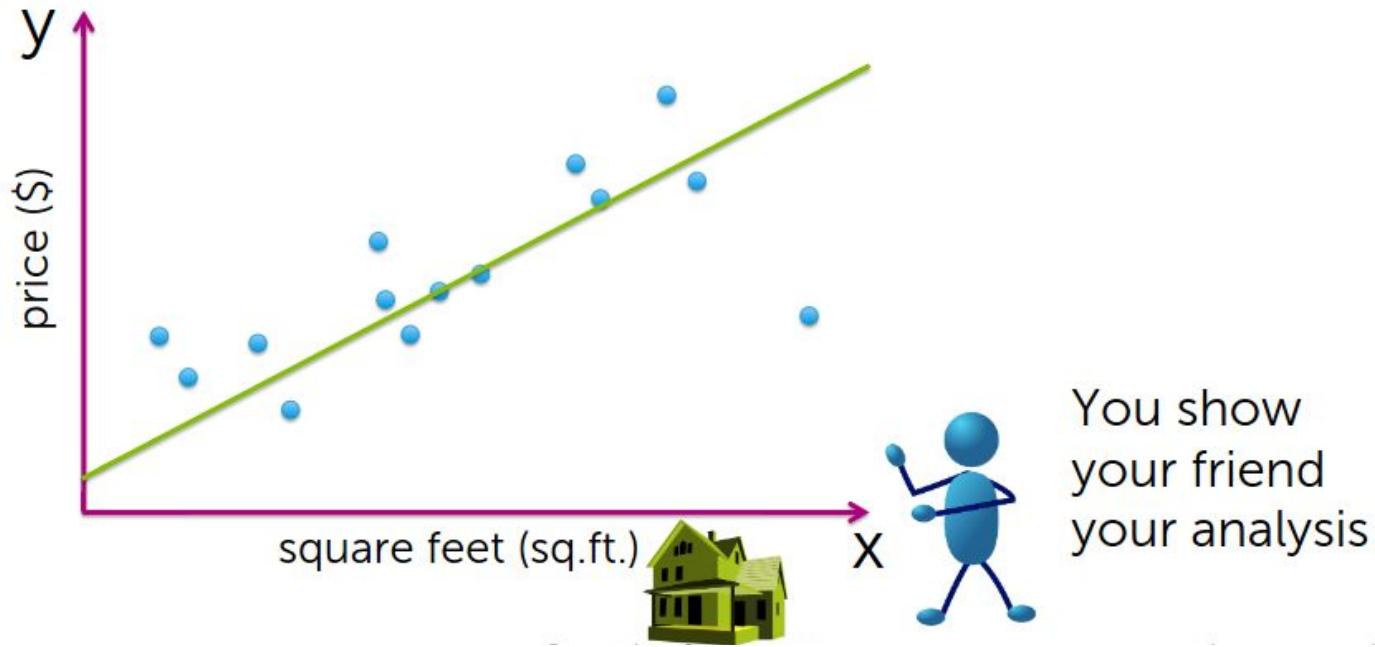
If overall, underpredicting \hat{y}_i , then $\sum [y_i - \hat{y}_i]$ is positive
→ w_0 is going to increase
similar intuition for w_1 , but multiply by x_i

Comparing the approaches

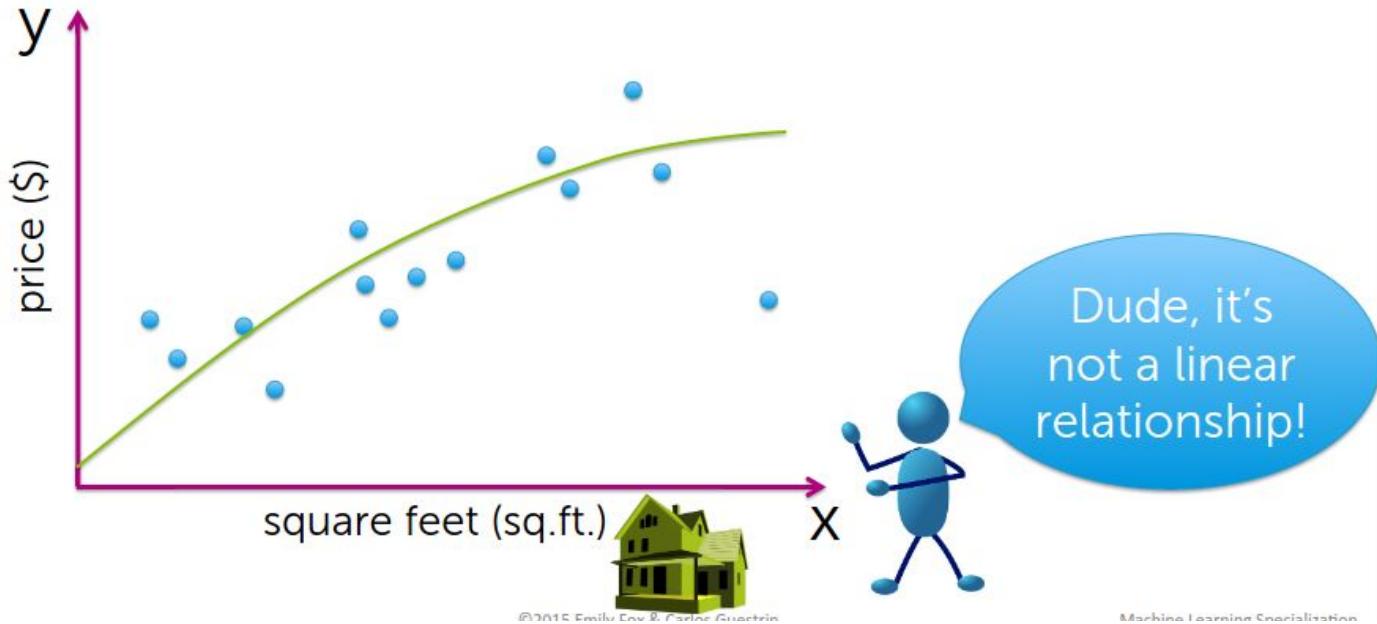
- For most ML problems,
cannot solve gradient = 0
- Even if solving gradient = 0
is feasible, gradient descent
can be more efficient
- Gradient descent relies on
choosing stepsize and
convergence criteria

Influence of high leverage points

Fit data with a line or ... ?



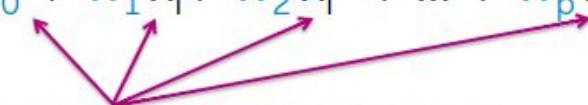
What about a quadratic function?



Polynomial regression

Model:

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$



feature 1 = 1 (constant) parameter 1 = w_0

feature 2 = x parameter 2 = w_1

feature 3 = x^2 parameter 3 = w_2

...

...

feature $p+1 = x^p$ parameter $p+1 = w_p$

Generic basis expansion

Model:

$$\begin{aligned}y_i &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i \\&= \sum_{j=0}^D w_j h_j(x_i) + \epsilon_i\end{aligned}$$

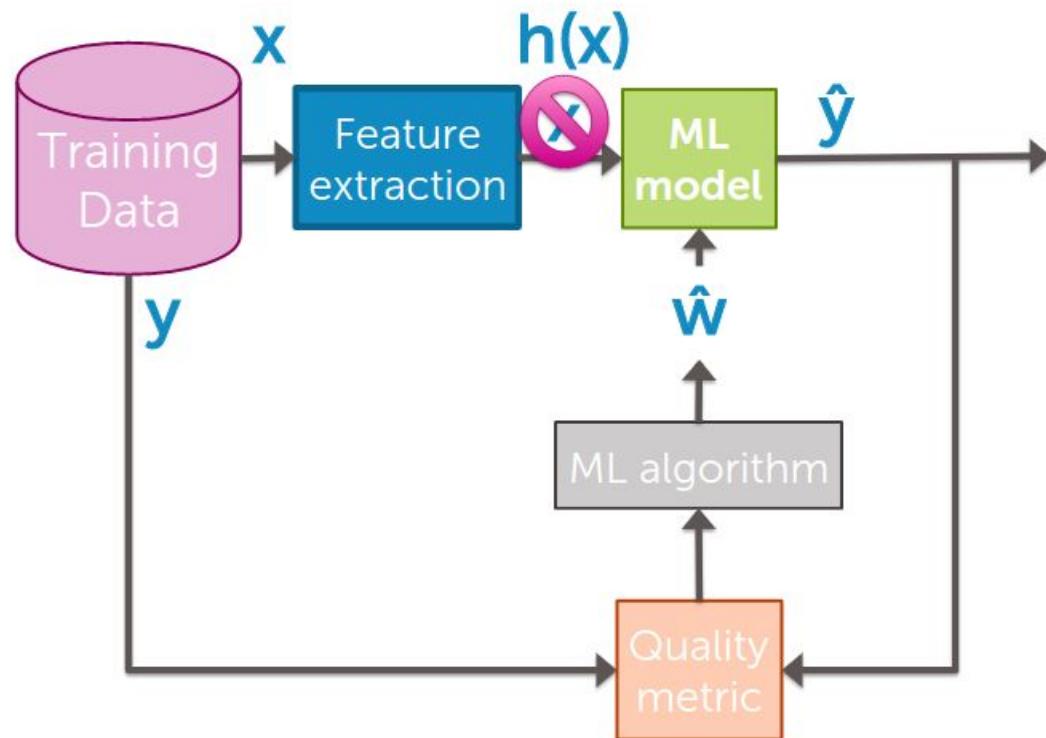
feature 1 = $h_0(x)$...often 1 (constant)

feature 2 = $h_1(x)$... e.g., x

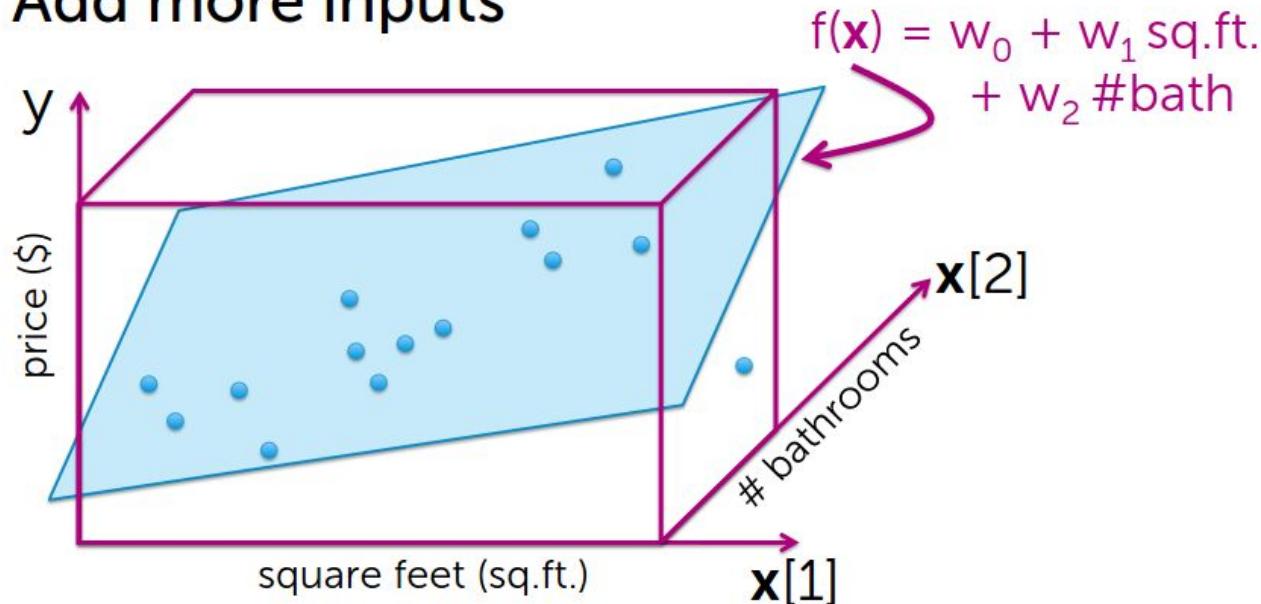
feature 3 = $h_2(x)$... e.g., x^2 or $\sin(2\pi x/12)$

...

feature $D+1 = h_D(x)$... e.g., x^p



Add more inputs



General notation

Output: $y \leftarrow$ scalar

Inputs: $\mathbf{x} = (\mathbf{x}[1], \mathbf{x}[2], \dots, \mathbf{x}[d])$
 \uparrow
d-dim vector

Notational conventions:

$\mathbf{x}[j] = j^{\text{th}}$ input (scalar)

$h_j(\mathbf{x}) = j^{\text{th}}$ feature (scalar)

$\mathbf{x}_i =$ input of i^{th} data point (vector)

$\mathbf{x}_i[j] = j^{\text{th}}$ input of i^{th} data point (scalar)

More generically... D-dimensional curve

Model:

$$\begin{aligned}y_i &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \varepsilon_i \\&= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i\end{aligned}$$

feature 1 = $h_0(\mathbf{x})$... e.g., 1

feature 2 = $h_1(\mathbf{x})$... e.g., $\mathbf{x}[1]$ = sq. ft.

feature 3 = $h_2(\mathbf{x})$... e.g., $\mathbf{x}[2]$ = #bath
or, $\log(\mathbf{x}[7]) \mathbf{x}[2] = \log(\#bed) \times \#bath$

...

feature $D+1 = h_D(\mathbf{x})$... some other function of $\mathbf{x}[1], \dots, \mathbf{x}[d]$

More on notation

observations (\mathbf{x}_i, y_i) : N

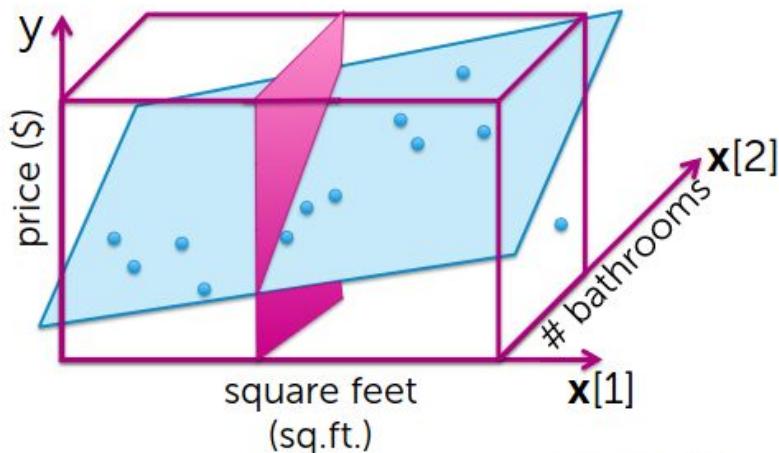
inputs $\mathbf{x}[j]$: d

features $h_j(\mathbf{x})$: D

Interpreting the coefficients – Two linear features

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x[1] + \hat{w}_2 x[2]$$

fix



Rewrite in matrix notation

For observation i

$$y_i = \sum_{j=0}^D w_j h_j(x_i) + \varepsilon_i$$

$$\begin{aligned} y_i &= \underbrace{\begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_D \end{bmatrix}}_{w^T} \underbrace{h(x_i)}_{\begin{bmatrix} h_0(x_i) \\ h_1(x_i) \\ h_2(x_i) \\ \vdots \\ h_D(x_i) \end{bmatrix}} = \begin{bmatrix} h^T(x_i) \end{bmatrix} \underbrace{\begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_D \end{bmatrix}}_w + \varepsilon_i \\ &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \varepsilon_i \\ &\quad \text{simpler} \\ &= w^T h(x_i) + \varepsilon_i \end{aligned}$$

Rewrite in matrix notation

For all observations together

$$\begin{matrix} \mathbf{y} \\ y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{matrix} = \mathbf{H} \begin{matrix} \mathbf{w} \\ w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{matrix} + \begin{matrix} \mathbf{\epsilon} \\ \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_N \end{matrix}$$

where \mathbf{H} is the matrix of features, \mathbf{w} is the vector of weights, and $\mathbf{\epsilon}$ is the vector of errors.

$$\Rightarrow \boxed{\mathbf{y} = \mathbf{H} \mathbf{w} + \mathbf{\epsilon}}$$

RSS in matrix notation

$$\begin{aligned}\text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2 \\ &= (\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

Why? (part 2)

residual ₁	residual ₂	residual ₃	...	residual _N	residual ₁	residual ₂	residual ₃	...	residual _N

$$\begin{aligned}& (\text{residual}_1^2 + \text{residual}_2^2 + \dots + \text{residual}_N^2) \\ &= \sum_{i=1}^N \text{residual}_i^2 \\ &\triangleq \text{RSS}(\mathbf{w})\end{aligned}$$

Gradient of RSS

$$\begin{aligned}\nabla \text{RSS}(\mathbf{w}) &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})] \\ &= -2\mathbf{H}^\top (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

Why? By analogy to 1D case:

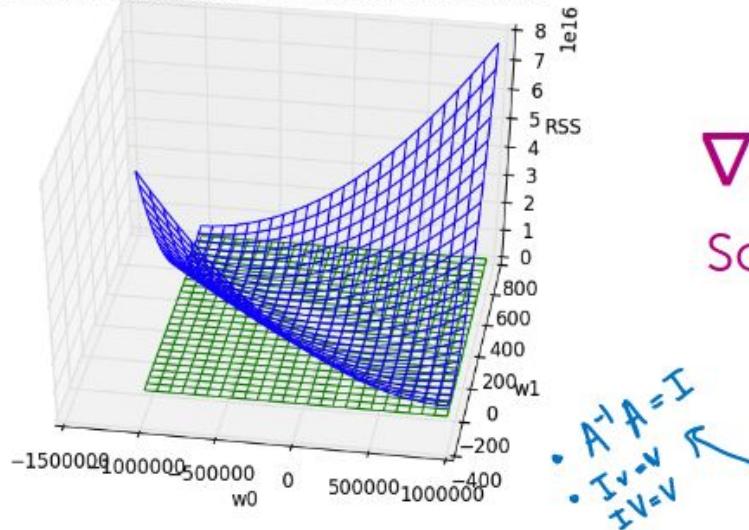
$$\frac{d}{dw} (y-hw)(y-hw) = \frac{d}{dw} (y-hw)^2 = 2 \cdot (y-hw)' (-h)$$

$\uparrow \uparrow$
scalars

$$= -2h(y-hw)$$

Closed-form solution

3D plot of RSS with tangent plane at minimum



$$\begin{aligned} & \bullet A^T A = I \\ & \bullet I^{-1} = V \\ & I = V \\ & \bullet \end{aligned}$$

$$\nabla \text{RSS}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) = 0$$

Solve for \mathbf{w} :

$$-\cancel{2\mathbf{H}^T y} + \cancel{2\mathbf{H}^T \mathbf{H} \hat{\mathbf{w}}} = 0$$

$$\mathbf{H}^T \mathbf{H} \hat{\mathbf{w}} = \mathbf{H}^T \mathbf{y}$$

$$\underbrace{(\mathbf{H}^T \mathbf{H})^{-1}}_I \mathbf{H}^T \hat{\mathbf{w}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$

$$\hat{\mathbf{w}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$

Closed-form solution

$$\hat{\mathbf{w}} = (\underbrace{\mathbf{H}^T \mathbf{H}}_{\text{# features } = D})^{-1} \mathbf{H}^T \mathbf{y}$$

The diagram illustrates the dimensions of the matrices involved in the closed-form solution. On the left, a vertical vector y is shown with a dimension of N indicated by a bracket below it. An arrow points from this vector to a matrix H , which is labeled with $\# \text{features} = D$ above it. The matrix H is shown with a dimension of $N \times D$ indicated by brackets. The product $H^T y$ results in a vector of size $D \times 1$. This vector is then multiplied by the inverse of the matrix $H^T H$, which is a square matrix of size $D \times D$ labeled with $\# \text{features}$ above it. The final result is a vector $\hat{\mathbf{w}}$ of size $D \times 1$.

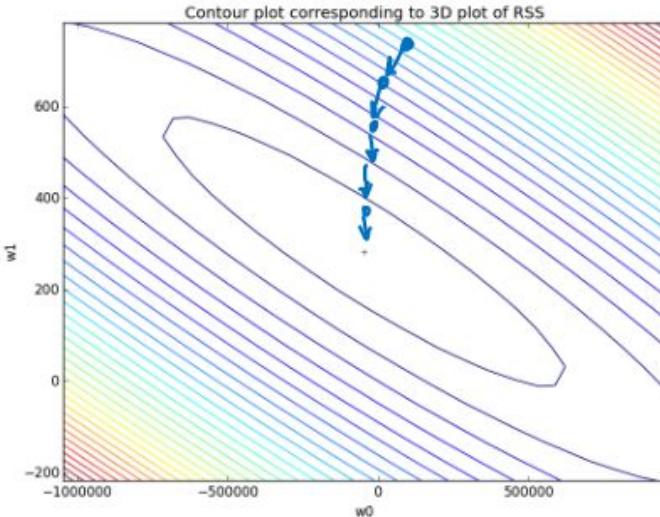
Invertible if:
In most cases is $N > D$

really,
 $\#$ of linearly
ind. observations

Complexity of inverse:

$$O(D^3)$$

Gradient descent



while not converged

$$\begin{aligned} \mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}^{(t)}) \\ &\quad - 2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) \\ &\leftarrow \mathbf{w}^{(t)} + 2\eta \mathbf{H}^T \underbrace{(\mathbf{y} - \mathbf{H}\mathbf{w}^{(t)})}_{\hat{\mathbf{y}}(\mathbf{w}^{(t)})} \end{aligned}$$

Feature-by-feature update

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^\top \mathbf{w})^2$$
$$= \sum_{i=1}^N (y_i - w_0 h_0(x_i) - w_1 h_1(x_i) - \dots - w_D h_D(x_i))^2$$

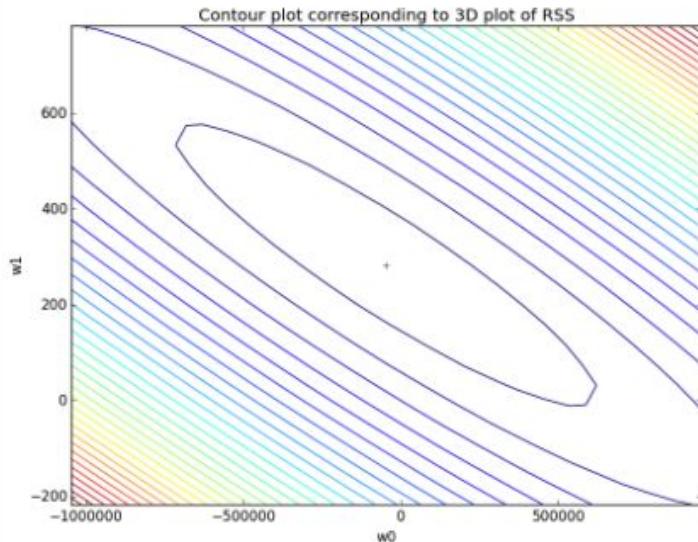
Partial with respect to w_j

$$\begin{aligned} & \sum_{i=1}^N 2(y_i - w_0 h_0(x_i) - w_1 h_1(x_i) - \dots - w_D h_D(x_i)) \\ & \quad \cdot (-\underline{h_j(x_i)}) \end{aligned}$$
$$= -2 \sum_{i=1}^N h_j(x_i) (y_i - h(\mathbf{x}_i)^\top \mathbf{w})$$

Update to j^{th} feature weight:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \left(-2 \sum_{i=1}^N h_j(x_i) (y_i - \underbrace{h^\top(\mathbf{x}_i) \mathbf{w}^{(t)}}_{\hat{y}_i(\mathbf{w}^{(t)})}) \right)$$

Summary of gradient descent for multiple regression



```
init  $\mathbf{w}^{(1)} = 0$  (or randomly, or smartly),  $t = 1$ 
while  $\|\nabla \text{RSS}(\mathbf{w}^{(t)})\| > \epsilon$  tolerance
    for  $j = 0, \dots, D$ 
        partial[j] =  $-2 \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$ 
         $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} - \eta \text{partial}[j]$ 
    t  $\leftarrow t + 1$ 
```