



Simple Regression:

Linear regression with one input

Recall Task: Predicting house prices

How much is my house worth?



How much is my house worth?



Look at recent sales in my neighborhood

- How much did they sell for?



Regression fundamentals: data, model, task

Data

input output



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



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



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



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



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

:

Data

input output



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



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



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



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



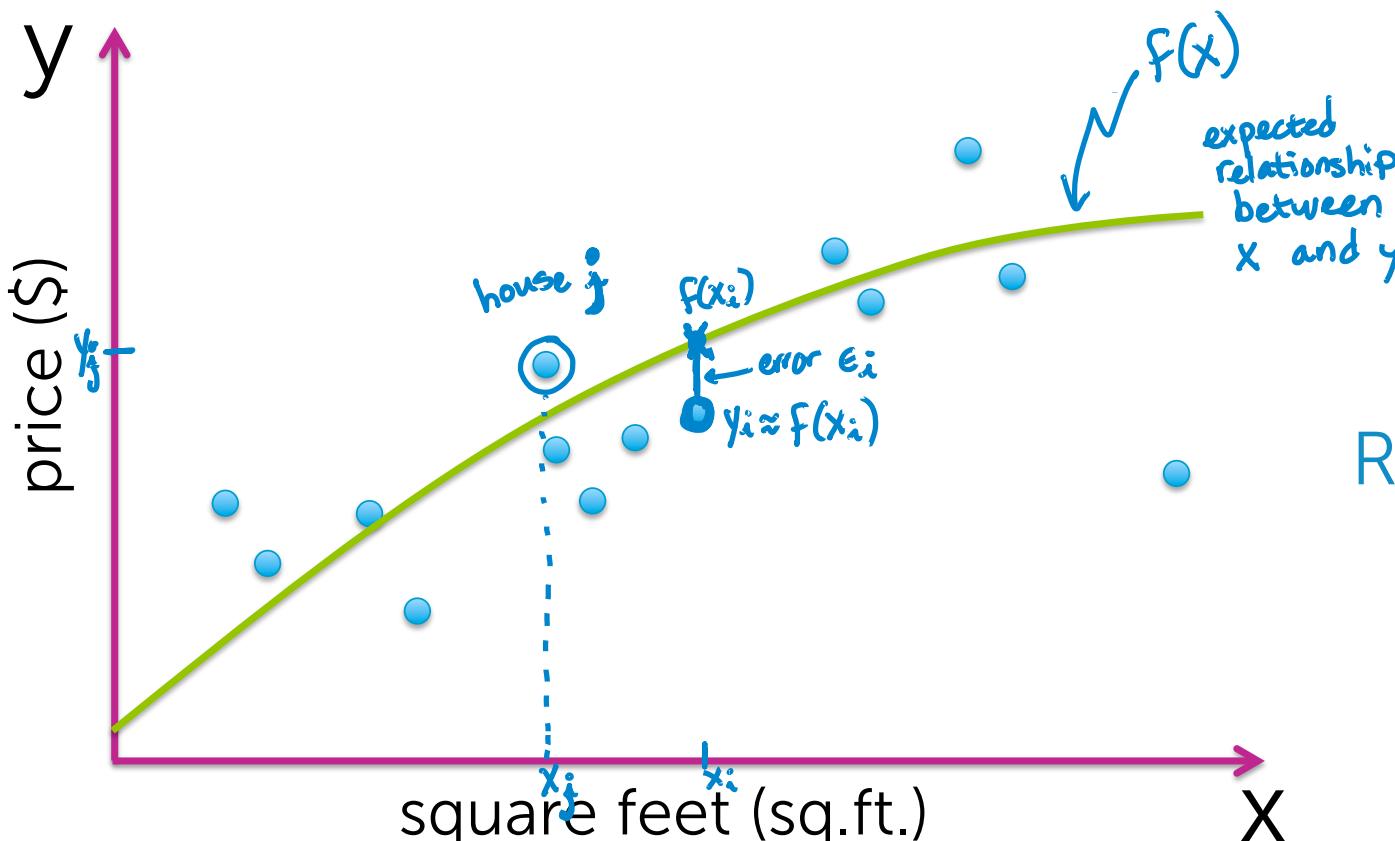
$(x_5 = \text{sq.ft.}, y_5 = \$)$

⋮

Input vs. Output:

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

Model – How we assume the world works



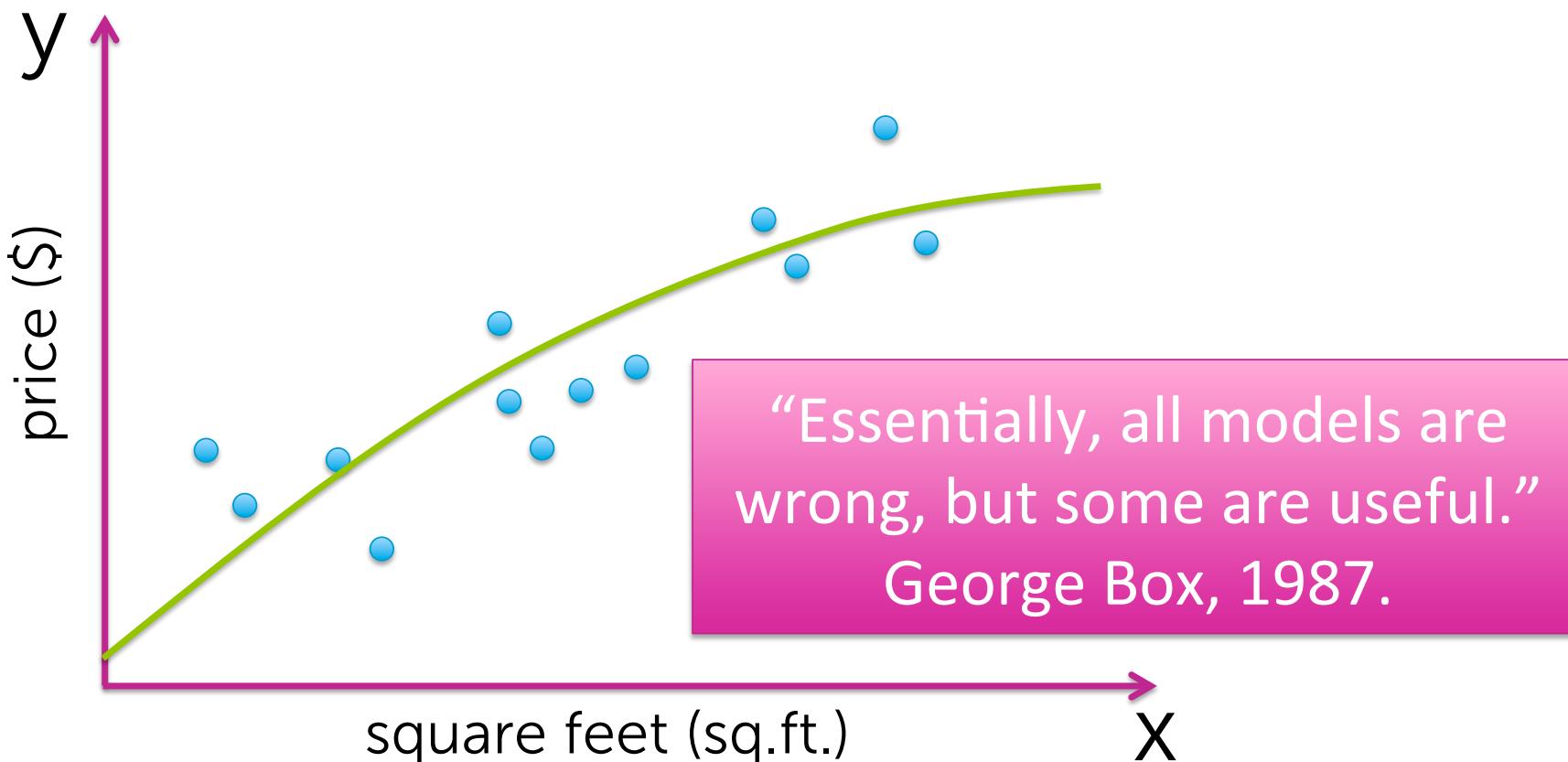
Regression model:

$$y_i = f(x_i) + e_i$$

$E[e_i] = 0$ ← equally likely
that error
is + or -
↑ expected value

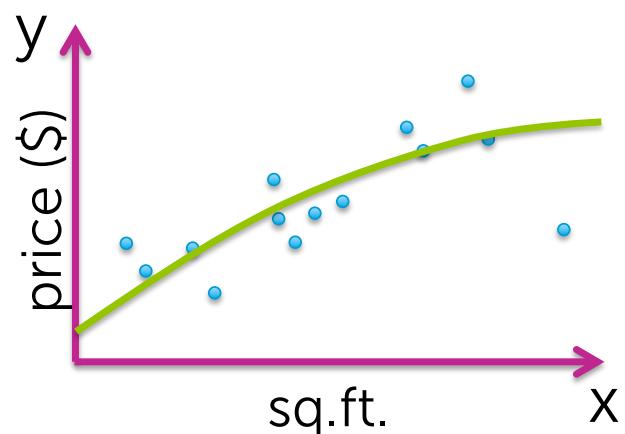
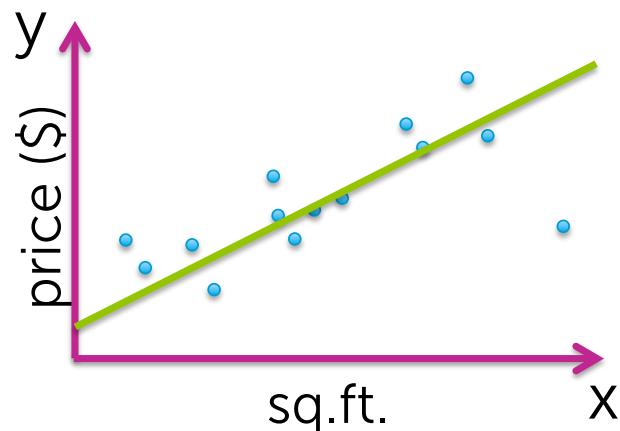
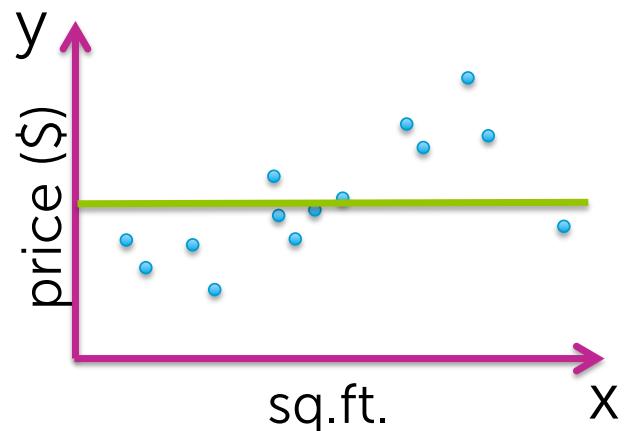
↓
 y_i is equally
likely to be above
or below $f(x_i)$

Model – How we *assume* the world works

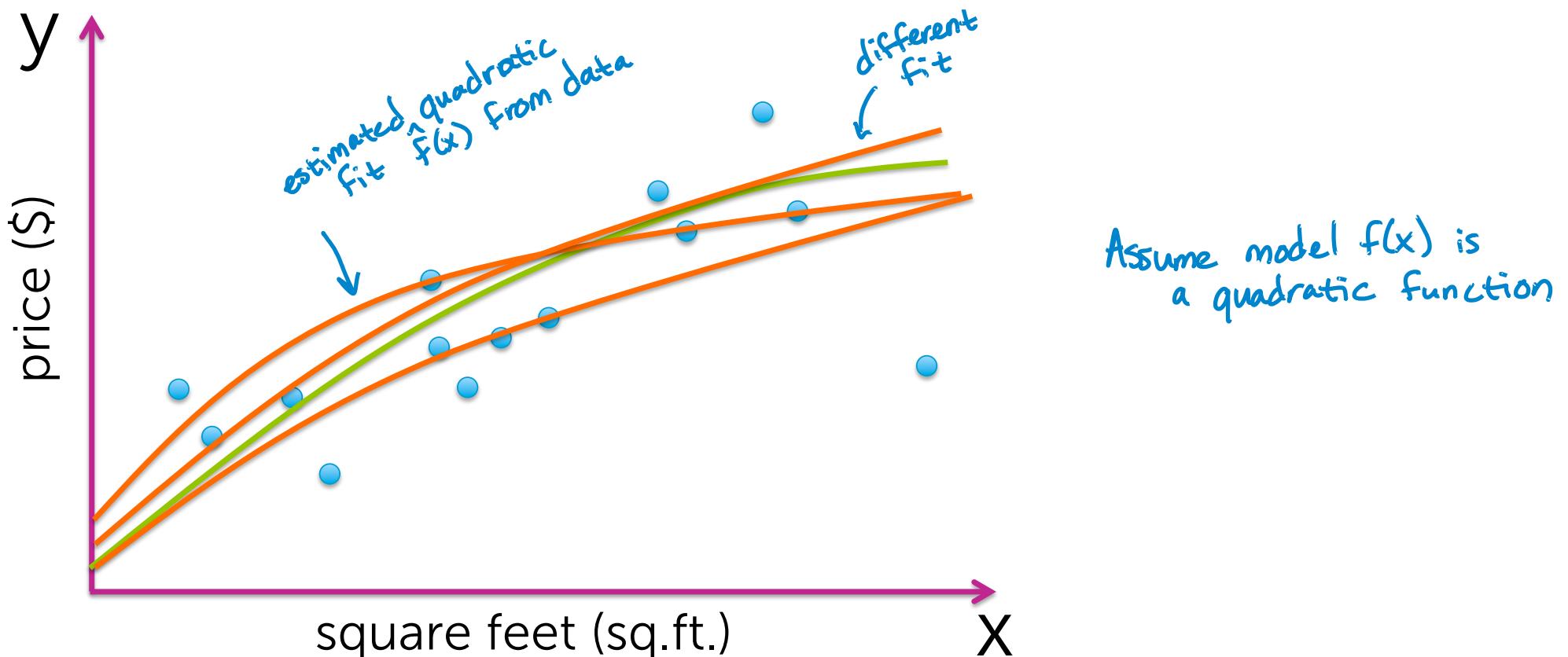


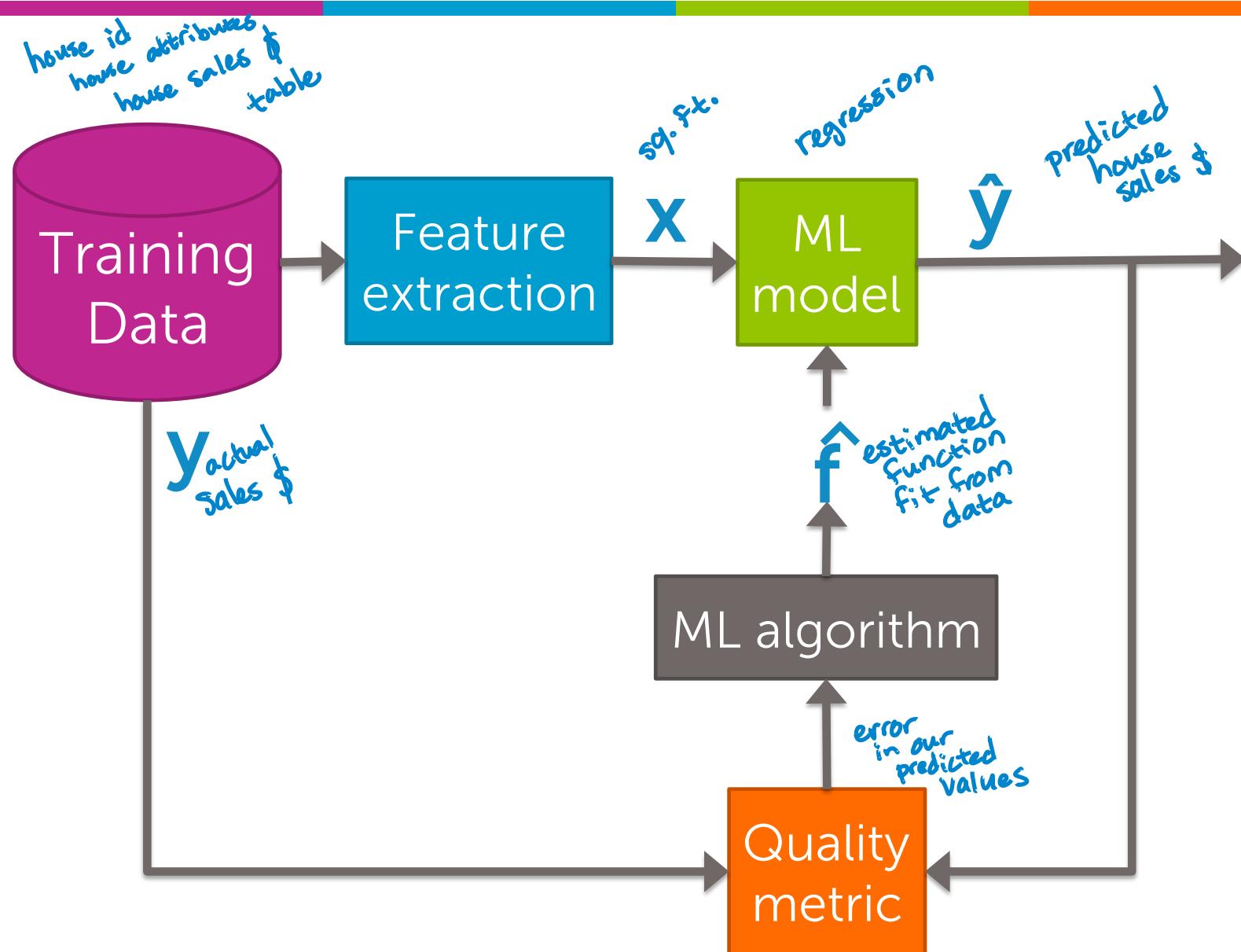
Task 1–

Which model $f(x)$?

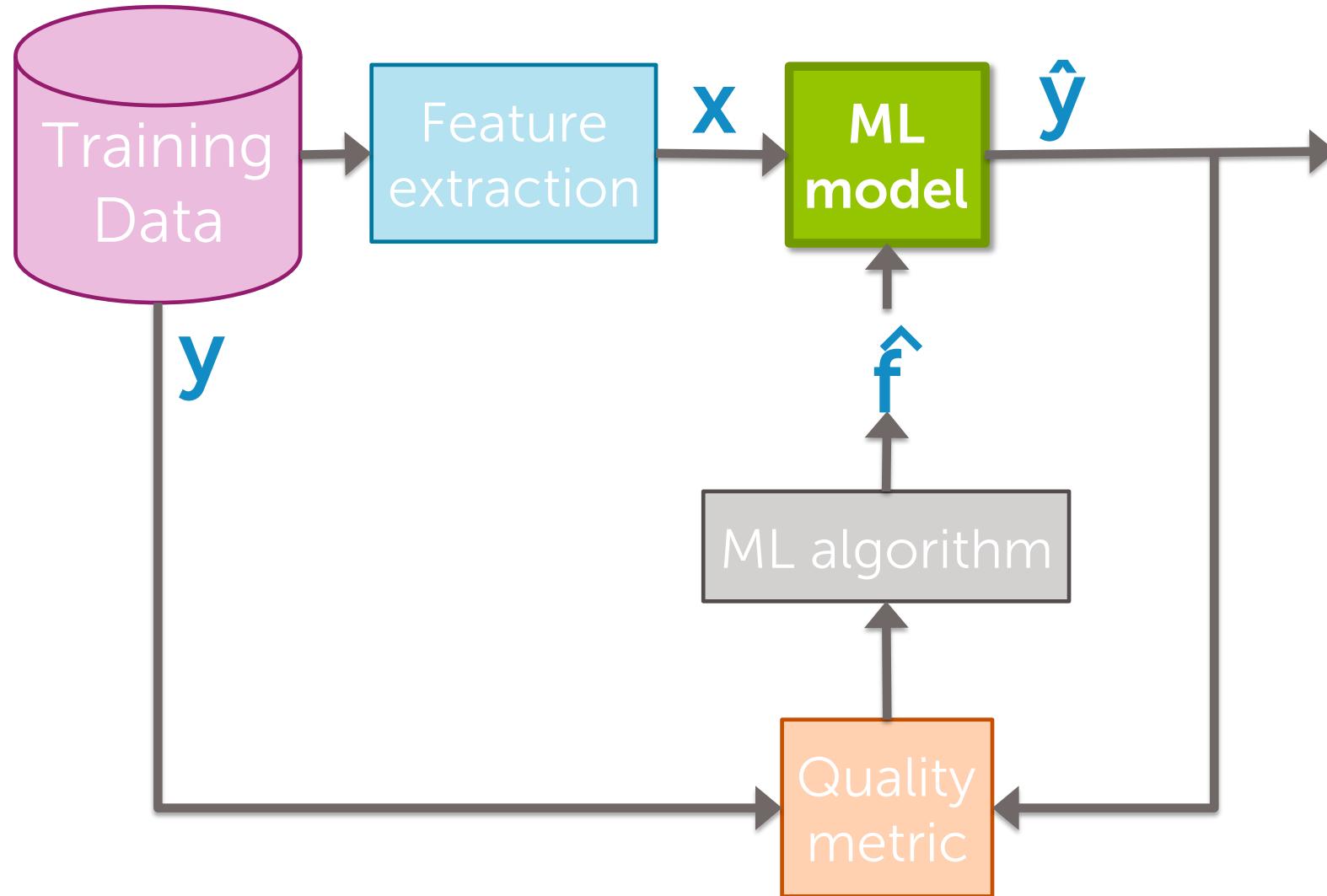


Task 2 – For a given model $f(x)$, estimate function $\hat{f}(x)$ from data

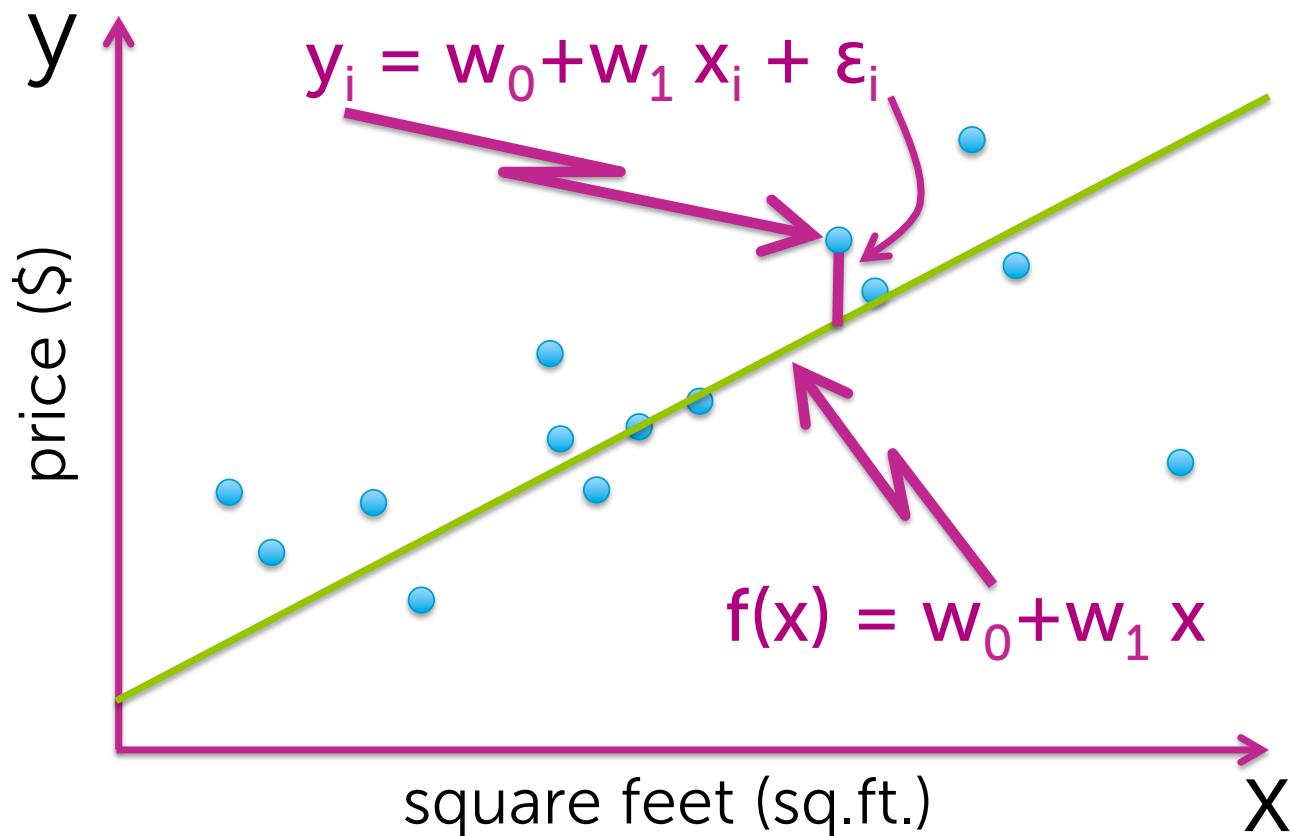




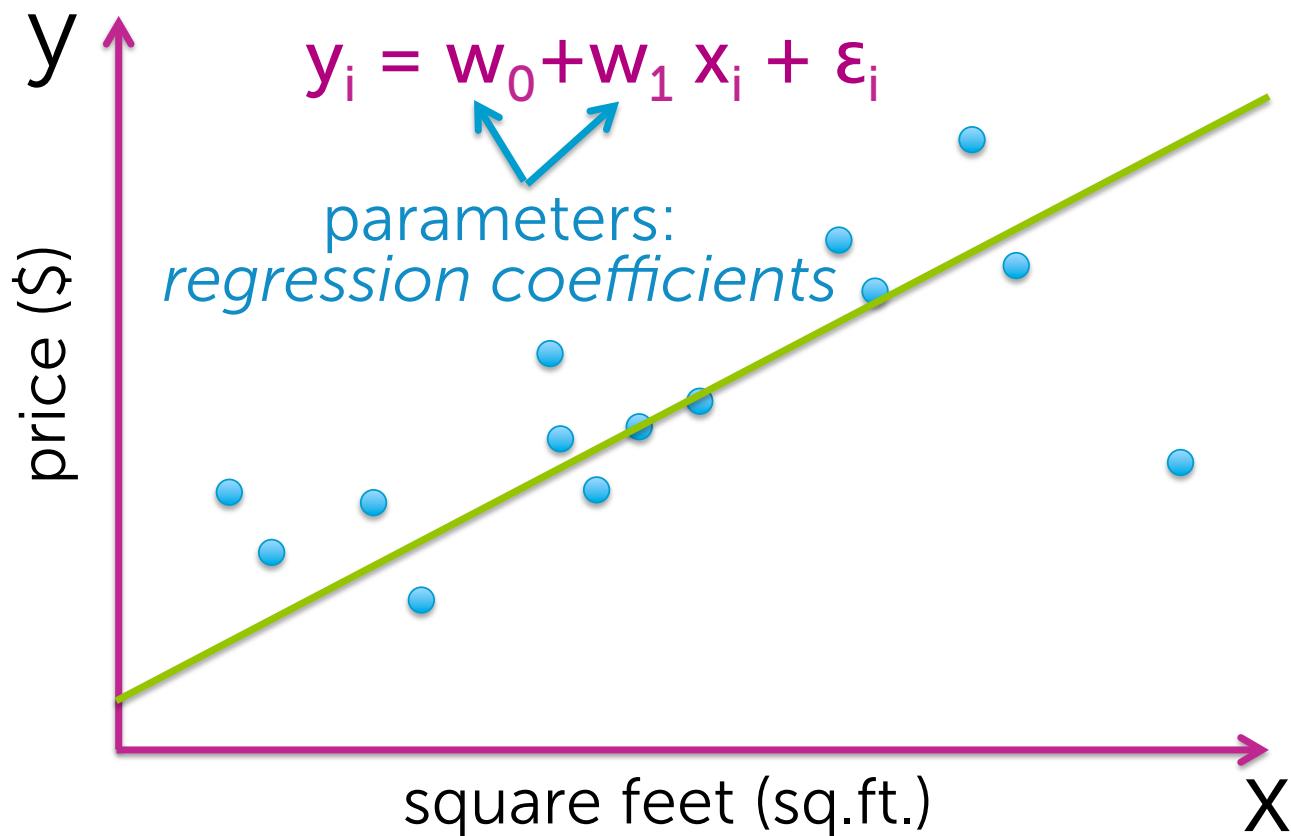
Simple linear regression



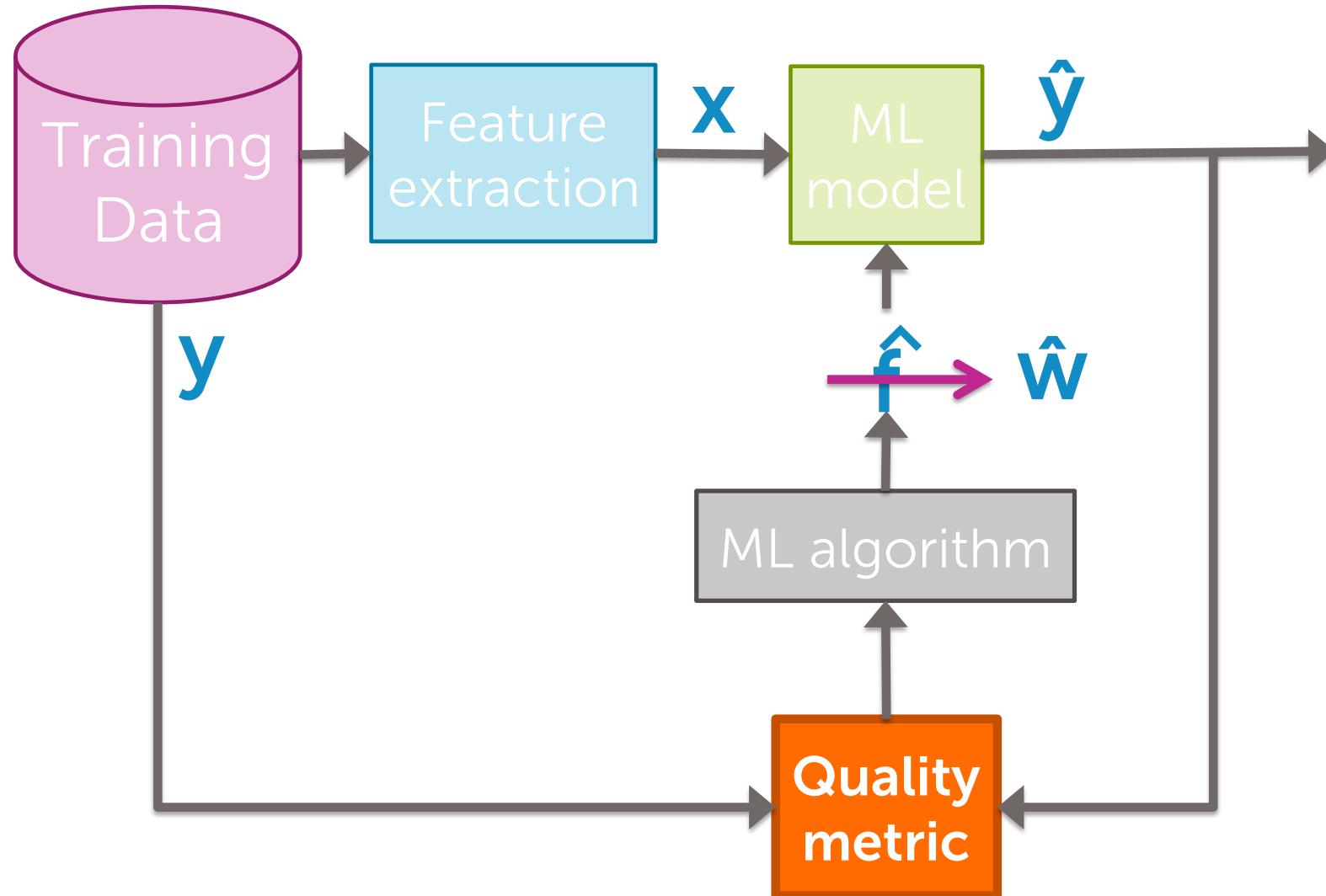
Simple linear regression model



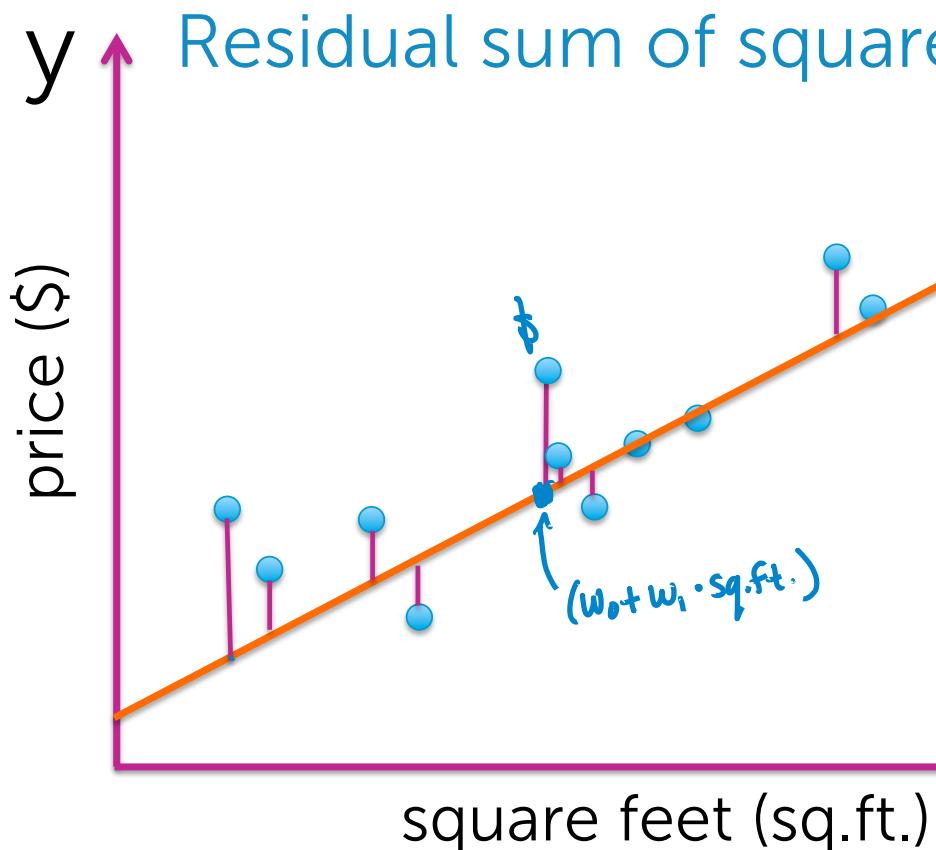
Simple linear regression model



Fitting a line to data



"Cost" of using a given line

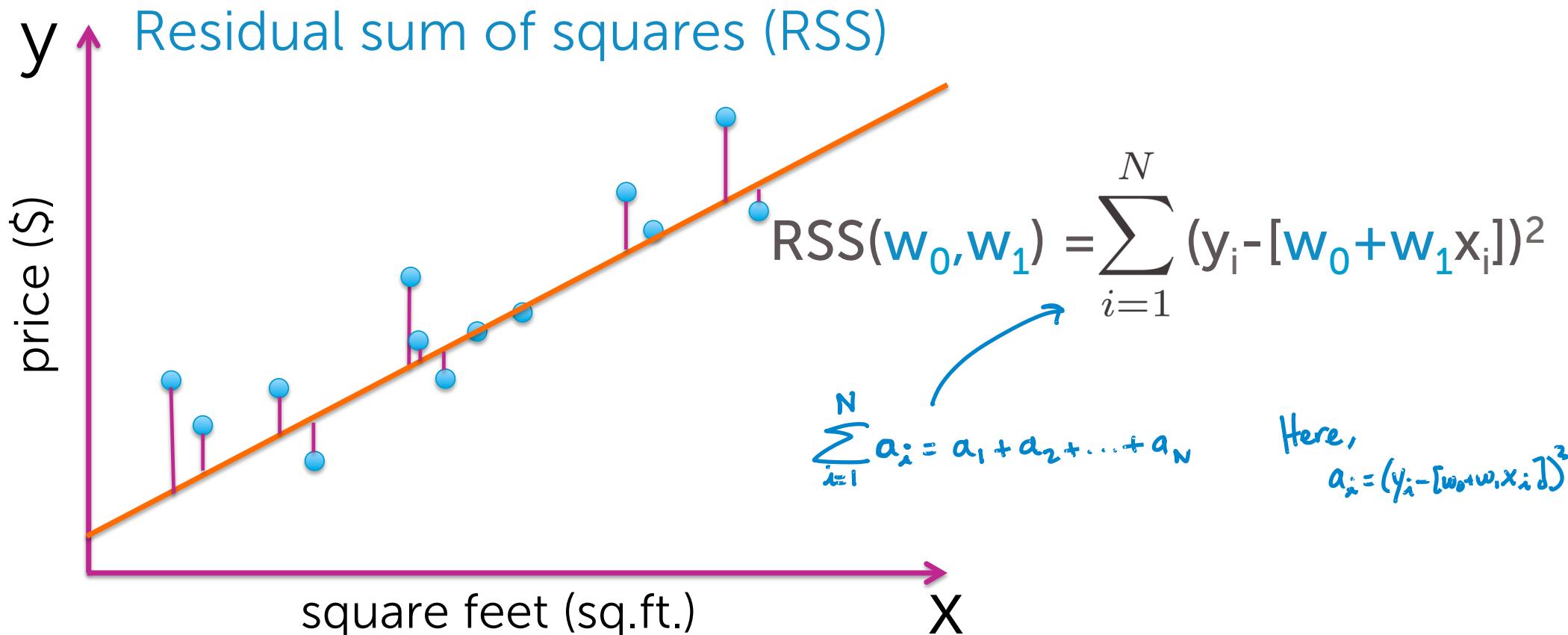


$\text{RSS}(w_0, w_1) =$

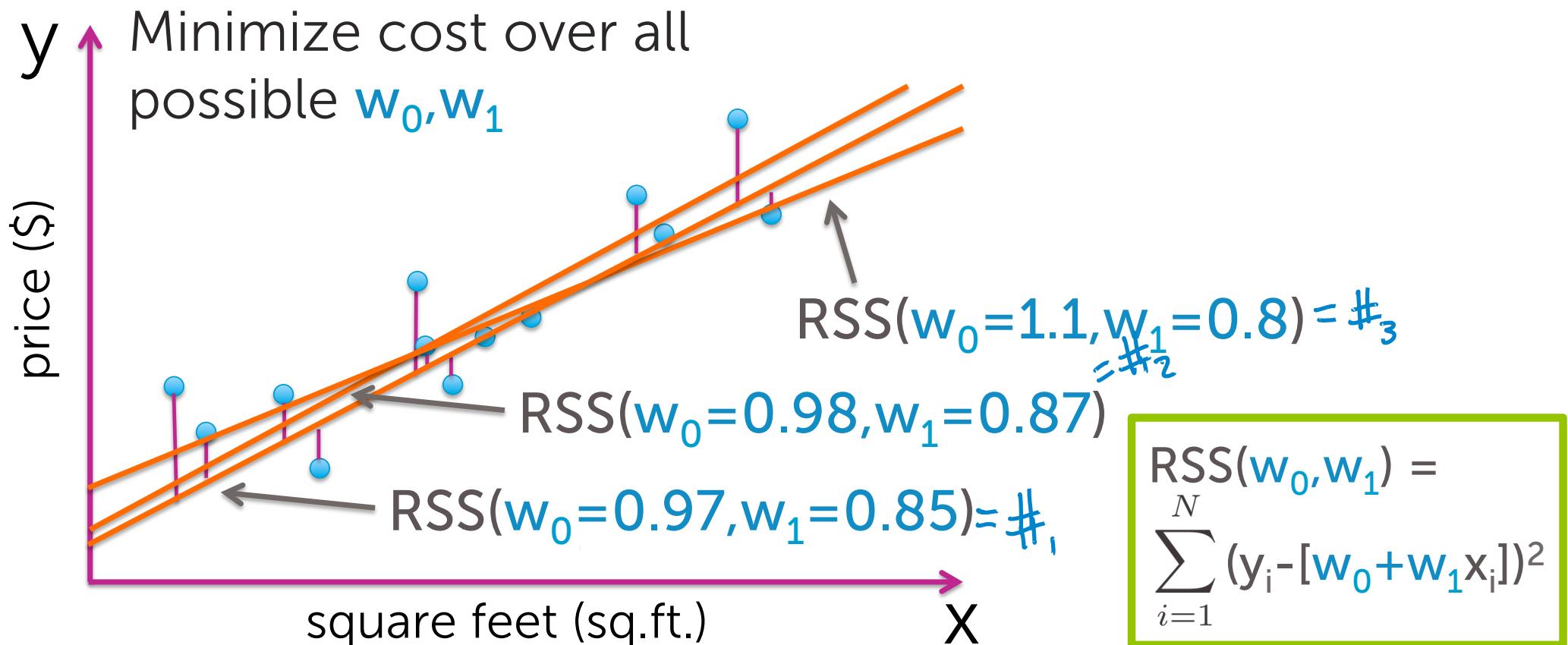
$$(\$_{\text{house 1}} - [w_0 + w_1 \text{sq.ft.}_{\text{house 1}}])^2$$
$$+ (\$_{\text{house 2}} - [w_0 + w_1 \text{sq.ft.}_{\text{house 2}}])^2$$
$$+ (\$_{\text{house 3}} - [w_0 + w_1 \text{sq.ft.}_{\text{house 3}}])^2$$

+ ...[include all training houses]

"Cost" of using a given line

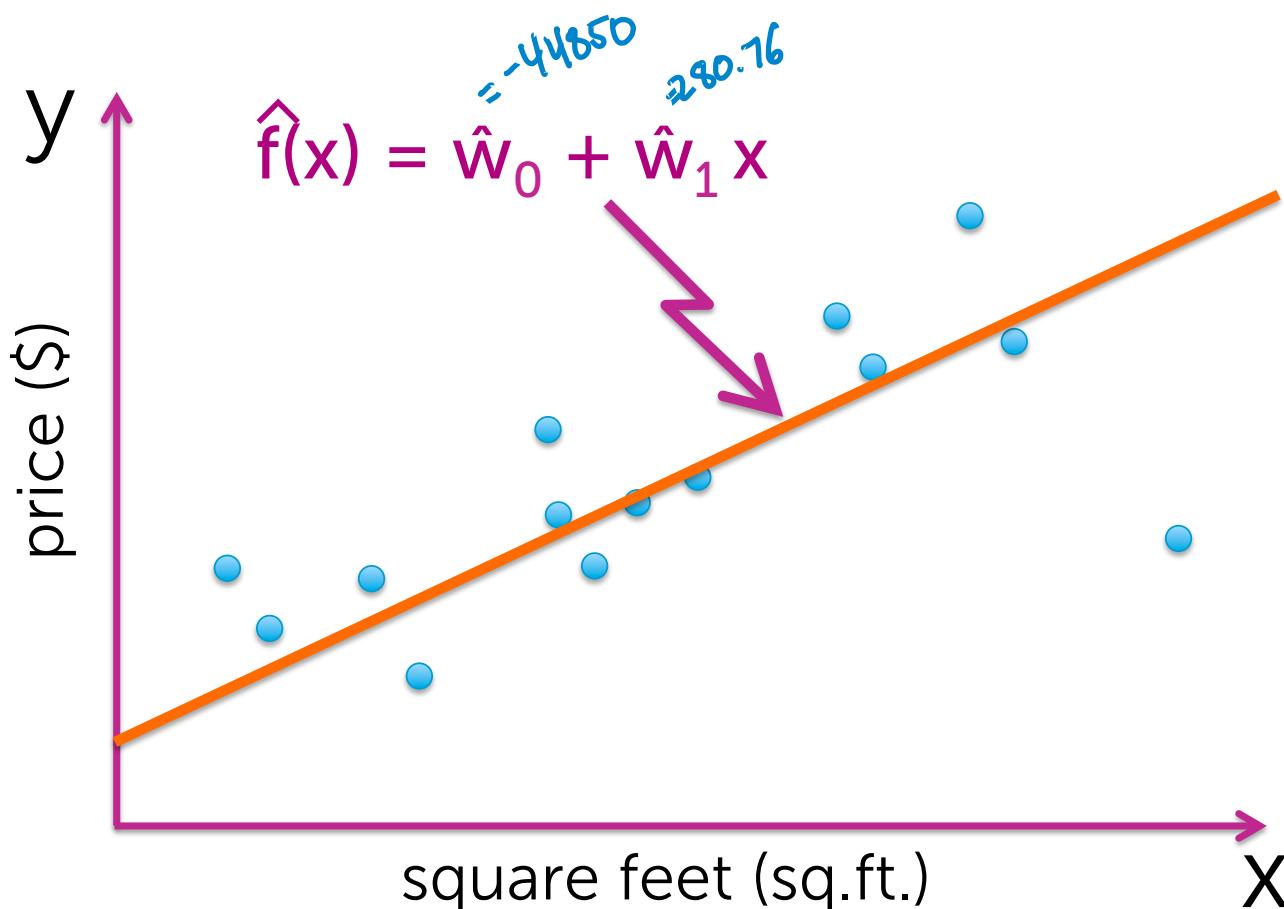


Find “best” line



The fitted line: use + interpretation

Model vs. fitted line



Regression model:

$$y_i = w_0 + w_1 x_i + \epsilon_i$$

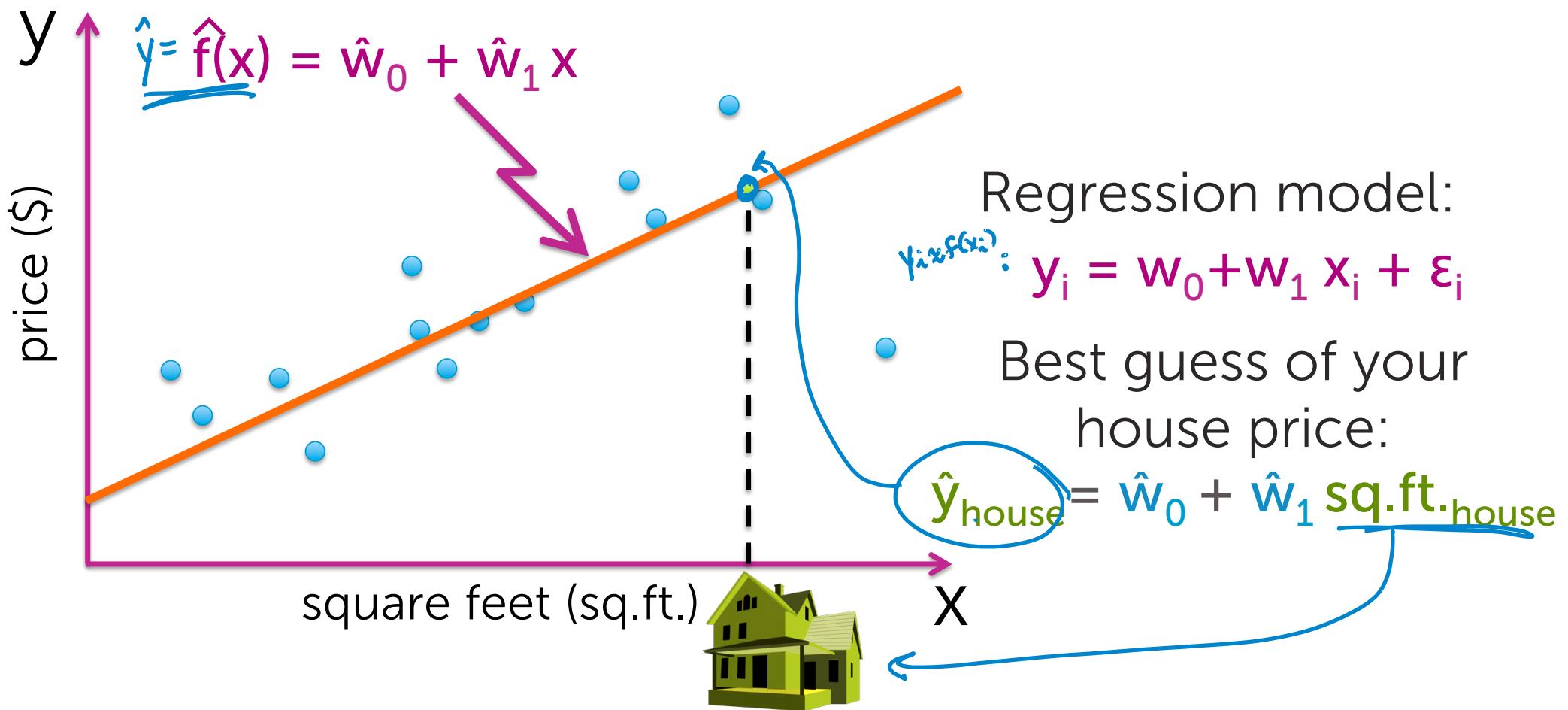
parameters (unknown variables)

Estimated parameters:

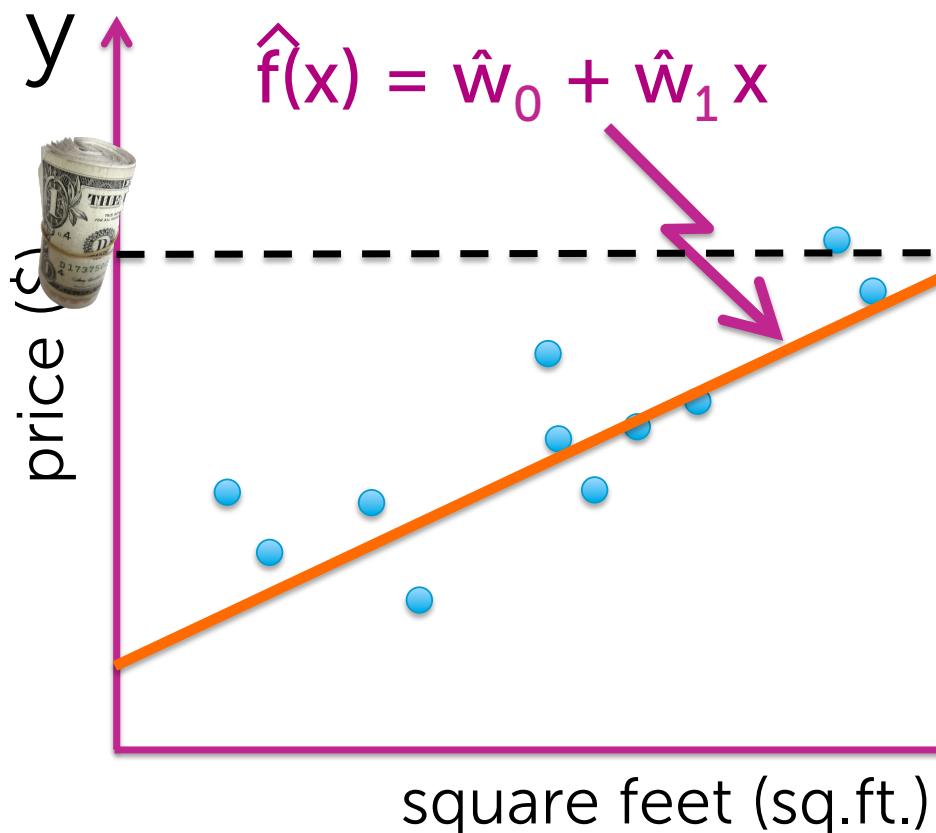
$$\hat{w}_0 = -44950, \hat{w}_1 = 280.76$$

take actual values

Seller: Predicting your house price



Buyer: Predicting size of house



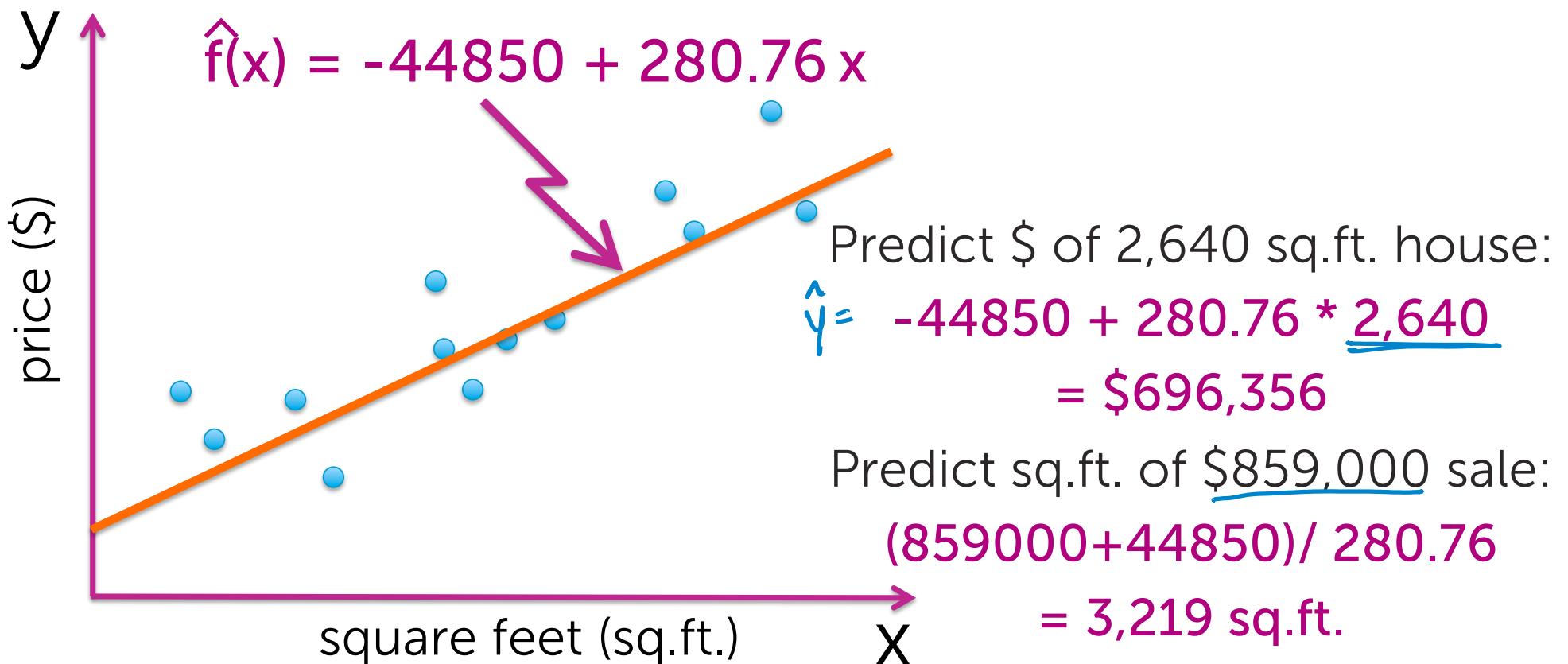
Regression model:

$$y_i = w_0 + w_1 x_i + \epsilon_i$$

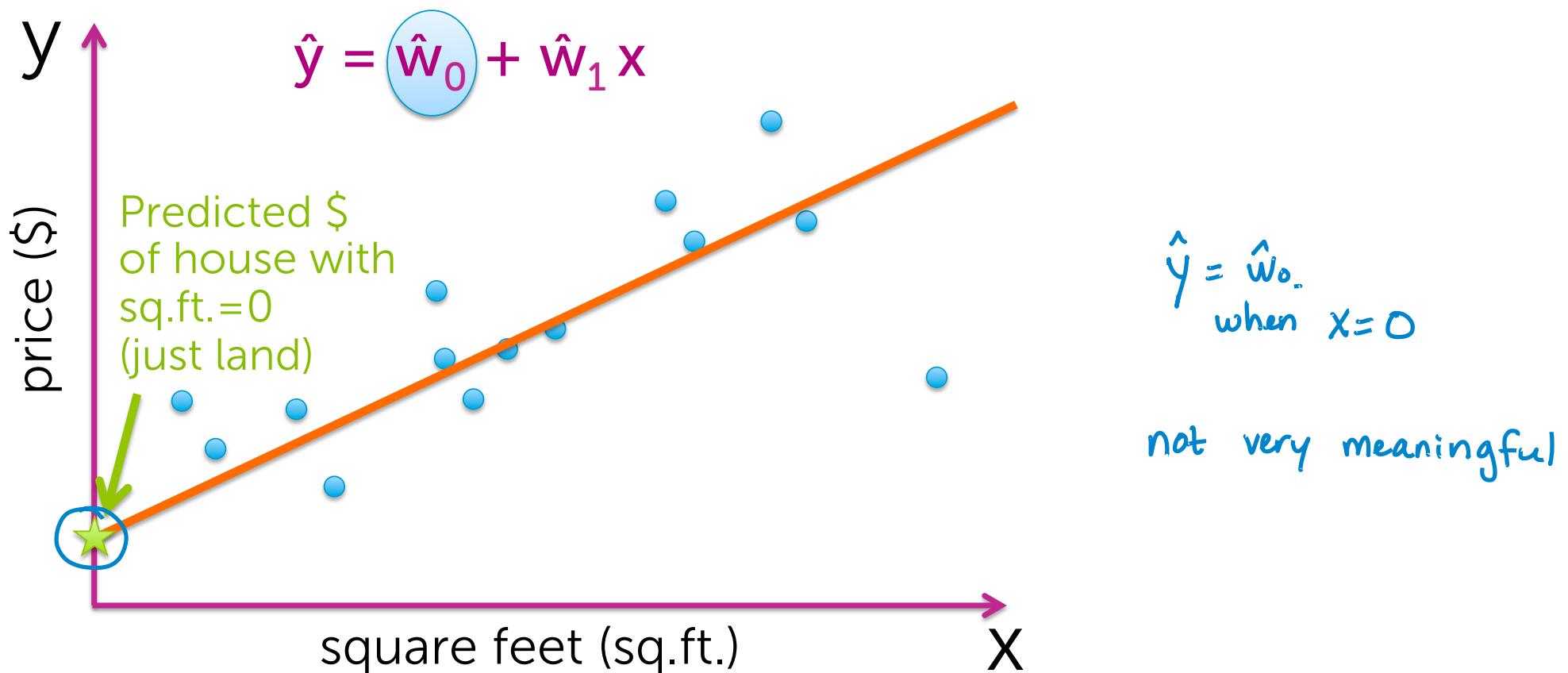
Best guess of size of house you can afford:

$$\text{\$ in bank} = \hat{w}_0 + \hat{w}_1 \text{sq.ft.}$$

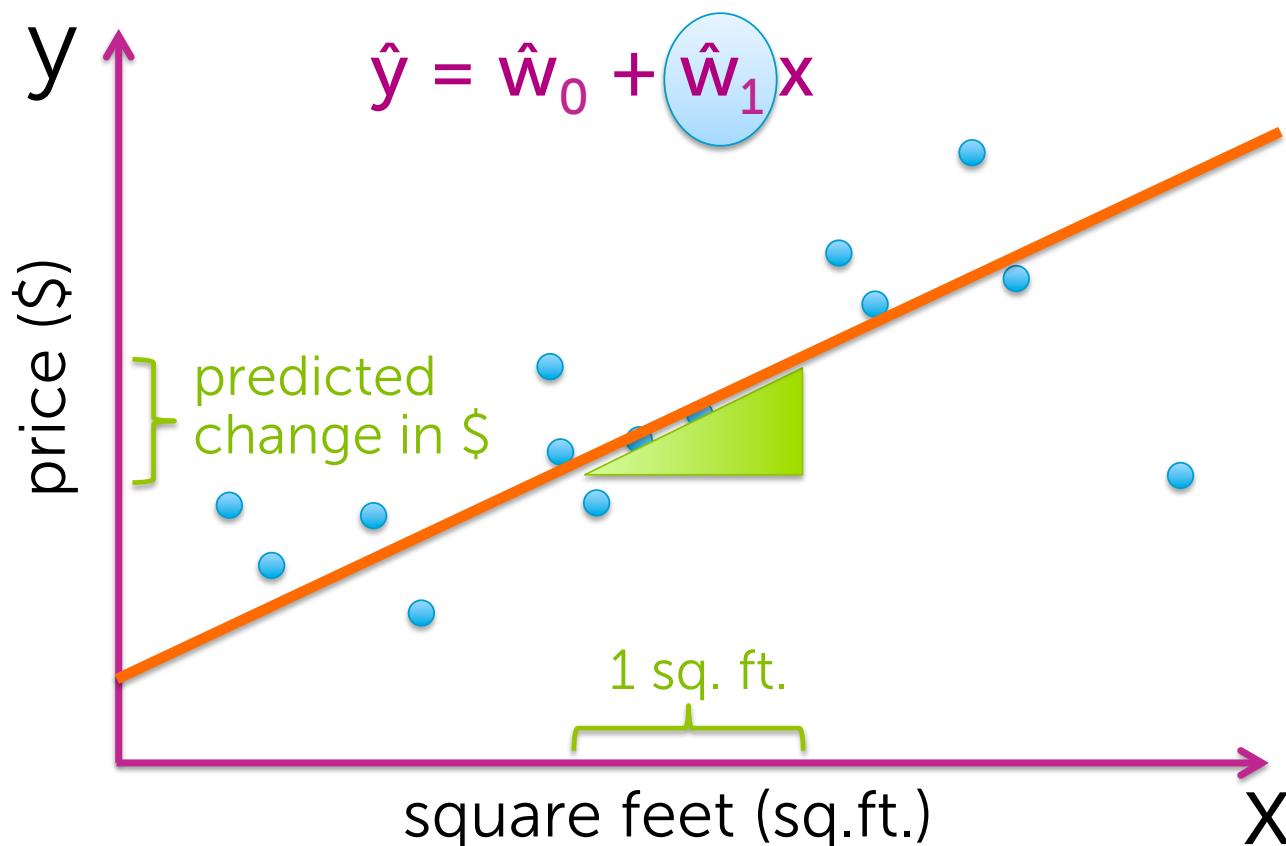
A concrete example



Interpreting the coefficients



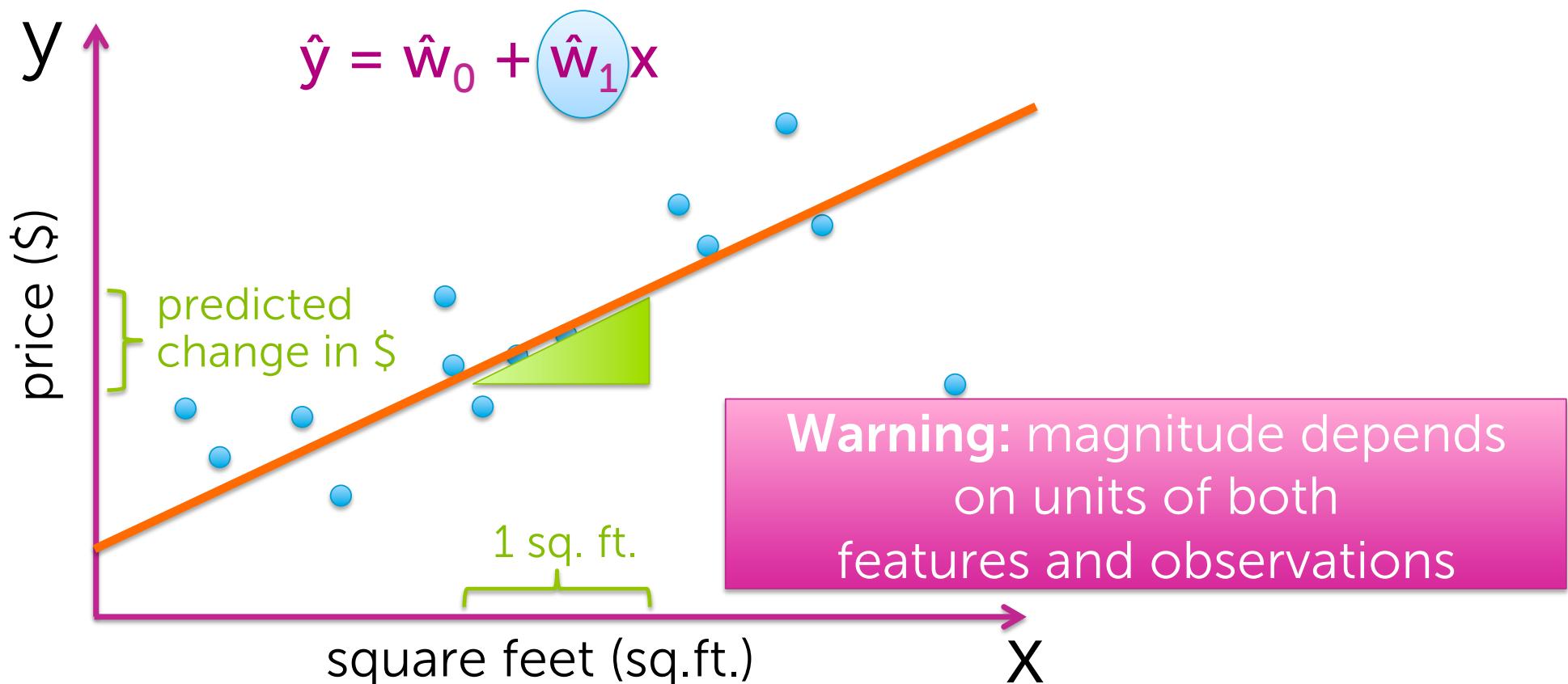
Interpreting the coefficients



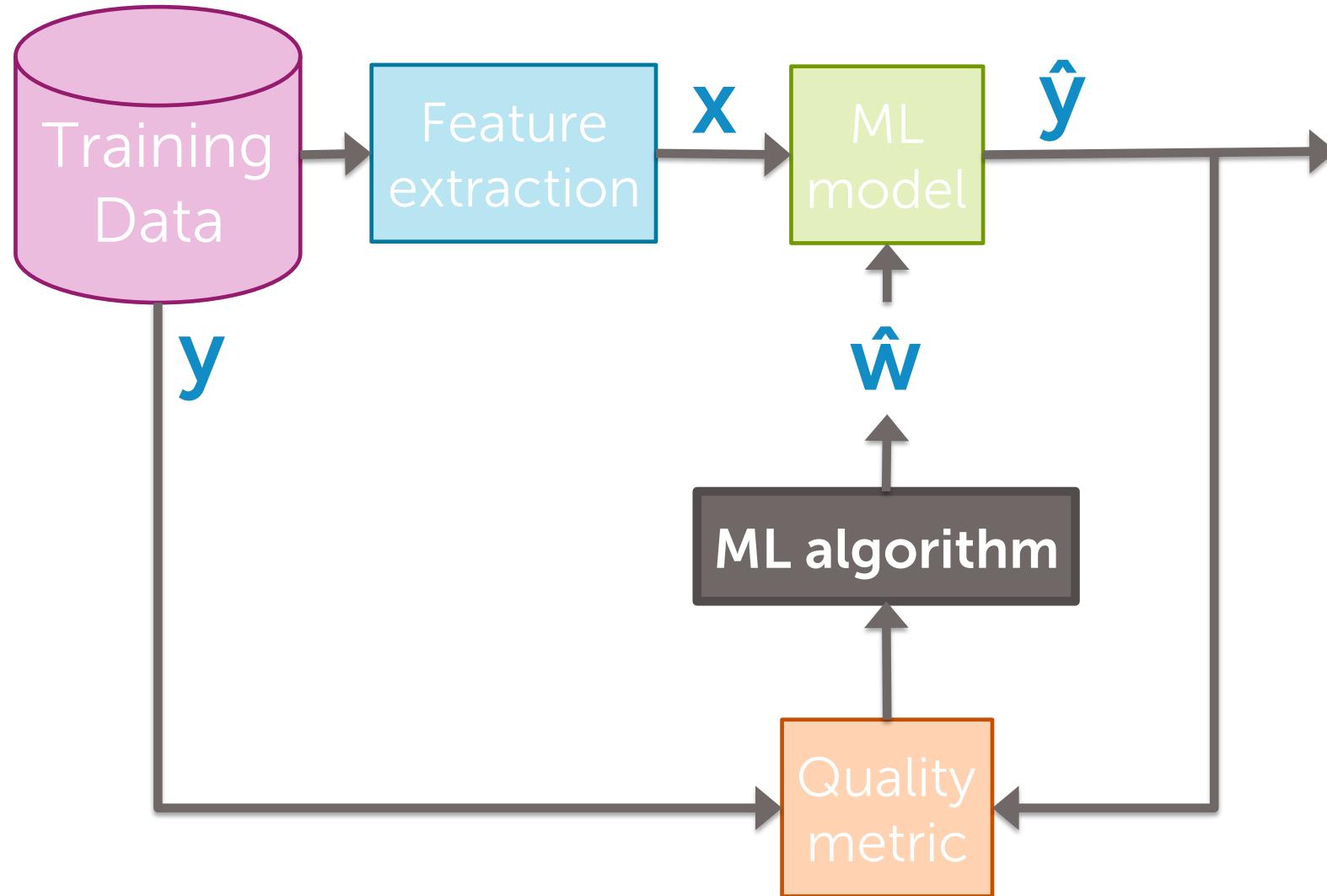
$$\begin{aligned}\hat{y}_{1001 \text{ sq.ft.}} - \hat{y}_{1000 \text{ sq.ft.}} &= \hat{w}_0 + \hat{w}_1 \cdot 1001 \text{ sq.ft.} \\ &\quad - (\hat{w}_0 + \hat{w}_1 \cdot 1000 \text{ sq.ft.}) \\ &= \hat{w}_1\end{aligned}$$

predicted change in the output
per unit change in input

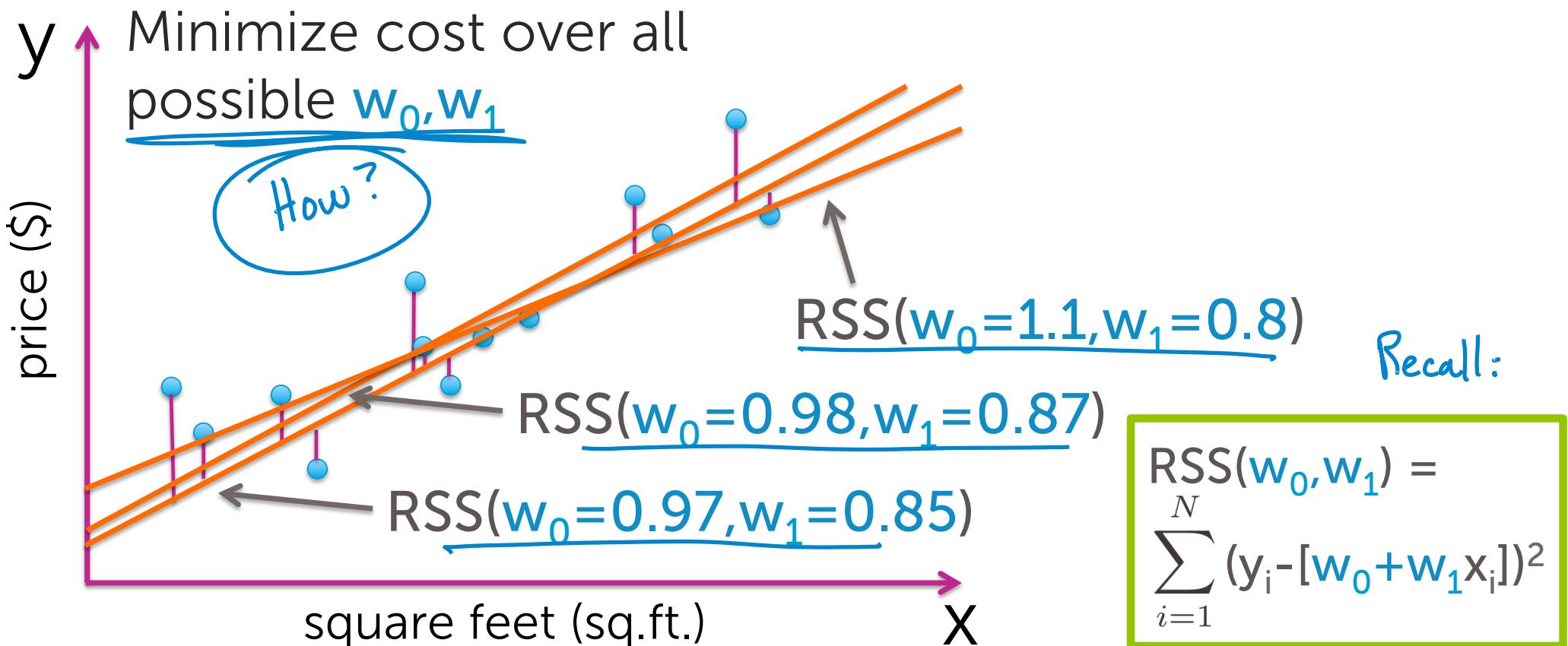
Interpreting the coefficients



Algorithms for fitting the model

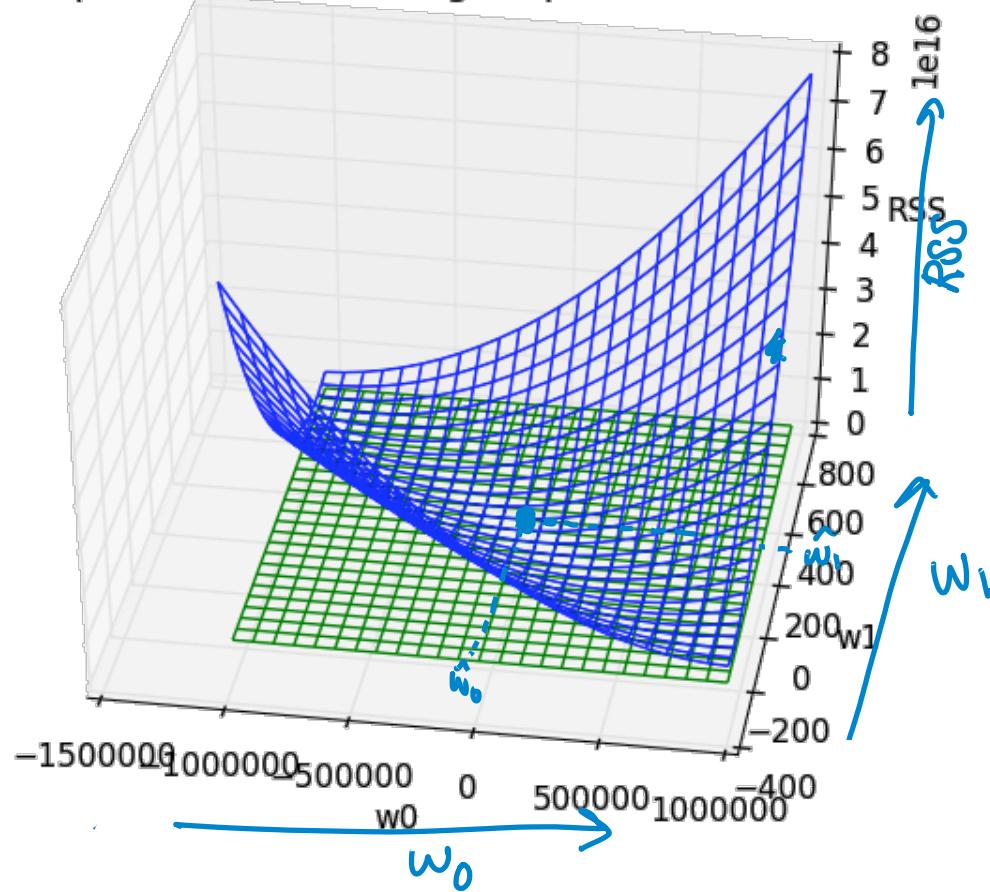


Find “best” line



Minimizing the cost

3D plot of RSS with tangent plane at minimum



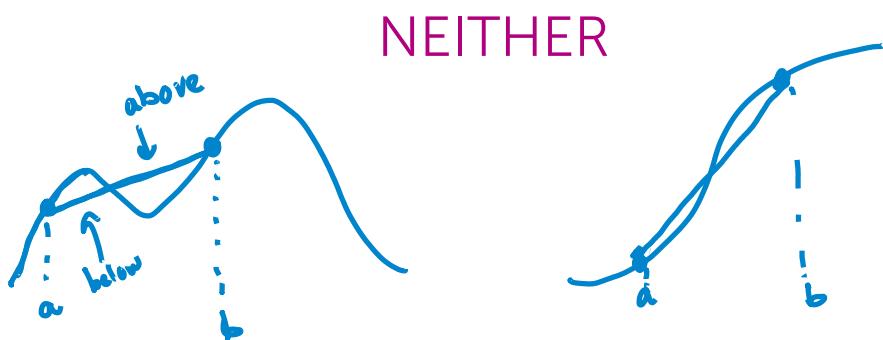
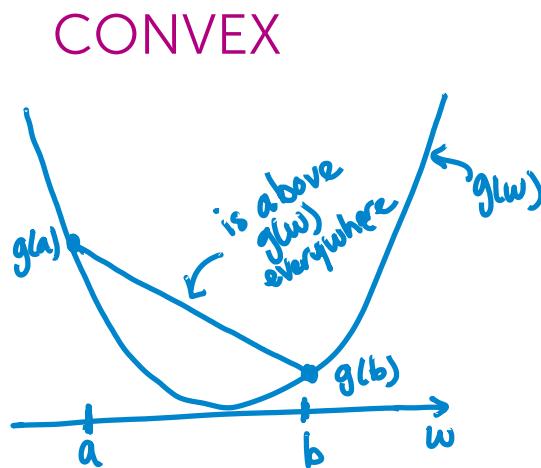
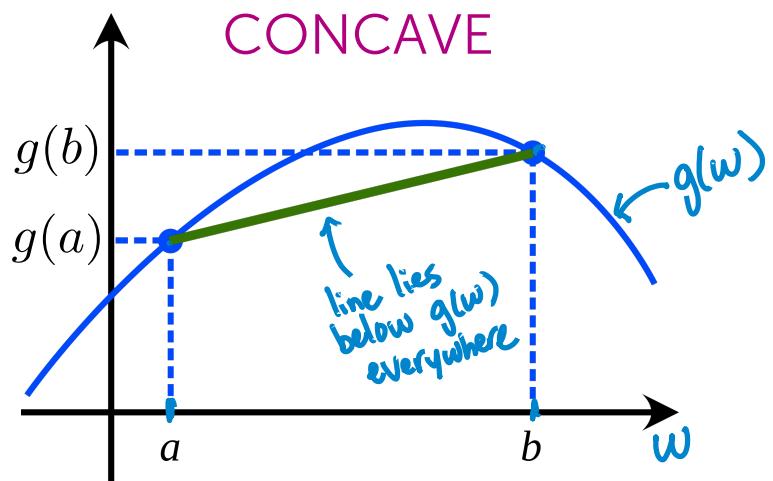
Minimize function
over all possible w_0, w_1

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

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

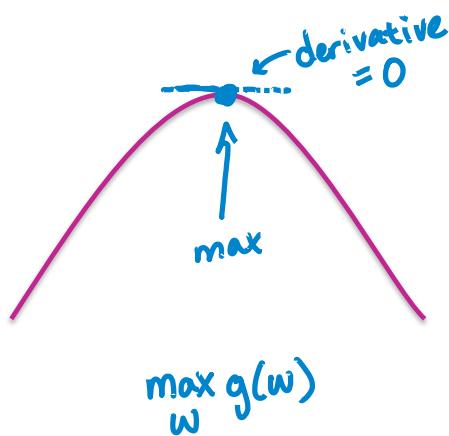
An aside on optimization

Convex/concave functions

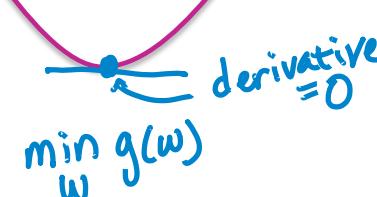


Finding the max or min analytically

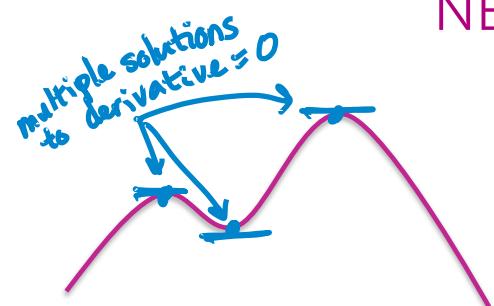
CONCAVE



CONVEX



NEITHER



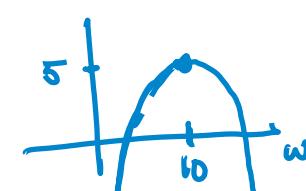
no solution
to derivative = 0

Example:

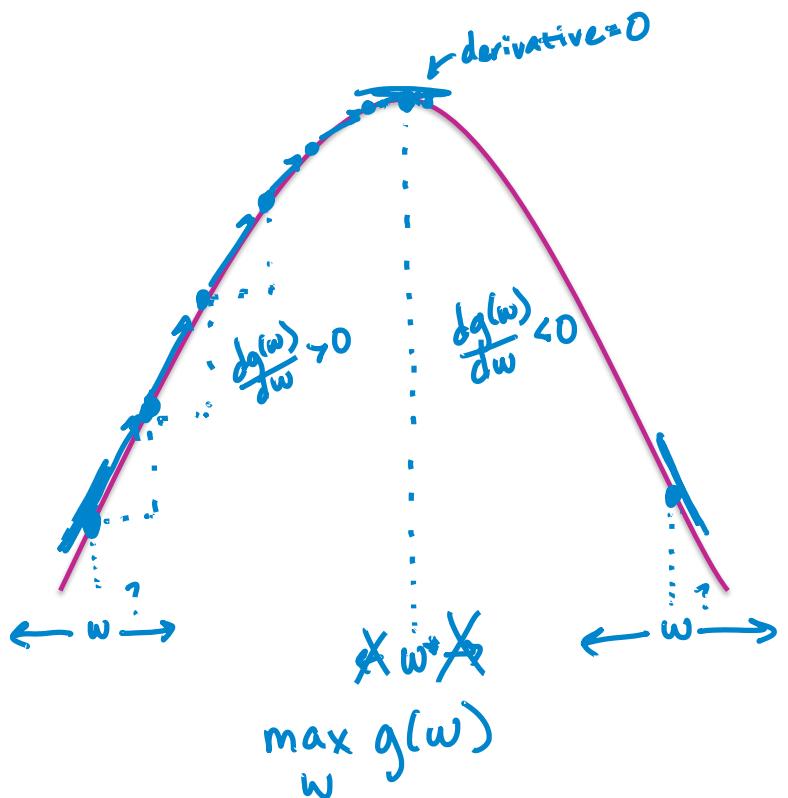
$$g(w) = 5 - (w-10)^2$$

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

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



Finding the max via hill climbing



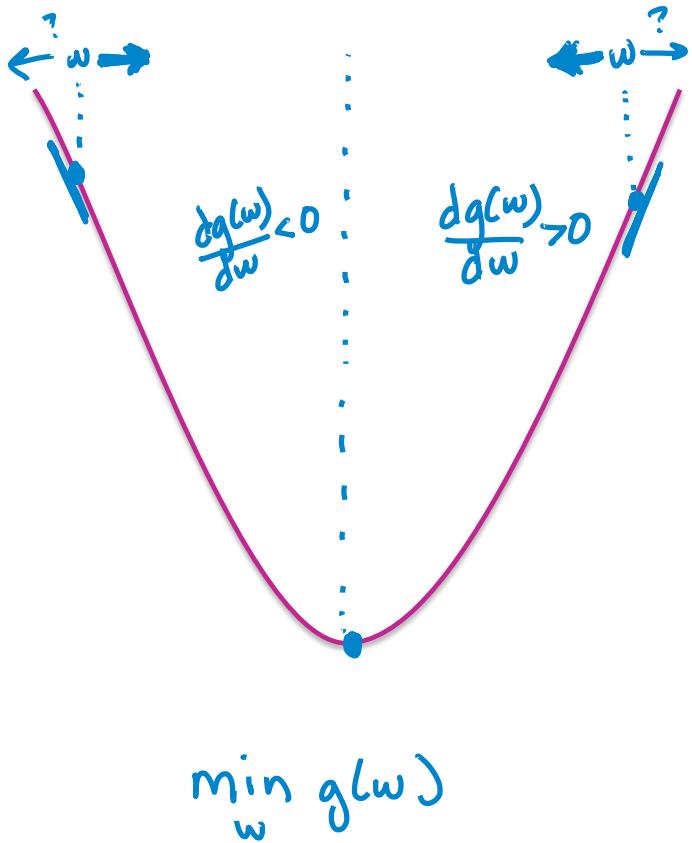
How do we know whether to move w to right or left?
(inc. or dec. the value of w ?)

while not converged

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

iteration t stepsize

Finding the min via hill descent



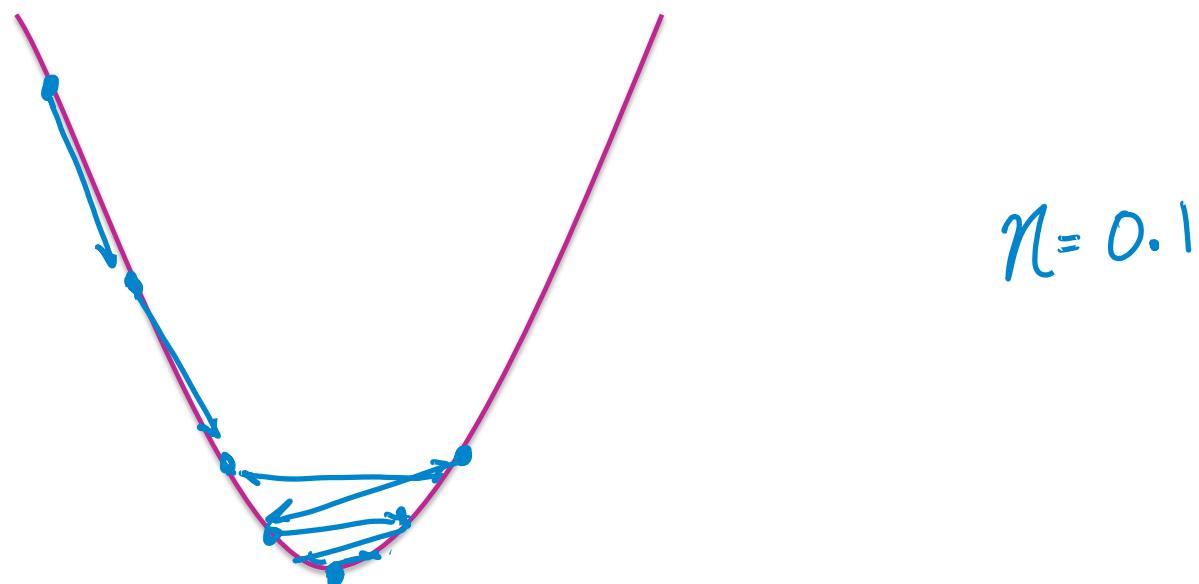
when derivative is positive, we want to decrease w
and when derivative is negative, we want to increase w

Algorithm:

while not converged
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} \Big|_{w^{(t)}}$

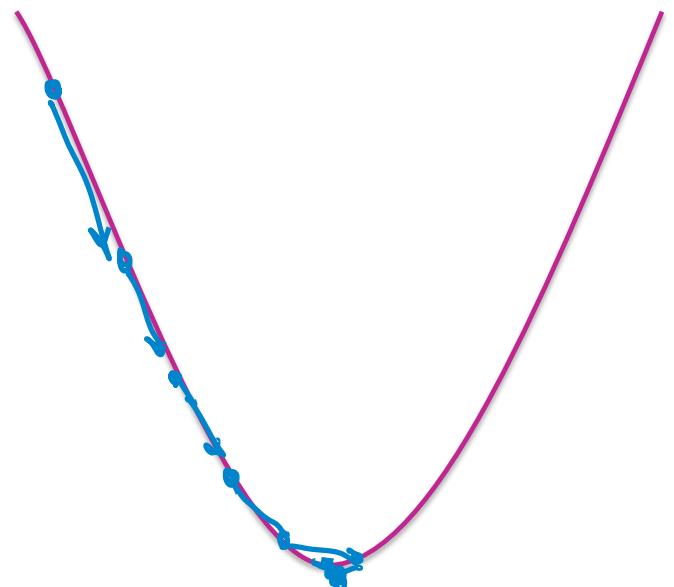
Choosing the stepsize— Fixed stepsize

η



Choosing the stepsize— Decreasing stepsize

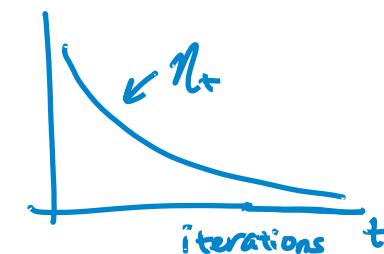
or stepsize schedule



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

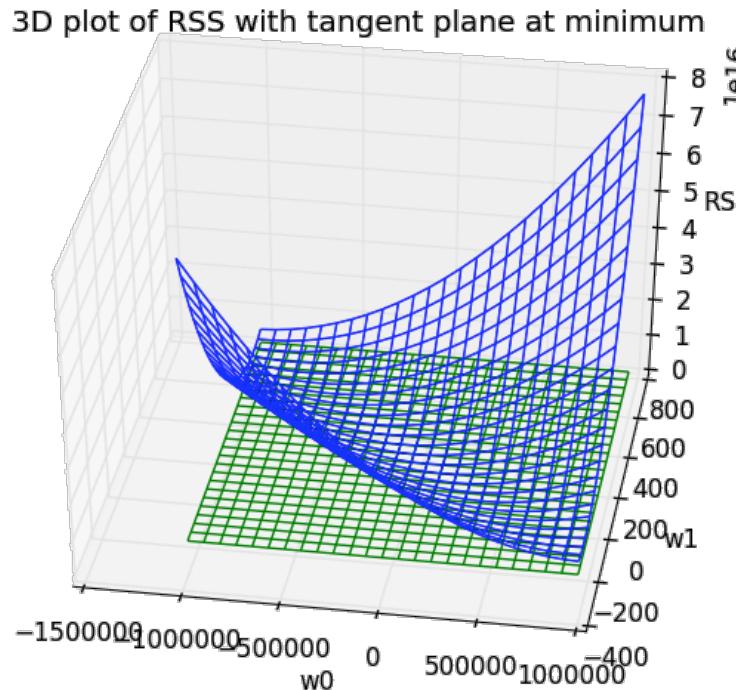
$$\left| \frac{dg(w)}{dw} \right| < \epsilon$$

↑
threshold
to be set

Algorithm:

while not converged
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} \Big|_{w^{(t)}}$

Moving to multiple dimensions: Gradients



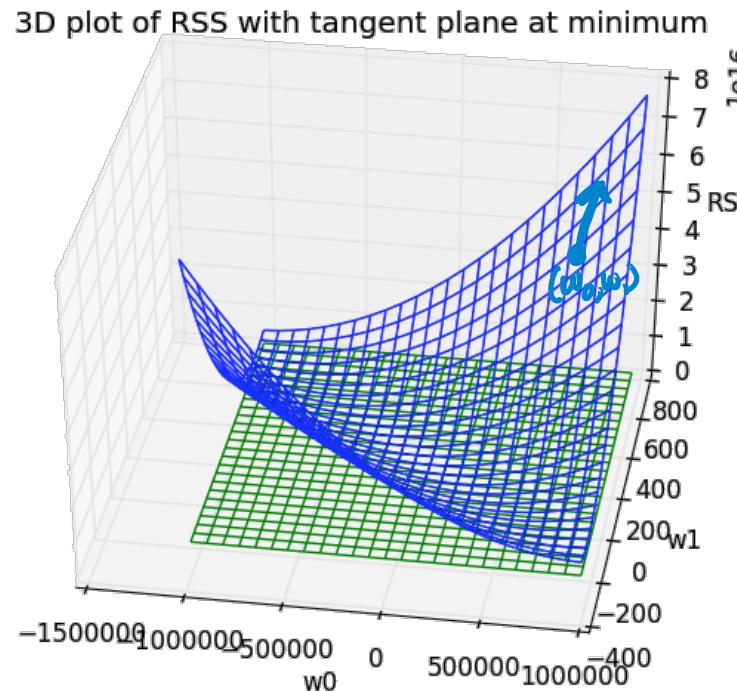
$$\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{bmatrix}$$

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

(p+1)-dimensional vector

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

Gradient example



$$g(w) = 5w_0 + 10w_0w_1 + 2w_1^2$$

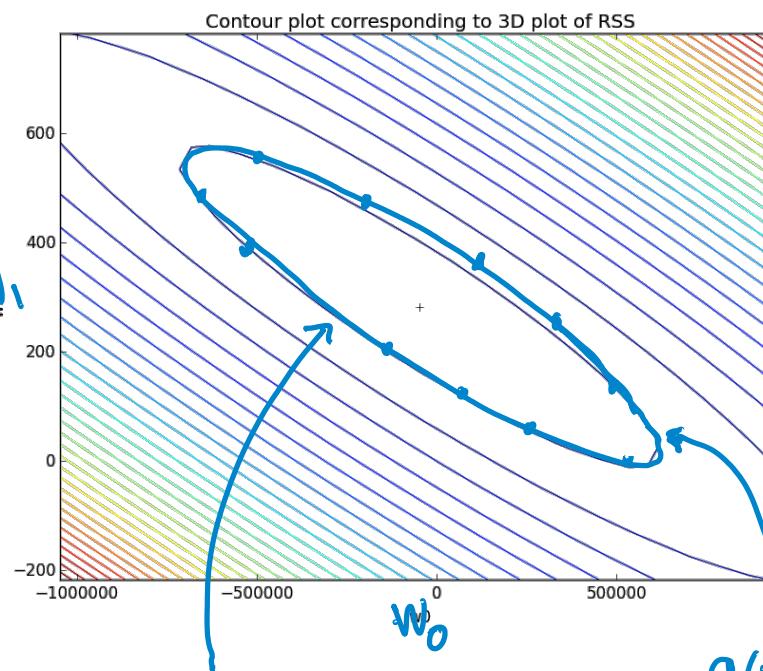
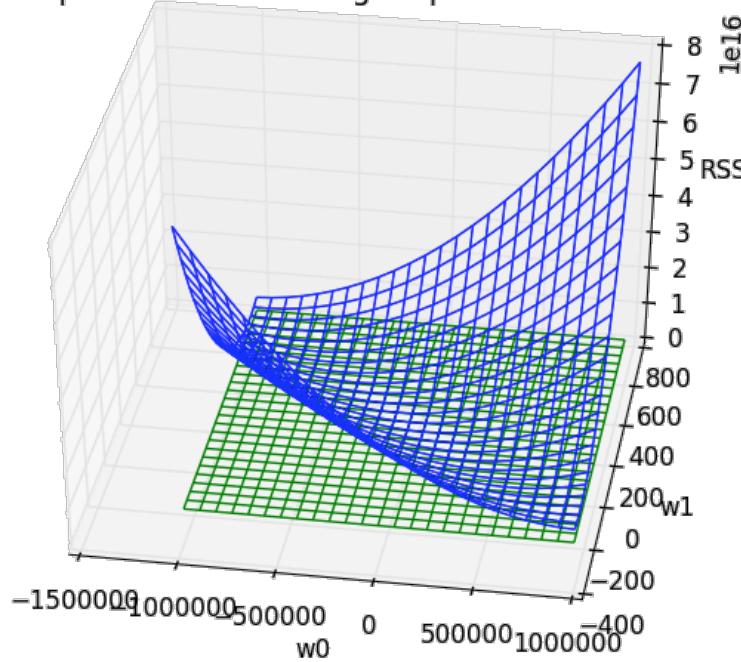
$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

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

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

Contour plots

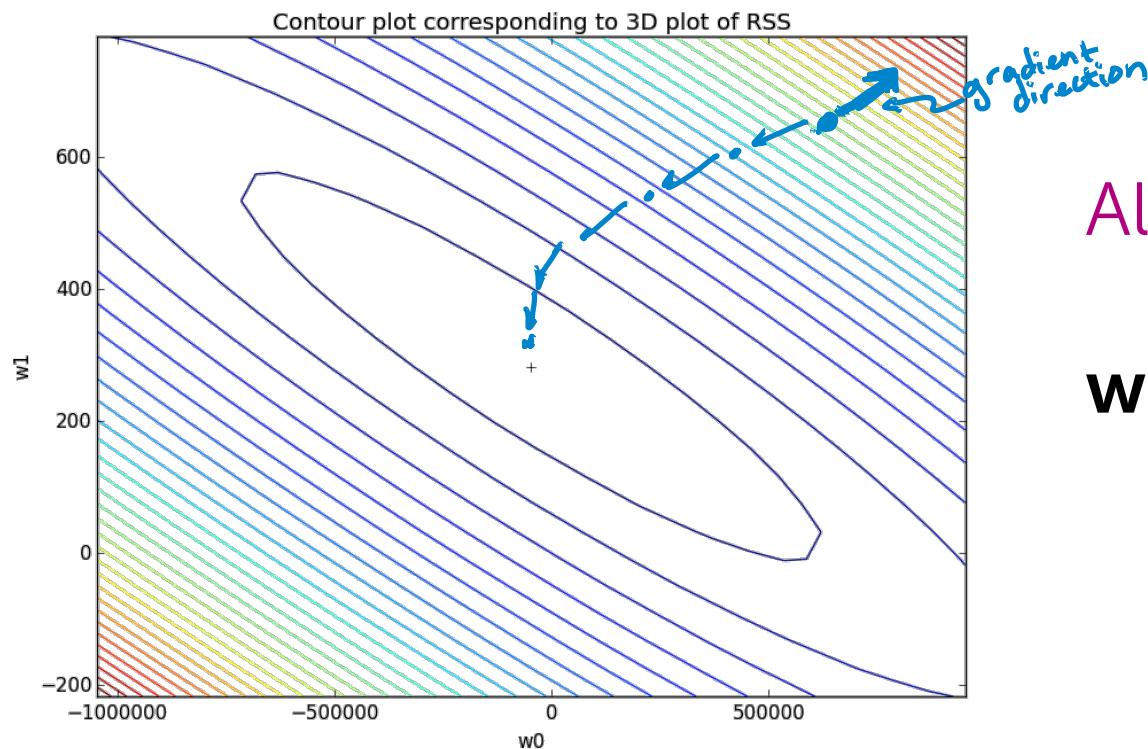
3D plot of RSS with tangent plane at minimum



a slice of the
3D surface

$g(w_0, w_1)$

Gradient descent



Algorithm:

while not converged
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla g(w^{(t)})$

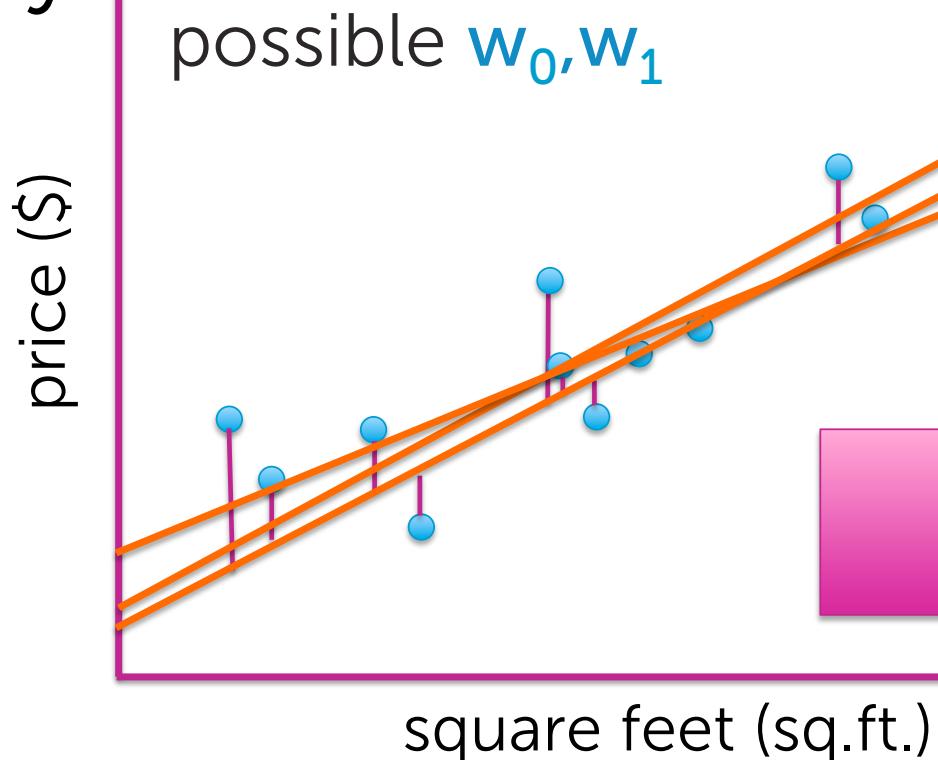
$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} - \eta \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

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

Finding the least squares line

Find “best” line

y
Minimize cost over all
possible w_0, w_1



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

CONVEX

⇒ 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:

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

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

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

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

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

$$\text{RSS}(\mathbf{w}_0, \mathbf{w}_1) = \sum_{i=1}^N (y_i - (\mathbf{w}_0 + \mathbf{w}_1 x_i))^2$$

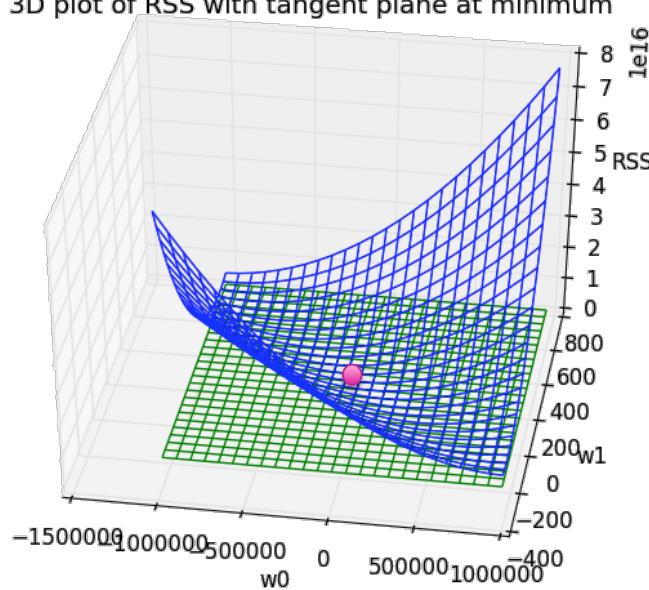
Putting it together:

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

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 house price
average sq-ft.

bottom term:

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

plug in

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

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:

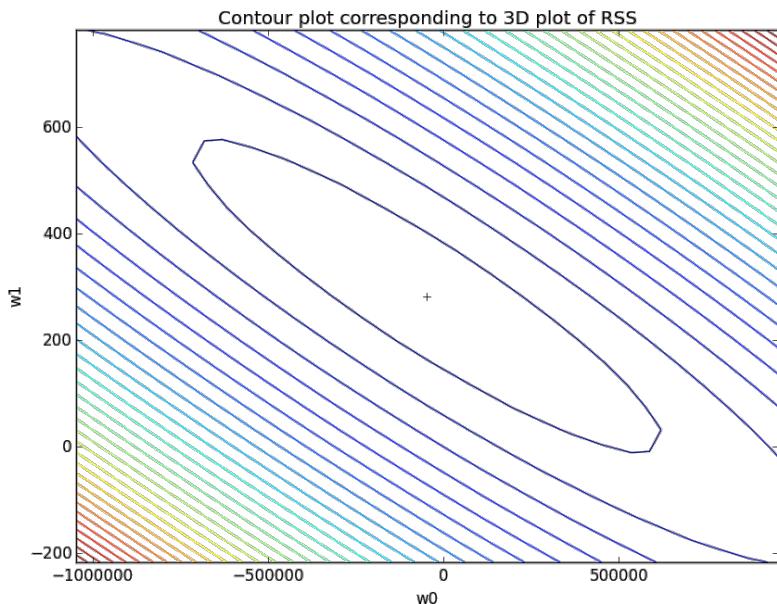
$$\nabla \text{RSS}(\mathbf{w}_0, \mathbf{w}_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (\underline{\mathbf{w}_0 + w_1 x_i})] \\ -2 \sum_{i=1}^N [y_i - (\underline{\mathbf{w}_0 + w_1 x_i})] x_i \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] x_i \end{bmatrix}$$

Annotations:

- y_i is labeled "actual house sales observation".
- $\hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)$ is labeled "predicted value".

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 $\leftarrow^{(w_0^{(t+1)}, w_1^{(t+1)})}$

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

If overall, under predicting \hat{y}_i , then $\sum [y_i - \hat{y}_i]$ is positive
 $\rightarrow w_0$ 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

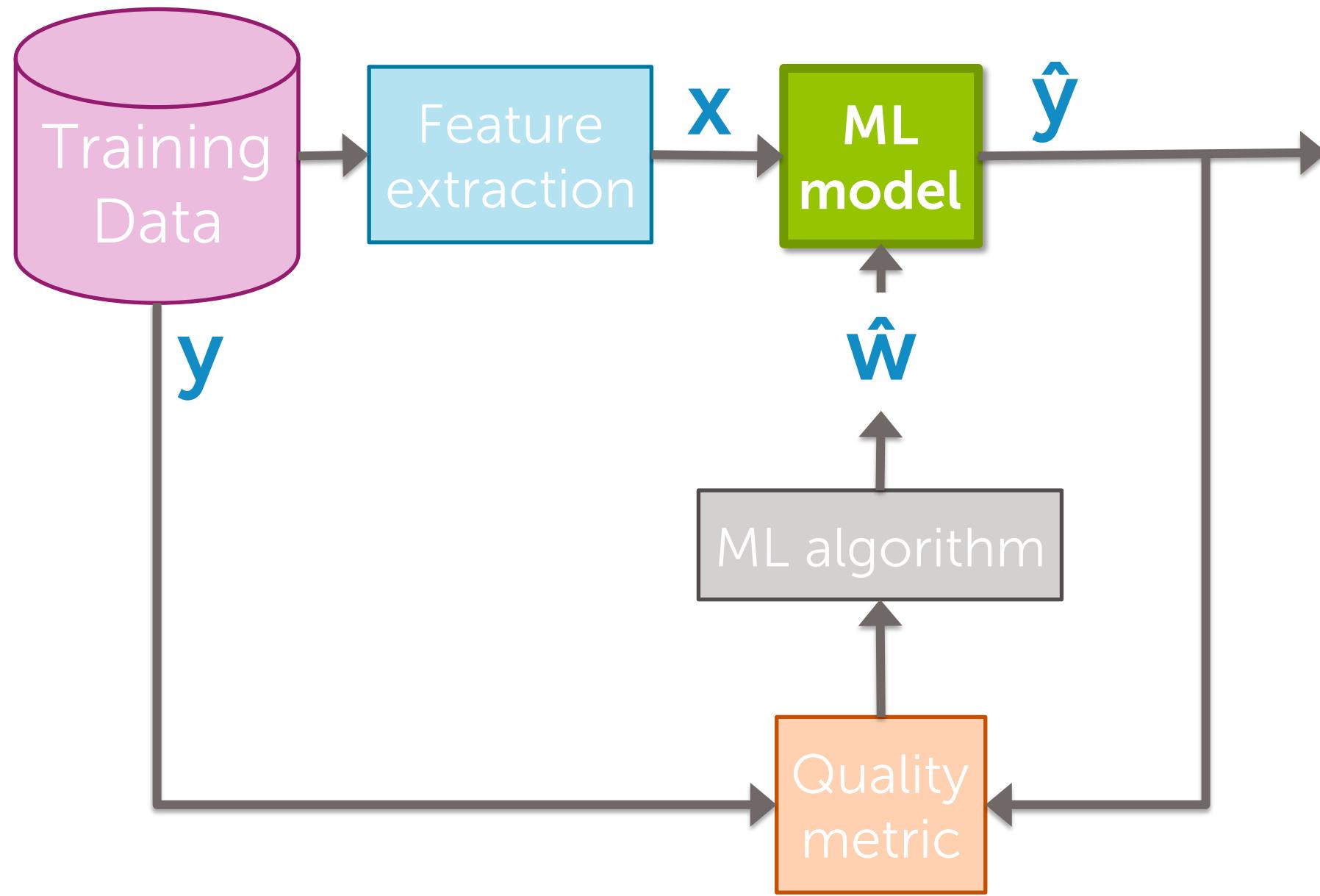
Summary for simple linear regression

What you can do now...

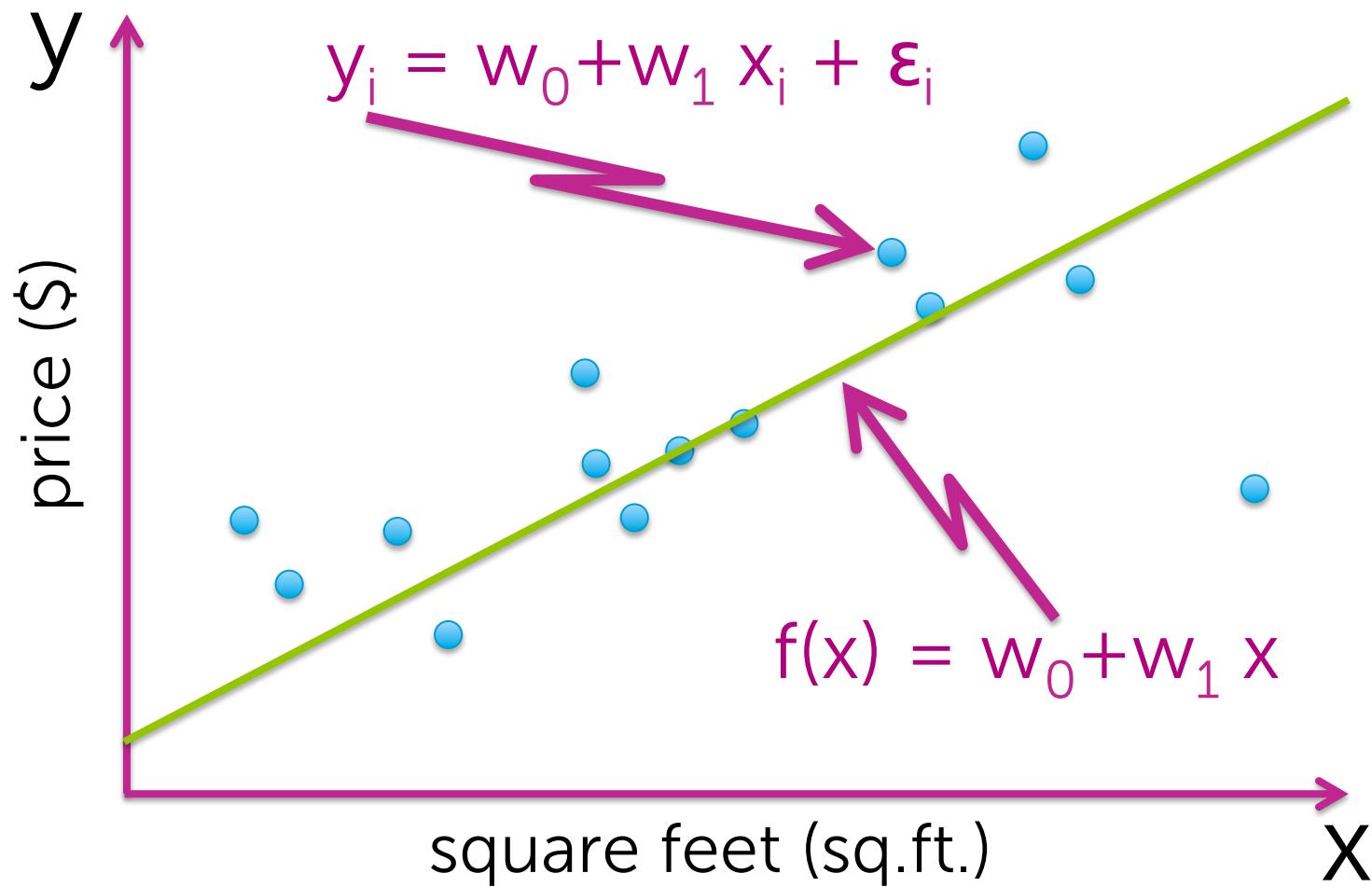
- Describe the input (features) and output (real-valued predictions) of a regression model
- Calculate a goodness-of-fit metric (e.g., RSS)
- Estimate model parameters to minimize RSS using gradient descent
- Interpret estimated model parameters
- Exploit the estimated model to form predictions
- Discuss the possible influence of high leverage points
- Describe intuitively how fitted line might change when assuming different goodness-of-fit metrics

Multiple Regression:

Linear regression with multiple features



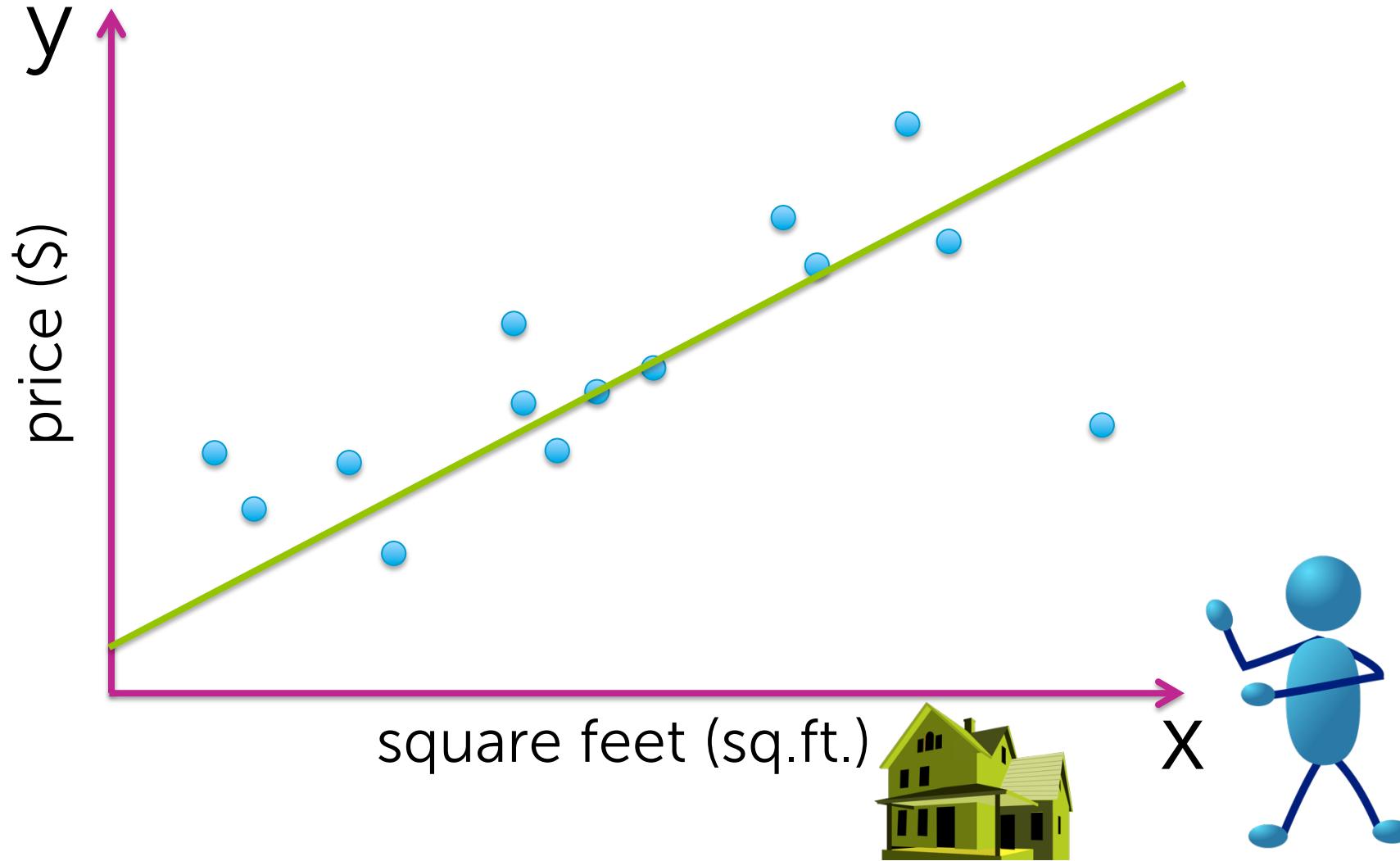
Simple linear regression model



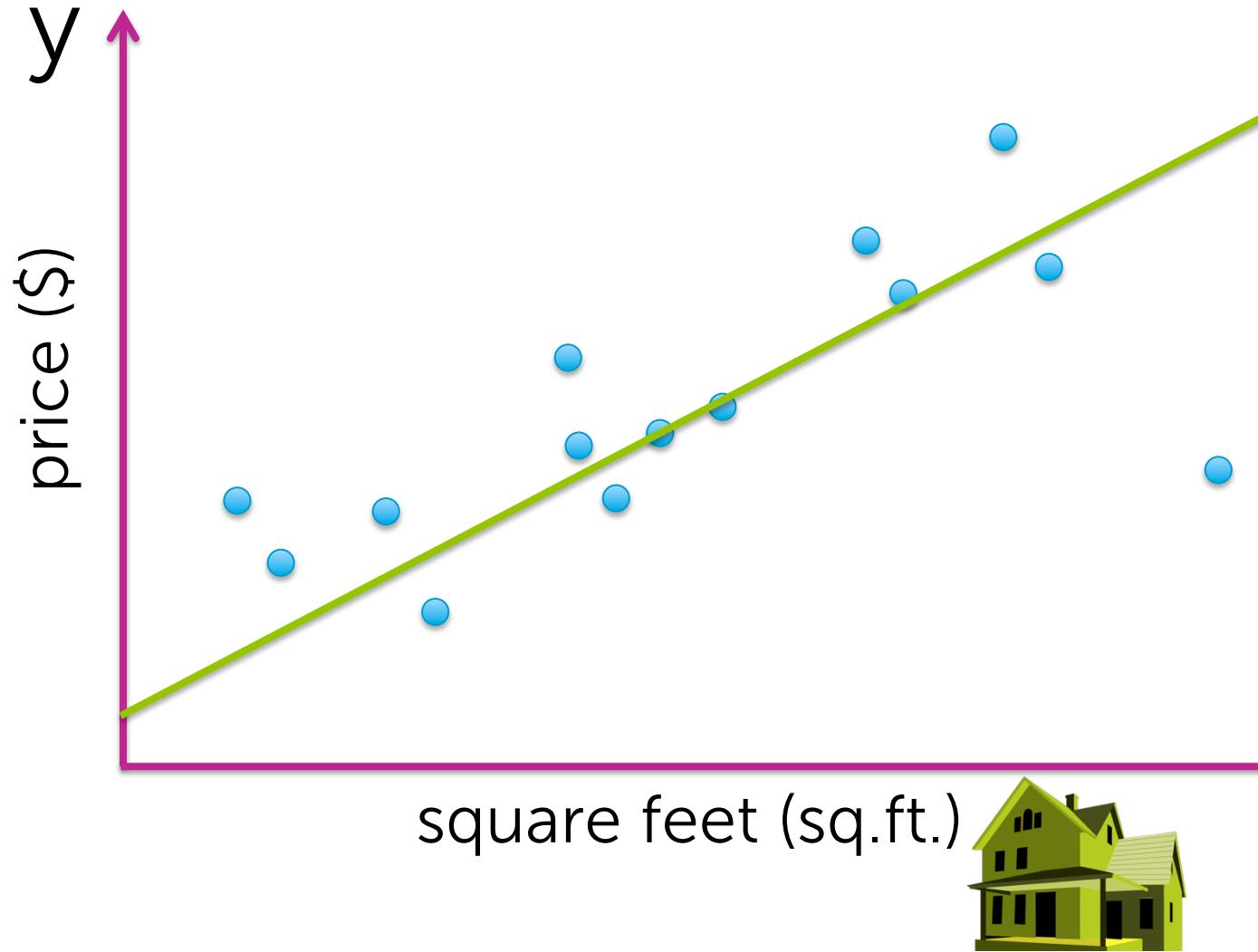
More complex functions of a single input

Polynomial regression

Fit data with a line or ... ?



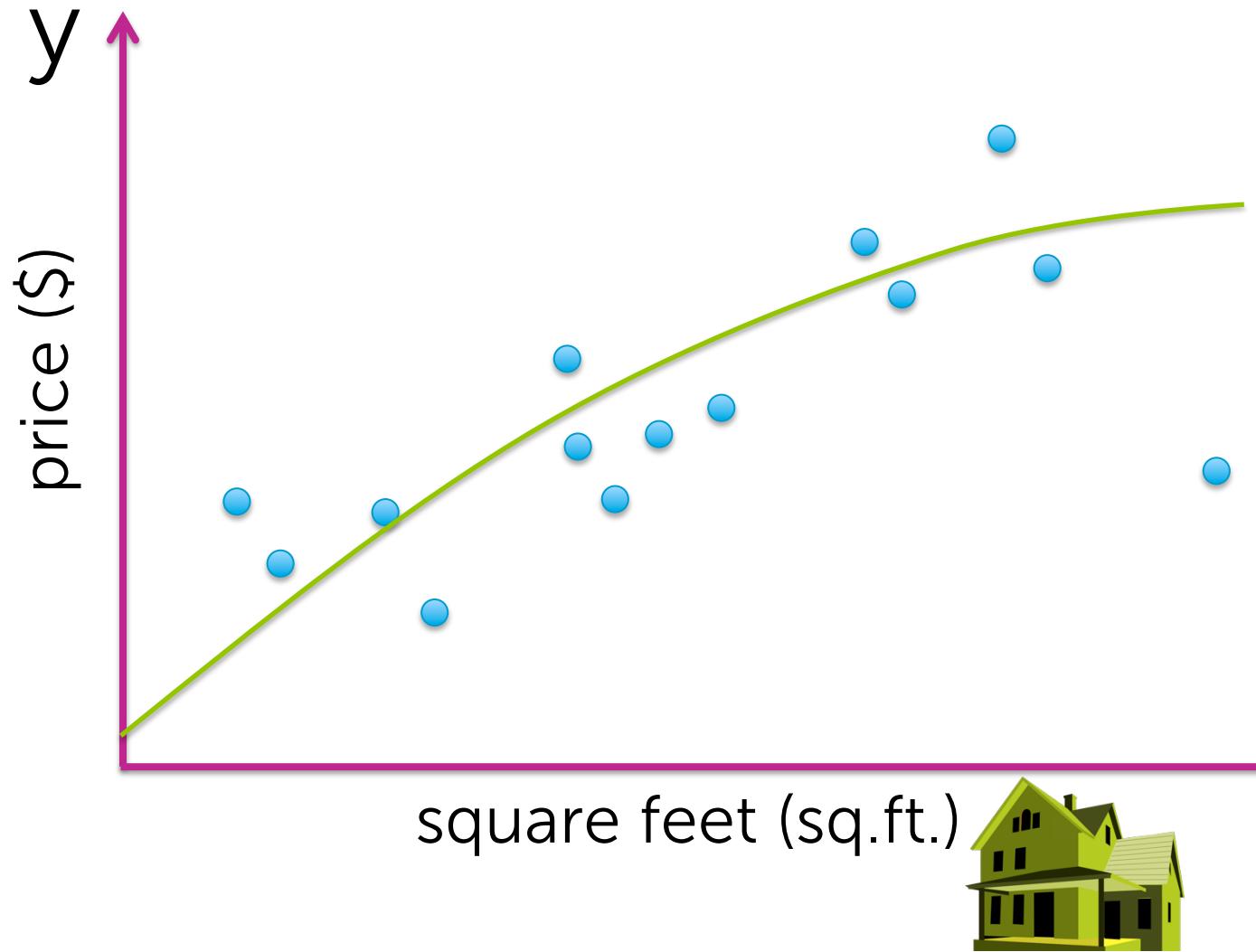
Fit data with a line or ... ?



Dude, it's
not a linear
relationship!



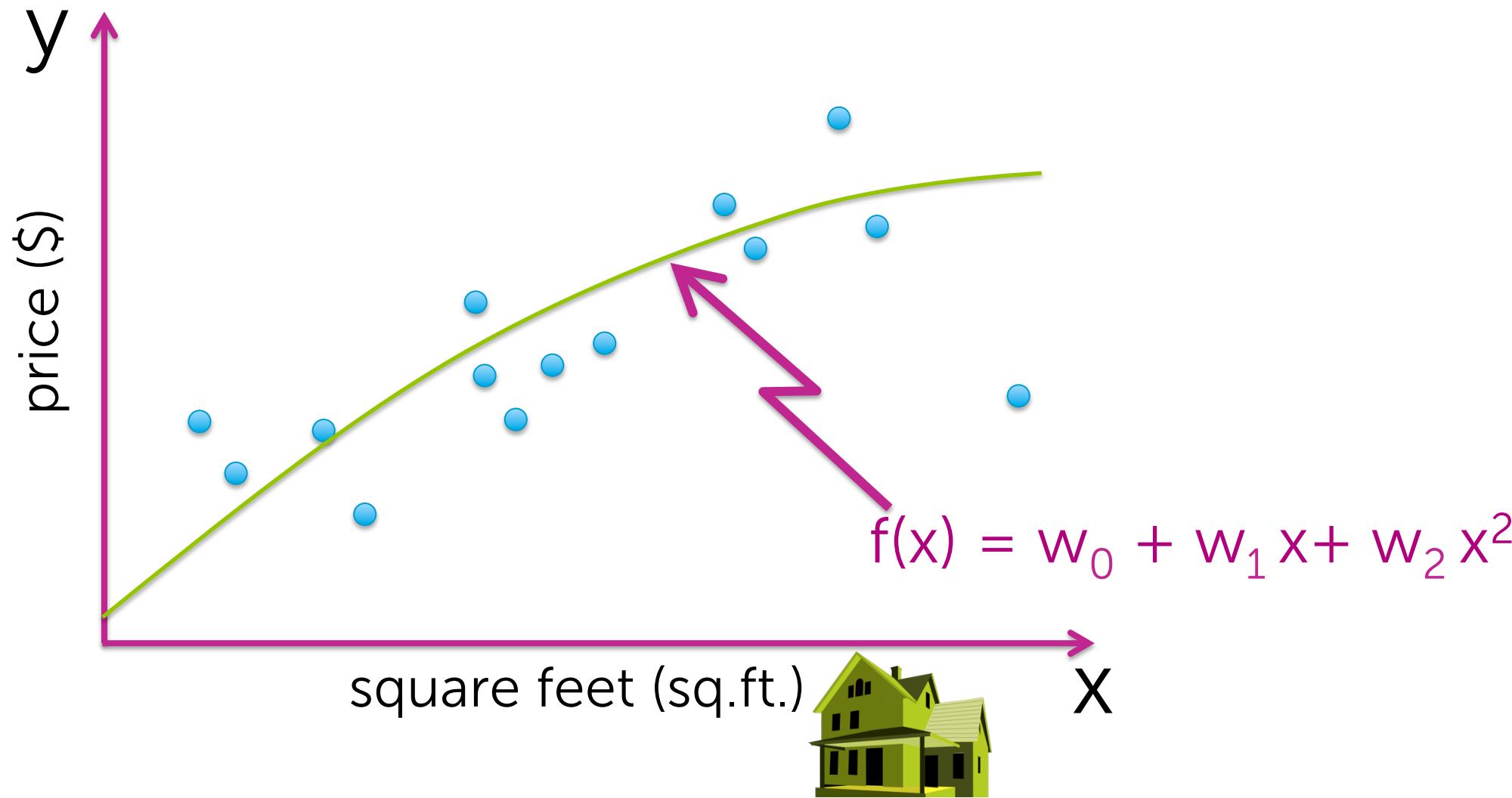
What about a quadratic function?



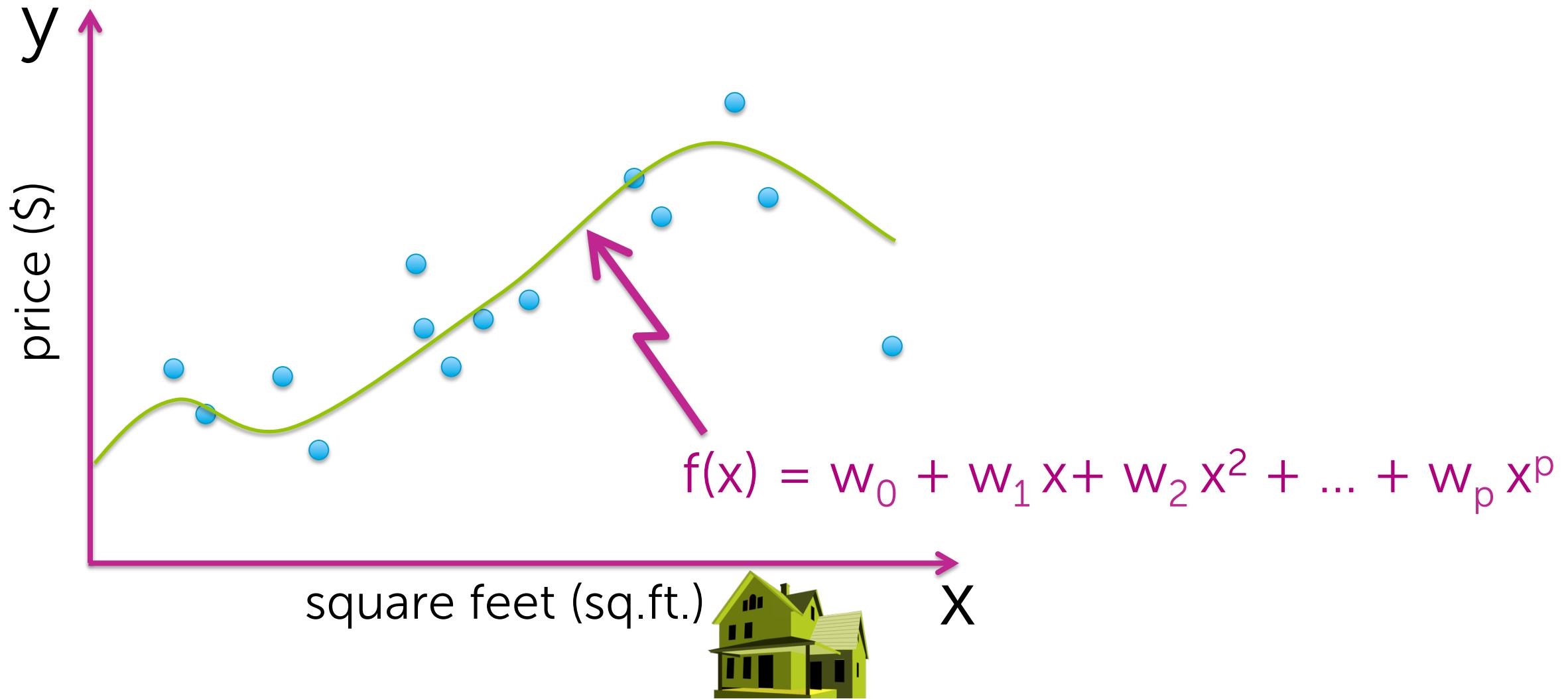
Dude, it's
not a linear
relationship!



What about a quadratic function?



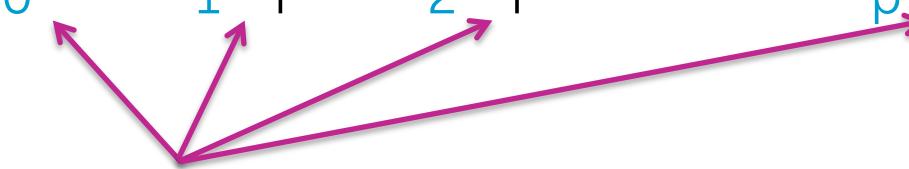
Even higher order polynomial



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

More generally...

Generic basis expansion

Model:

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

*jth regression coefficient
or weight*

jth feature

Generic basis expansion

Model:

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

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

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

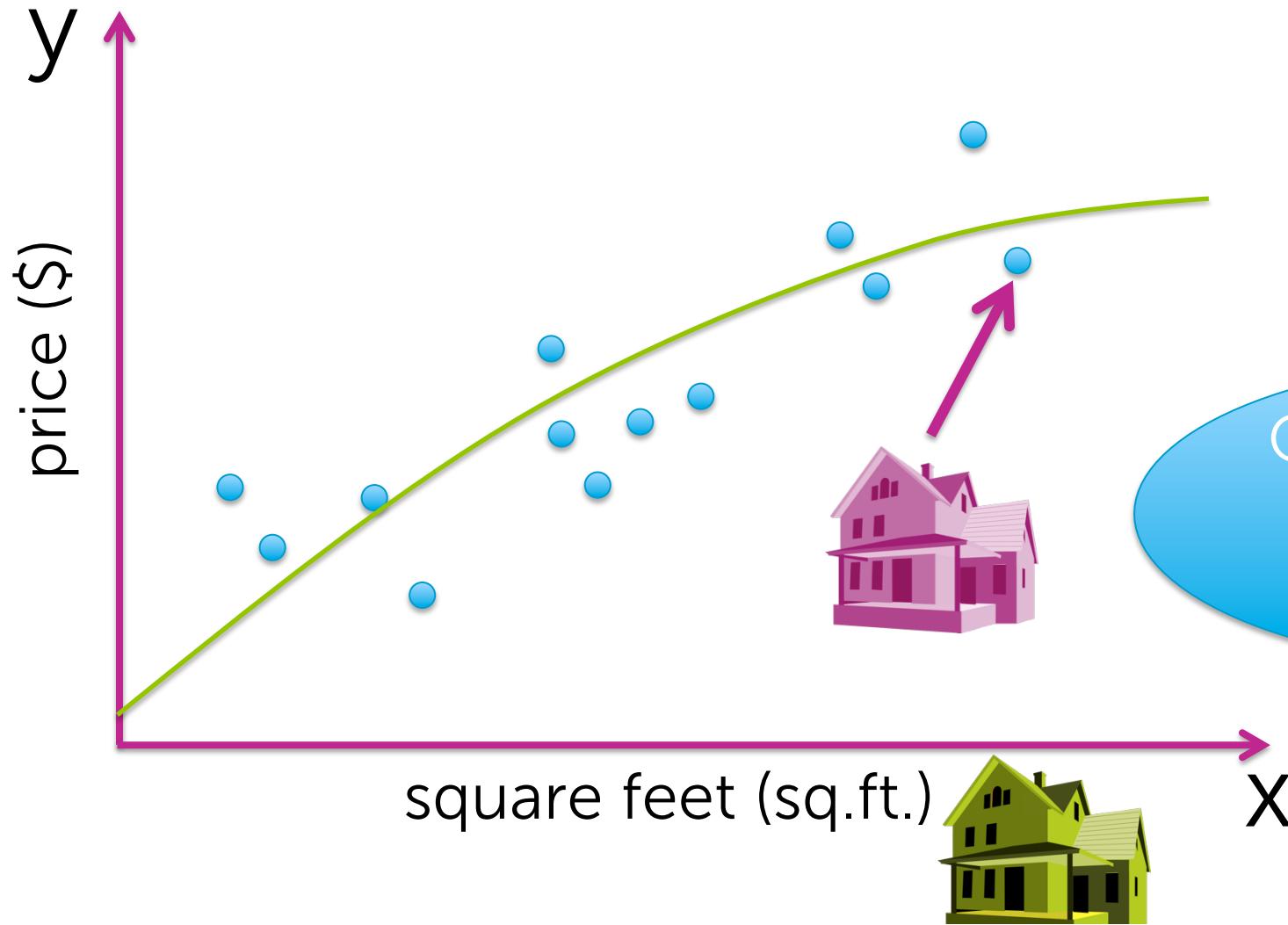
feature 3 = $h_2(x)$... e.g., x^2 or $\sin(2\pi x/12)$

...

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

Incorporating multiple inputs

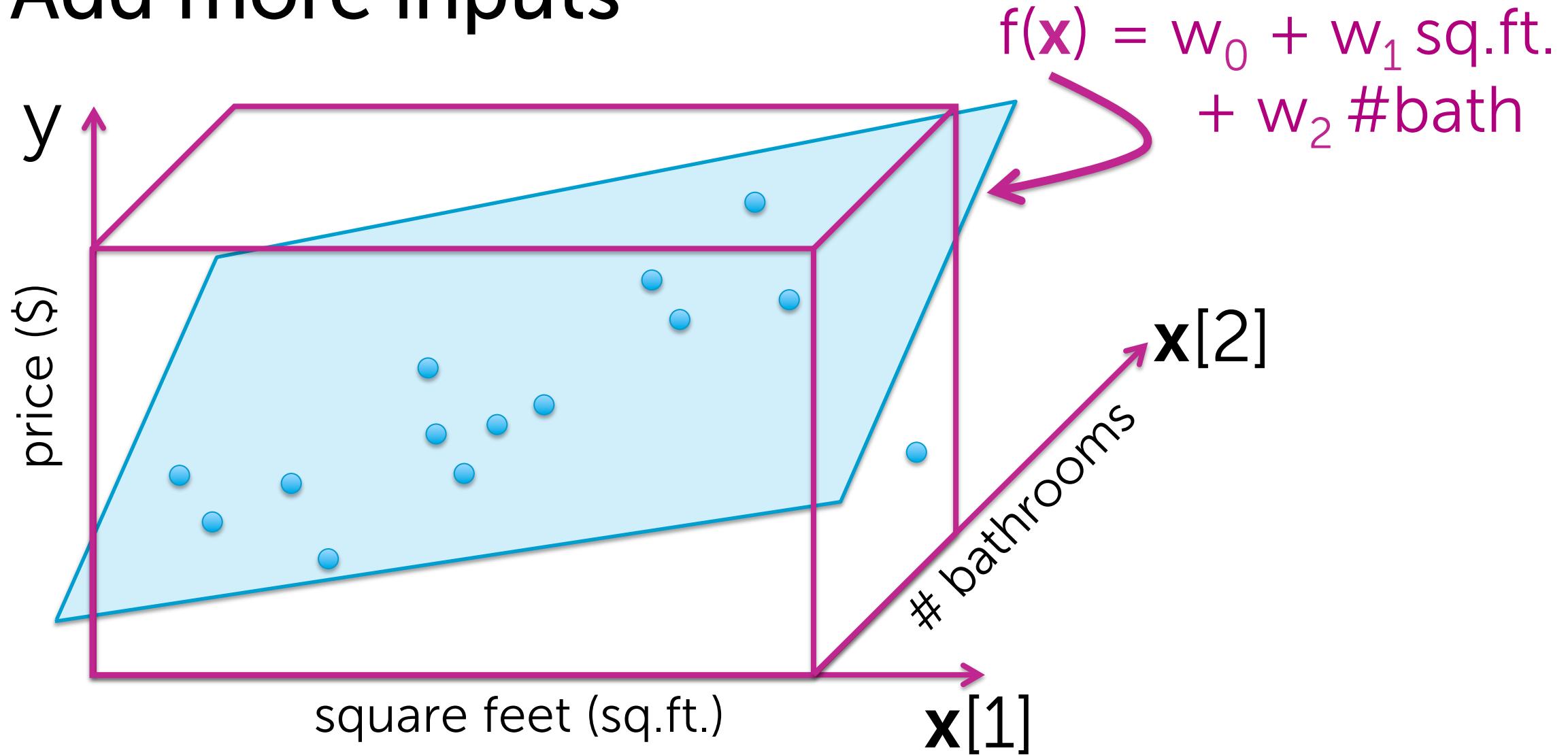
Predictions just based on house size



Only 1 bathroom!
Not same as my
3 bathrooms



Add more inputs



Many possible inputs

- Square feet
- # bathrooms
- # bedrooms
- Lot size
- Year built
- ...

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)

Simple hyperplane

Model:

$$y_i = w_0 + w_1 x_i[1] + \dots + w_d x_i[d] + \varepsilon_i$$

feature 1 = 1

feature 2 = $x[1]$... e.g., sq. ft.

feature 3 = $x[2]$... e.g., #bath

...

feature $d+1 = x[d]$... e.g., lot size

More generically... D-dimensional curve

Model:

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

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

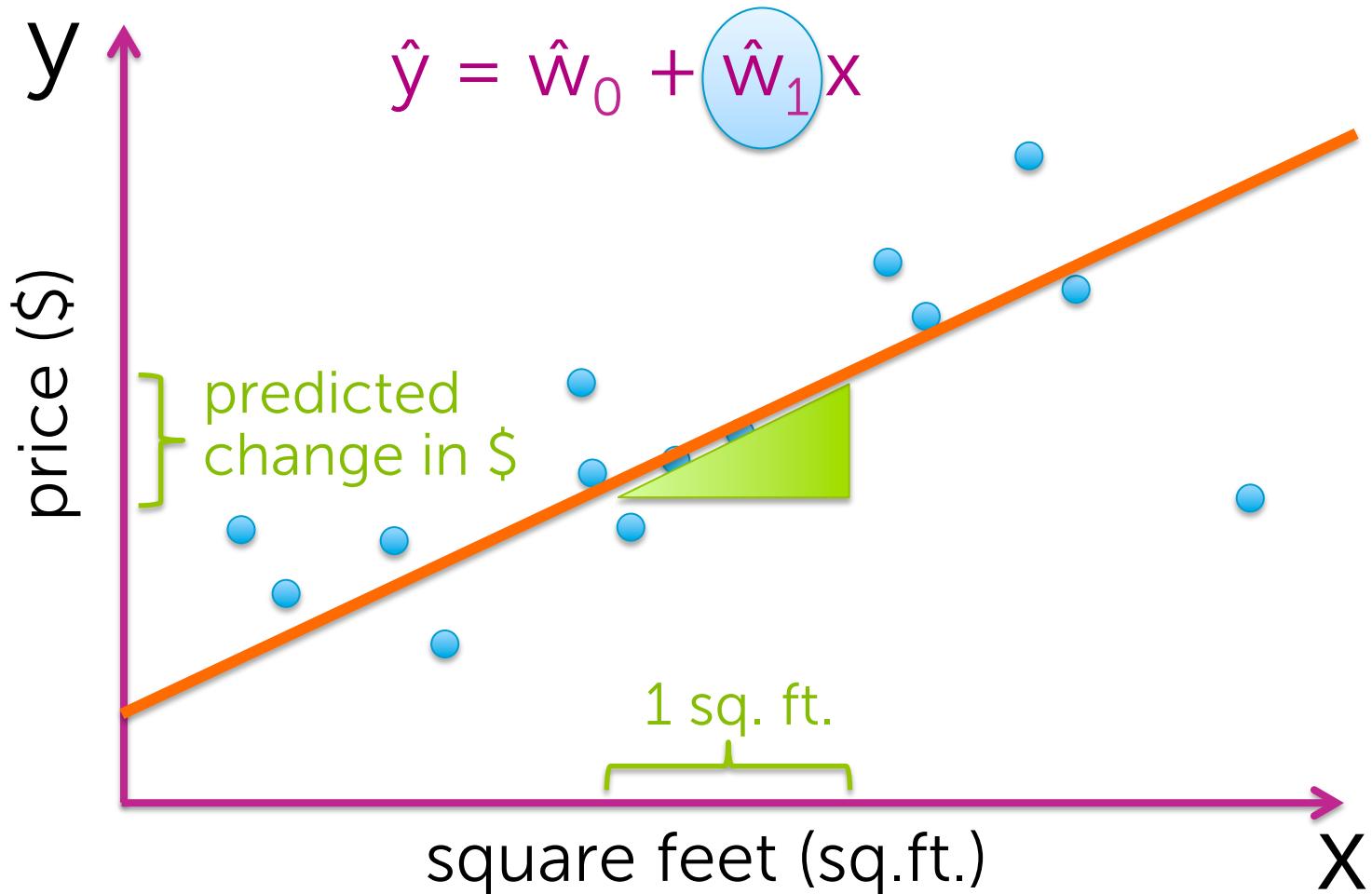
observations (\mathbf{x}_i, y_i) : N

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

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

Interpreting the fitted function

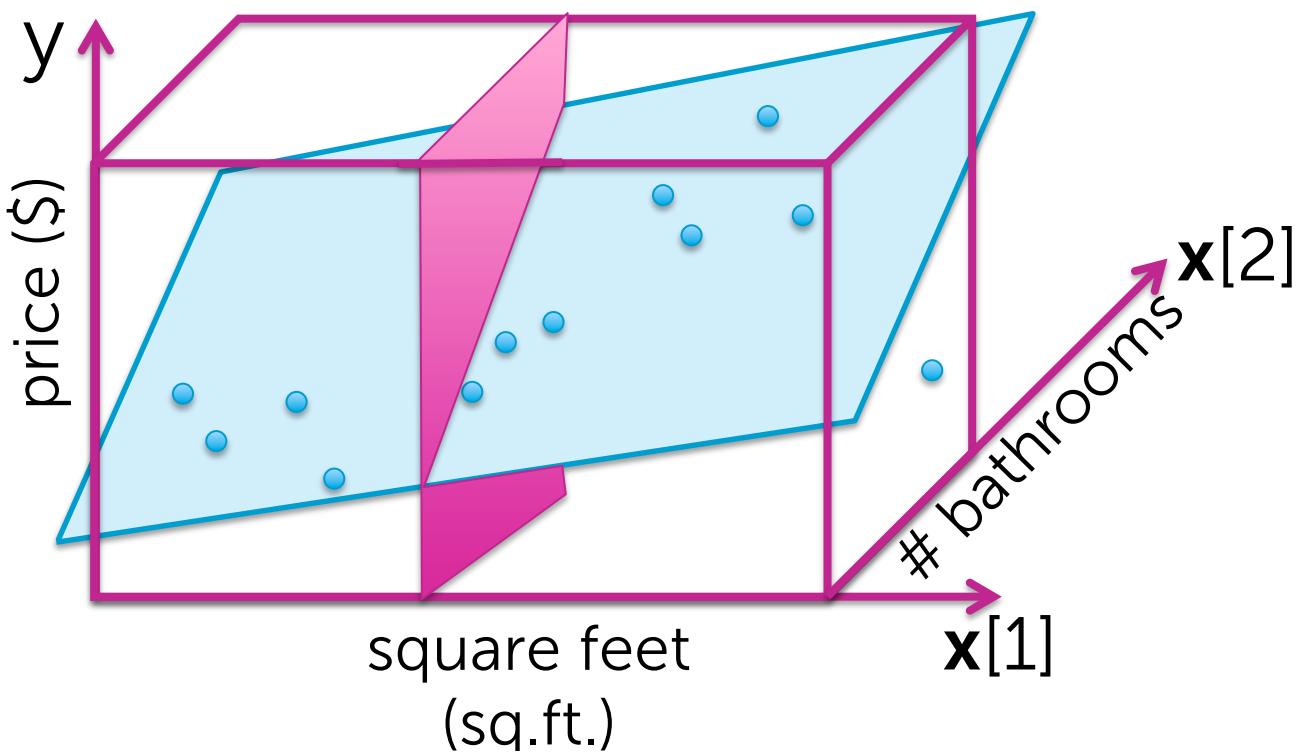
Interpreting the coefficients – Simple linear regression



Interpreting the coefficients – Two linear features

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

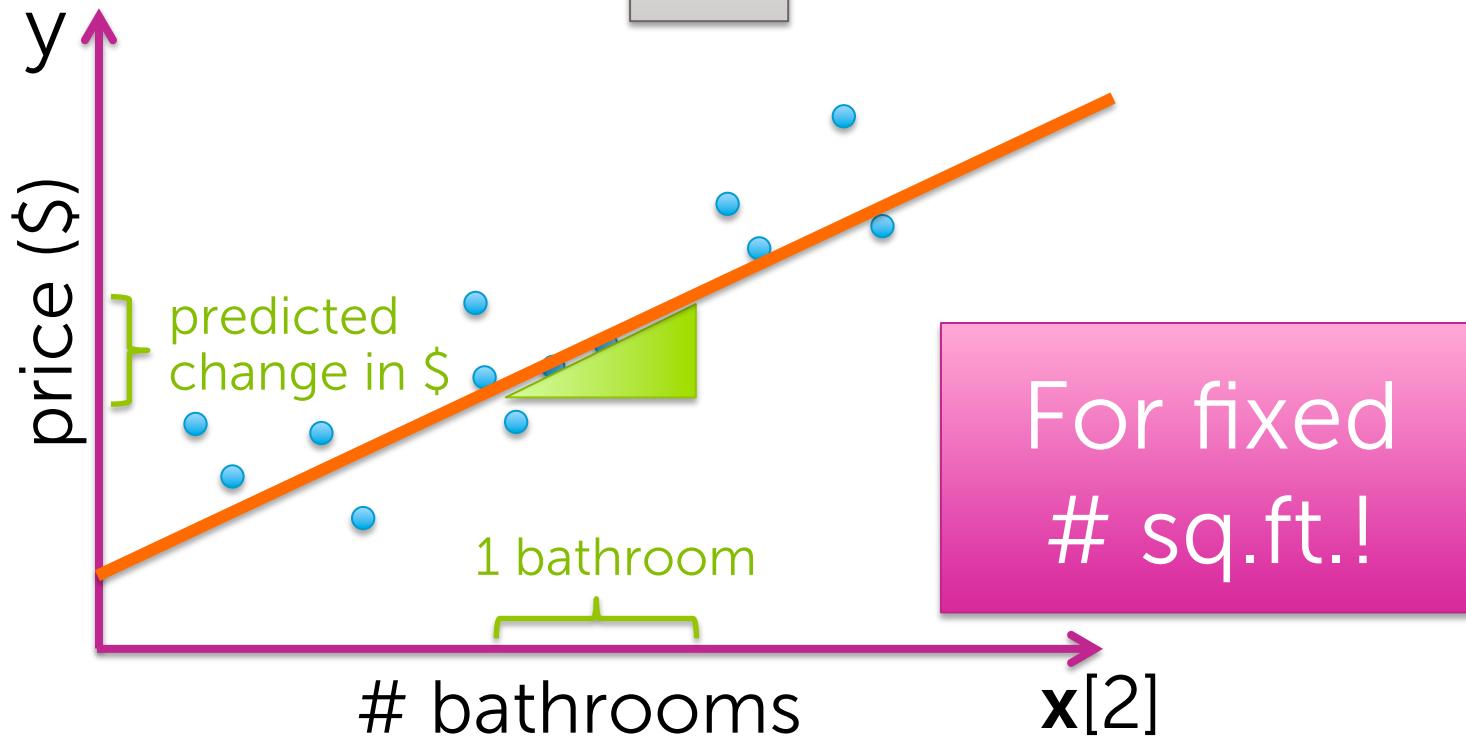
fix



Interpreting the coefficients – Two linear features

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

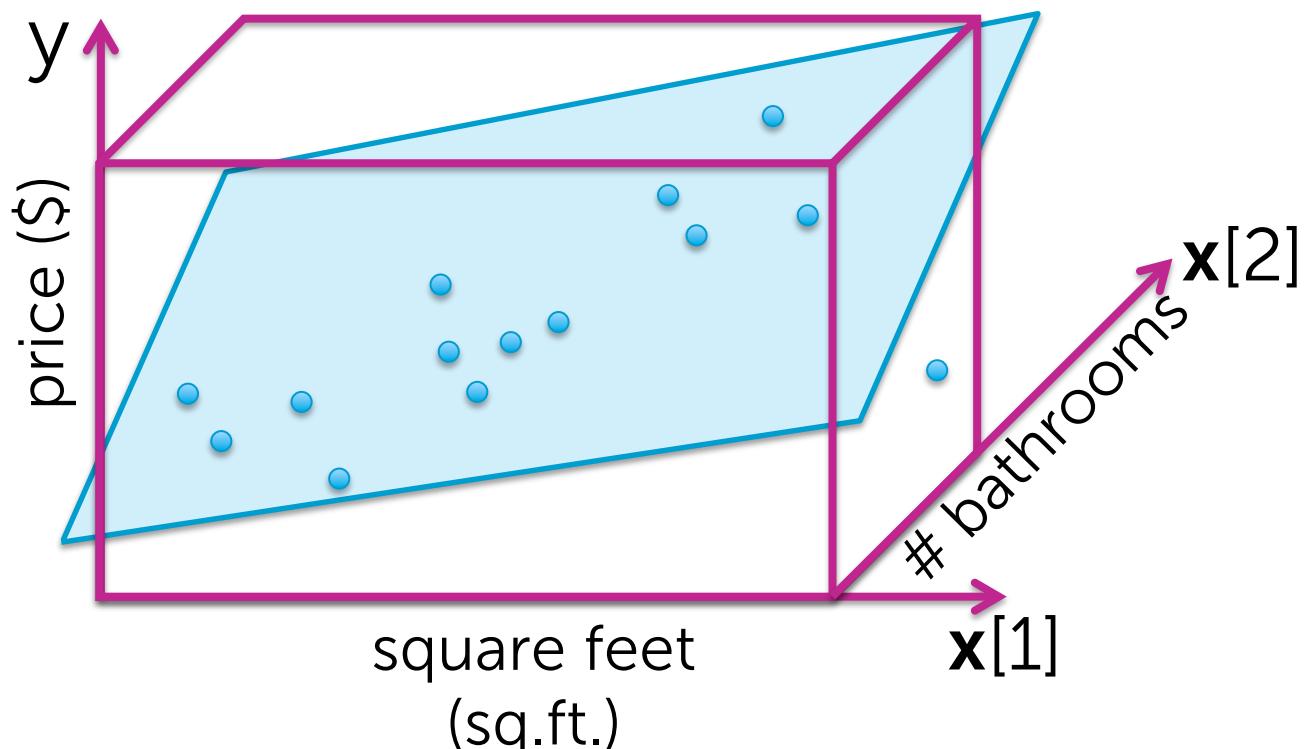
fix



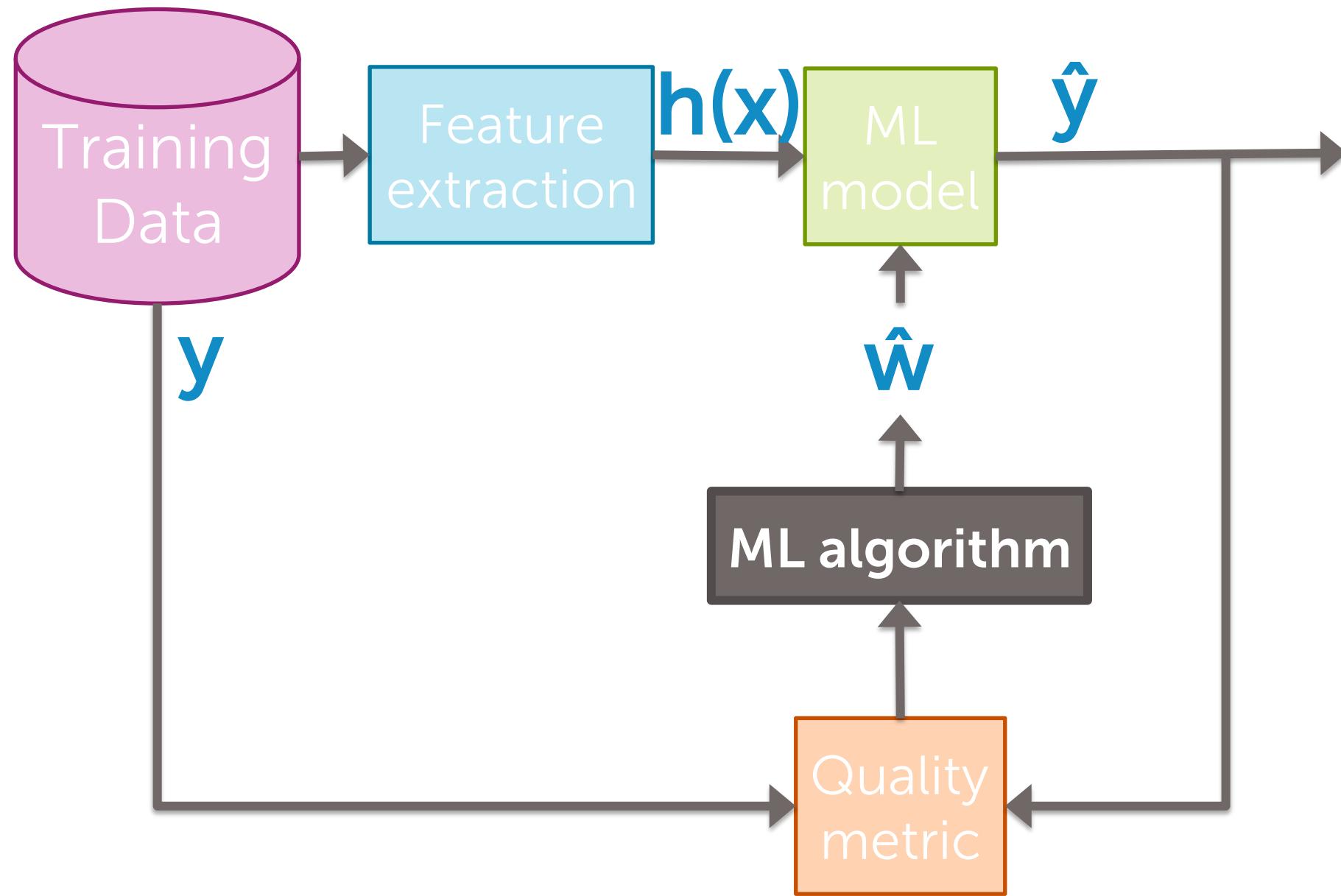
Interpreting the coefficients – Multiple linear features

$$\hat{y} = \hat{w}_0 + \hat{w}_1 \mathbf{x}[1] + \dots + \hat{w}_j \mathbf{x}[j] + \dots + \hat{w}_d \mathbf{x}[d]$$

fix fix fix fix



Fitting D-dimensional curves





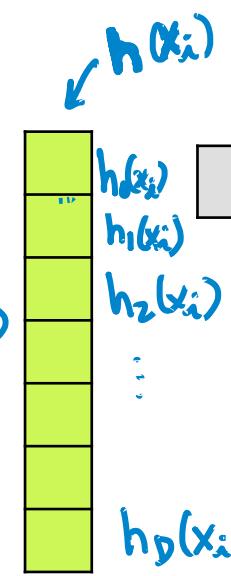
Step 1:

Rewrite the regression model

Rewrite in matrix notation

For observation i

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

y_i =  

$$\begin{aligned} y_i &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i \\ &= w^T h(x_i) + \epsilon_i \end{aligned}$$

$w^T h(x_i)$

$$\begin{aligned} y_i &= \begin{bmatrix} h_0(x_i) & h_1(x_i) & \dots & h_D(x_i) \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix} + \epsilon_i \\ &= h_0(x_i) w_0 + h_1(x_i) w_1 + \dots + h_D(x_i) w_D + \epsilon_i \end{aligned}$$

w

w_0 w_1 \vdots w_D

Rewrite in matrix notation

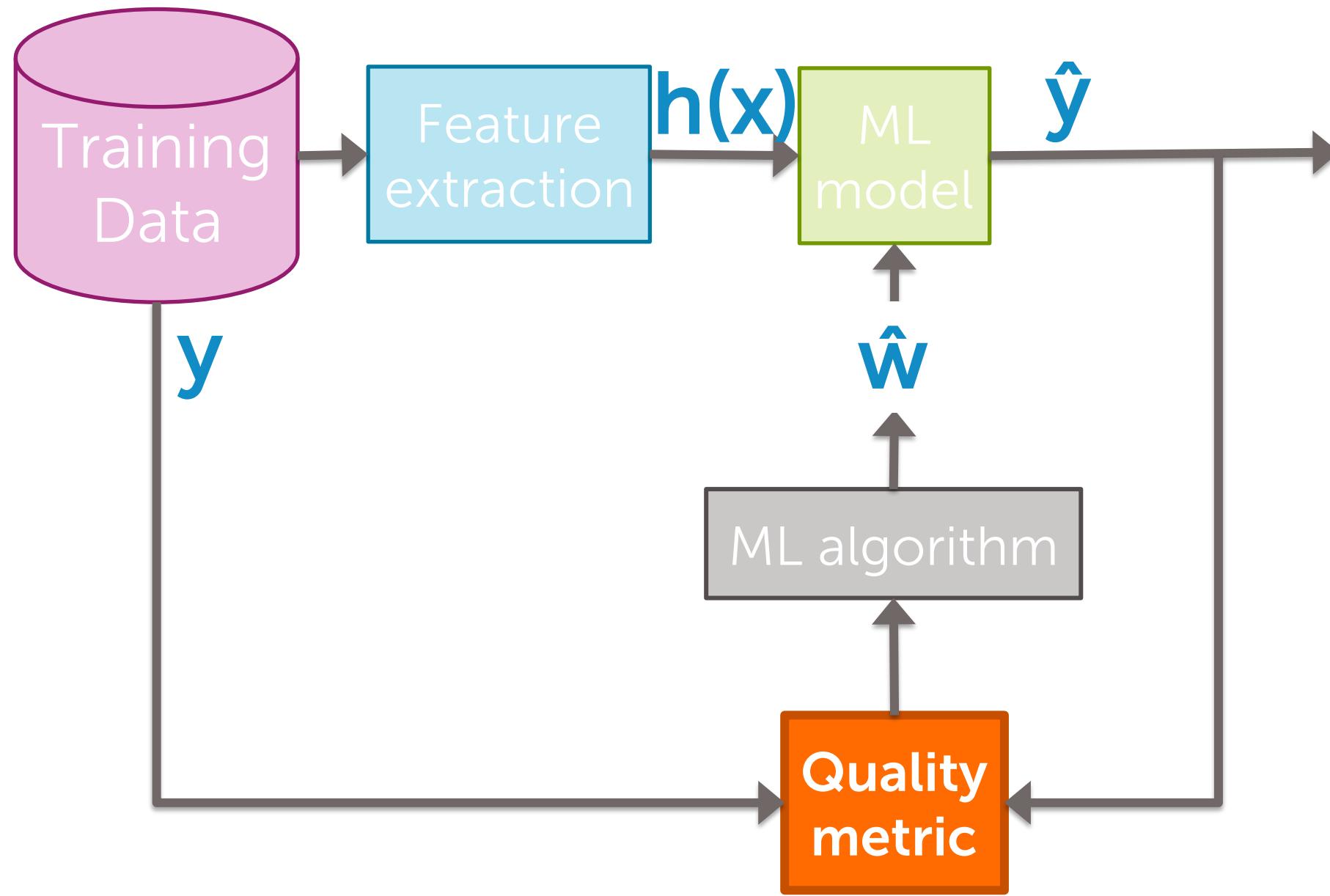
For all observations together

$$\begin{matrix} \mathbf{y} \\ \mathbf{y}_1 \\ \mathbf{y}_2 \\ \mathbf{y}_3 \\ \vdots \\ \vdots \\ \mathbf{y}_N \end{matrix} = \mathbf{H} \begin{matrix} \mathbf{w} \\ w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{matrix} + \begin{matrix} \mathbf{\epsilon} \\ \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \vdots \\ \epsilon_N \end{matrix} \Rightarrow \boxed{\mathbf{y} = \mathbf{H}\mathbf{w} + \mathbf{\epsilon}}$$

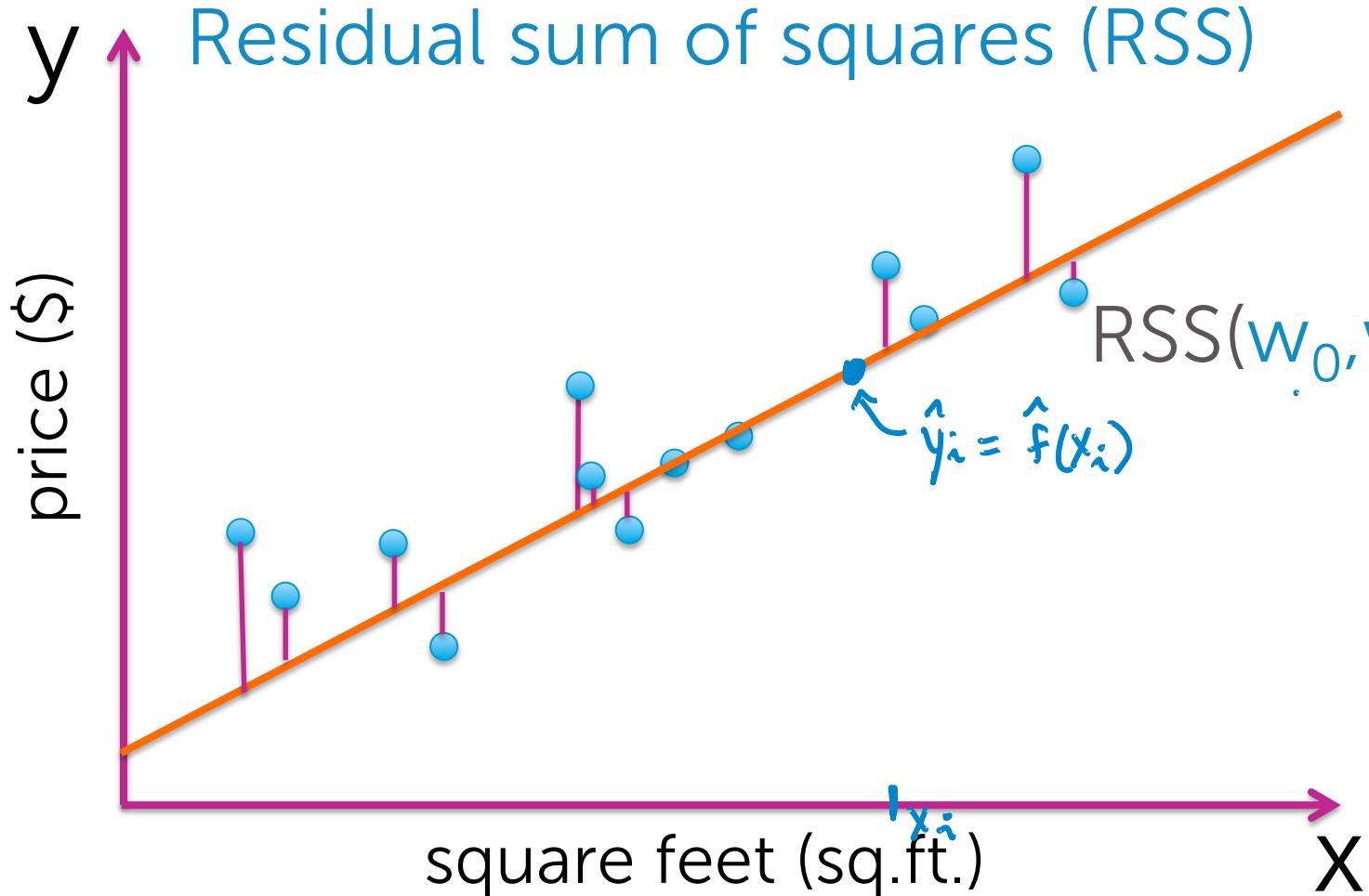
Diagram illustrating the matrix notation for a linear regression model. A vertical vector \mathbf{y} on the left is shown as a stack of observations $y_1, y_2, y_3, \dots, y_N$. To its right is a large green matrix labeled \mathbf{H} , representing the feature matrix. Above \mathbf{H} , a blue arrow points from \mathbf{y} to the first column of \mathbf{H} with the label $h^T(x_1)$. To the right of \mathbf{H} is a vertical vector \mathbf{w} representing the weight vector, with components $w_0, w_1, w_2, \dots, w_D$. To the right of \mathbf{w} is a plus sign (+). To the right of the plus sign is a vertical vector $\mathbf{\epsilon}$ representing the error term, with components $\epsilon_1, \epsilon_2, \epsilon_3, \dots, \epsilon_N$. An arrow points from the equation to the final boxed result $\boxed{\mathbf{y} = \mathbf{H}\mathbf{w} + \mathbf{\epsilon}}$.

Step 2:

Compute the cost



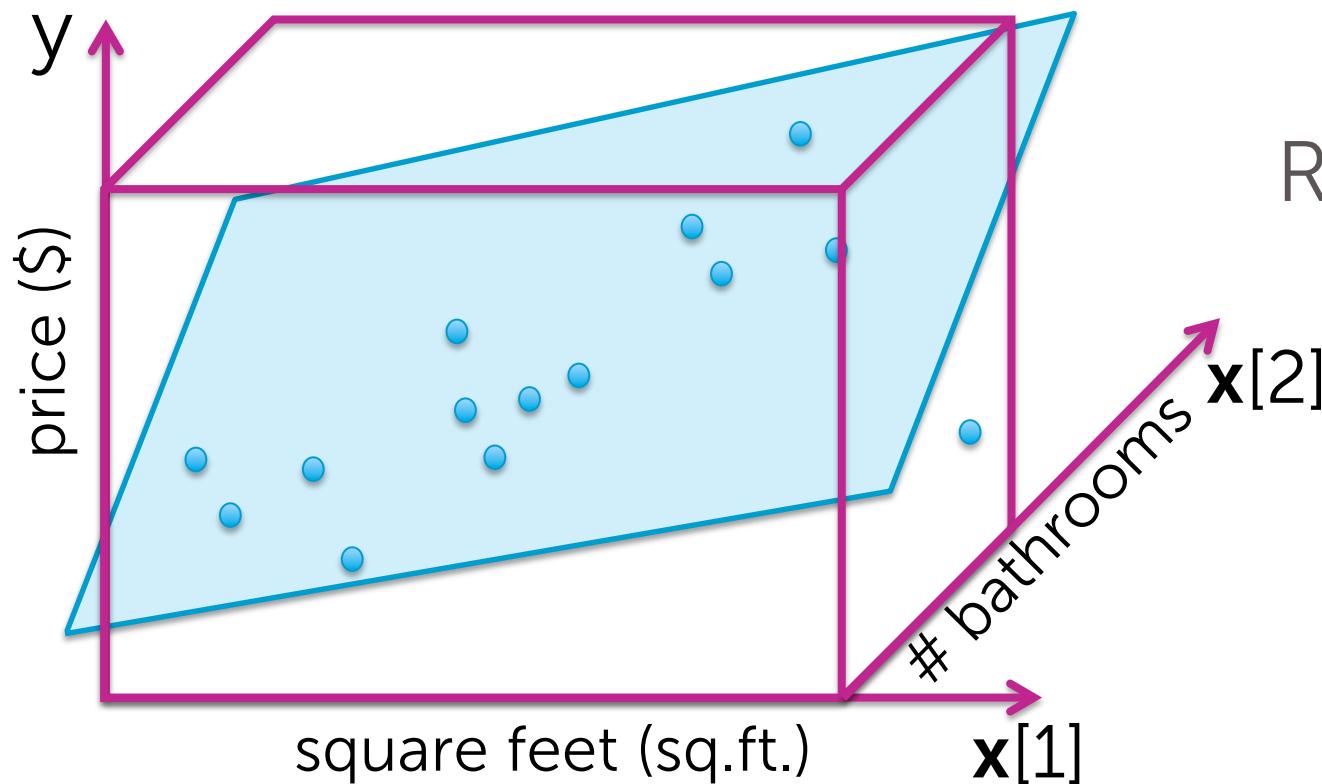
"Cost" of using a given line



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

$\hat{y}_i(w_0, w_1)$

RSS for multiple regression



$$\text{RSS}(\underline{\mathbf{w}}) = \sum_{i=1}^N (y_i - \hat{y}_i(\underline{\mathbf{w}}))^2$$
$$\hat{y}_i = \begin{bmatrix} h_0(x_i) & h_1(x_i) & \dots & h_D(x_i) \end{bmatrix} \underline{\mathbf{w}}$$
$$\underline{\mathbf{w}} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

RSS in matrix notation

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

Why? (part 1)

$$\begin{matrix} \hat{\mathbf{y}} \\ \vdots \\ \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{matrix} = \mathbf{H} \begin{matrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{matrix}$$

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{w}$$
$$(\mathbf{y} - \tilde{\mathbf{H}}\mathbf{w}) = (\mathbf{y} - \hat{\mathbf{y}}) = \begin{bmatrix} \text{residual}_1 \\ \text{residual}_2 \\ \vdots \\ \text{residual}_N \end{bmatrix}$$

residual_i = $y_i - \hat{y}_i$

RSS in matrix notation

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

Why? (part 2)

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

residual ₁
residual ₂
residual ₃
...
residual _N

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

Step 3:

Take the gradient

Gradient of RSS

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

Why? By analogy to 1D case:

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

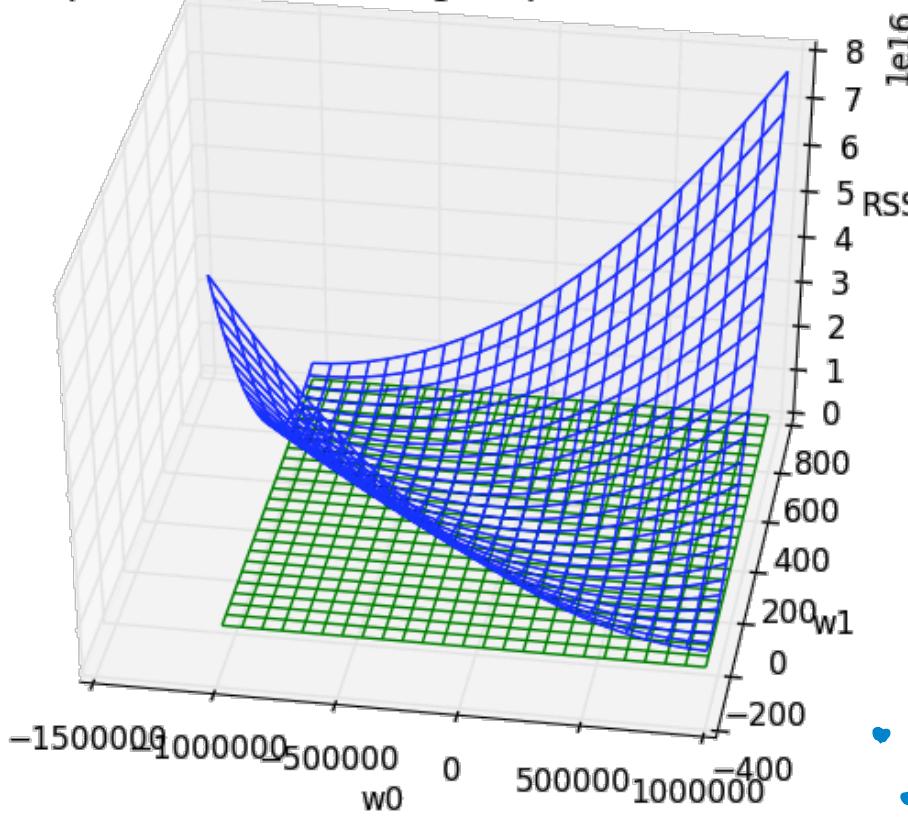
\nearrow
scalars

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

Step 4, Approach 1:
Set the gradient = 0

Closed-form solution

3D plot of RSS with tangent plane at minimum



$$\begin{matrix} \bullet & A^{-1}A = I \\ \bullet & IV = V \\ \bullet & IV = V \end{matrix}$$

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

Solve for \mathbf{w} :

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

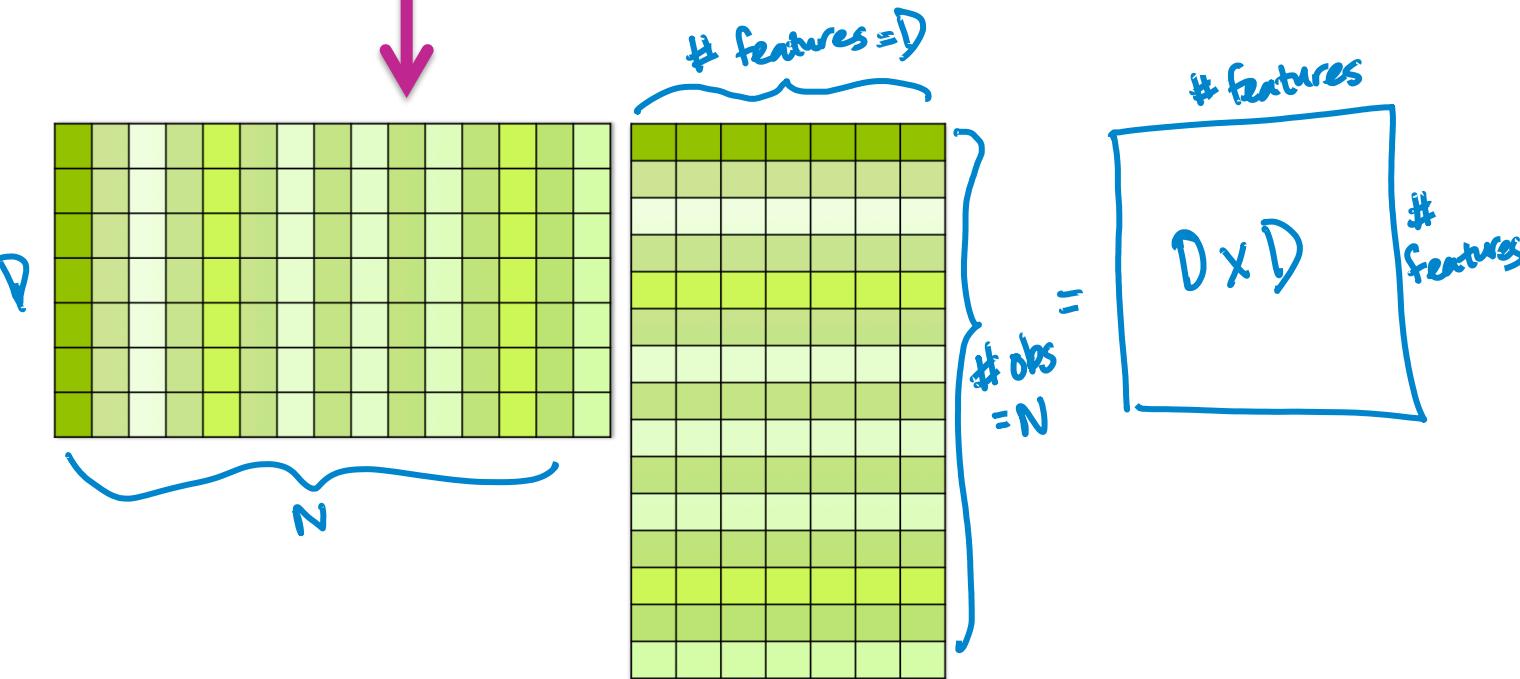
$$\mathbf{H}^T \mathbf{H} \hat{\mathbf{w}} = \mathbf{H}^T \mathbf{y}$$

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

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

Closed-form solution

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



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

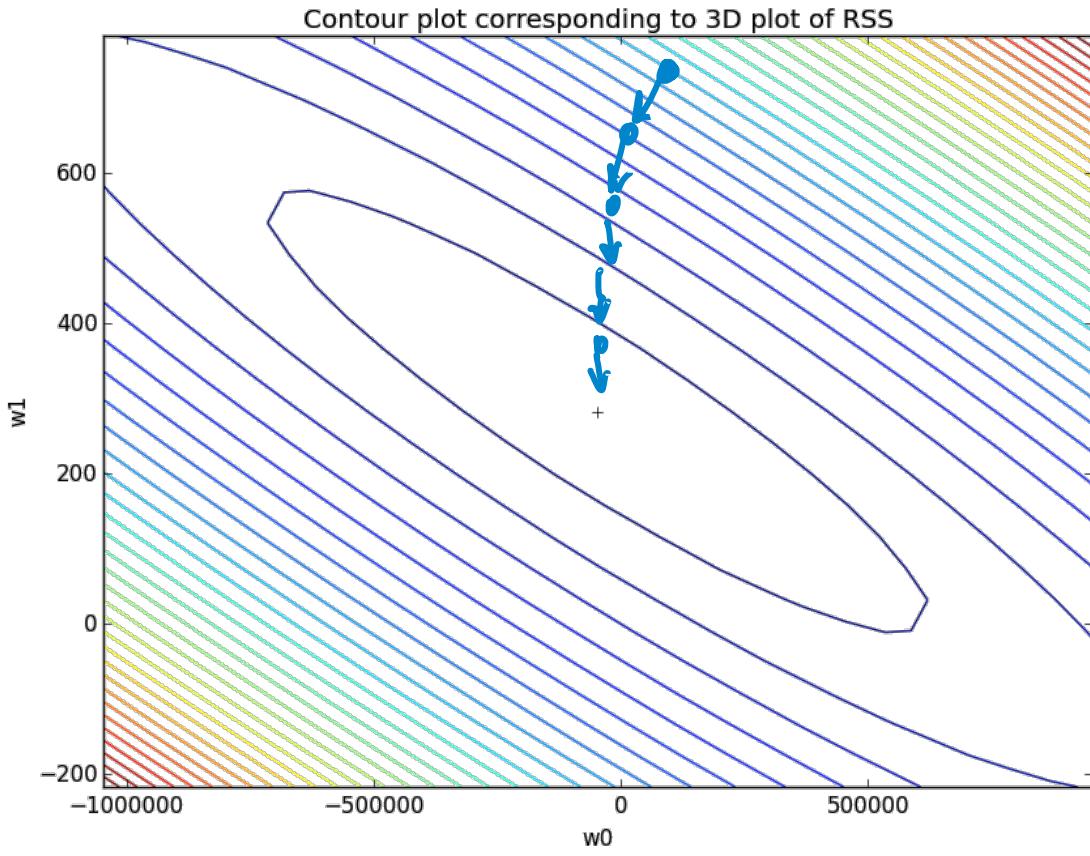
really,
of linearly
ind. observations

Complexity of inverse:

$O(D^3)$

Step 4, Approach 2: Gradient descent

Gradient descent



while not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}^{(t)})$$
$$= \mathbf{w}^{(t)} + 2\eta \mathbf{H}^T (\mathbf{y} - \mathbf{H}\mathbf{w}^{(t)})$$

$\hat{\mathbf{y}}(\mathbf{w}^{(t)})$

Feature-by-feature update

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

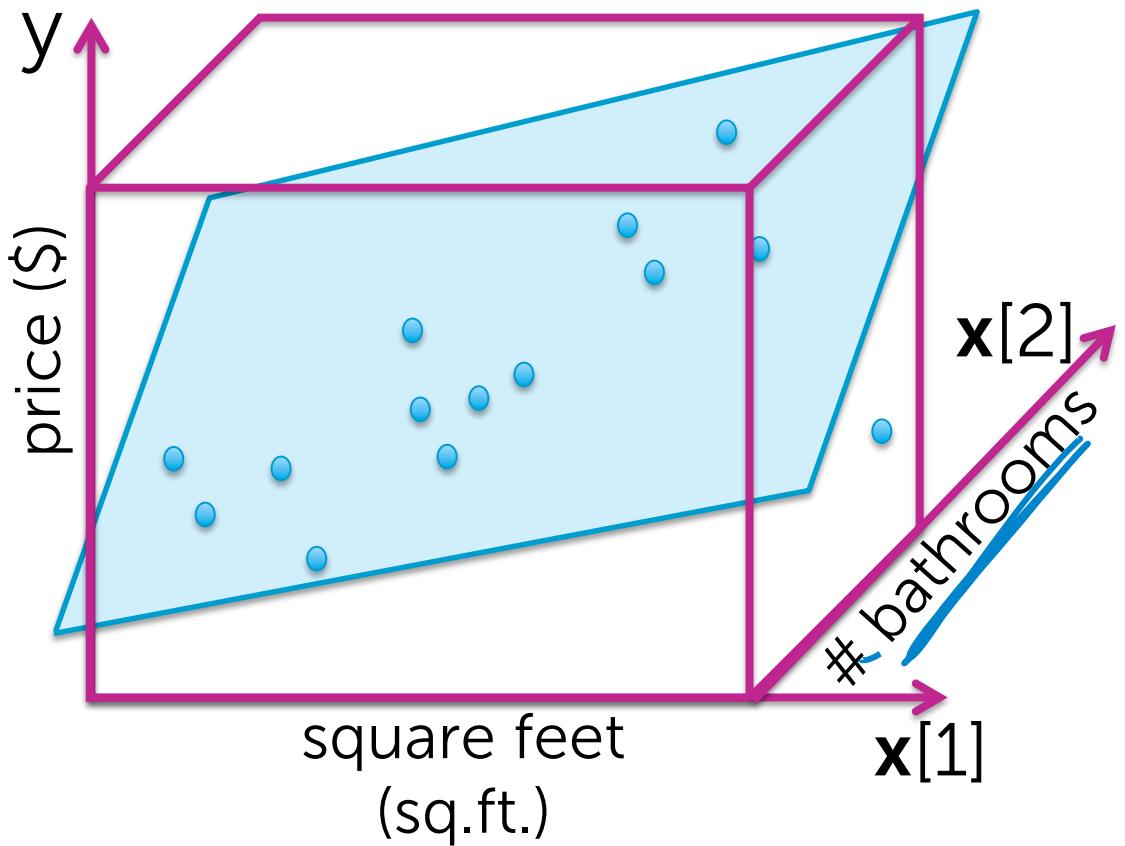
Partial with respect to w_j .

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

Update to j^{th} feature weight:

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

Interpreting elementwise

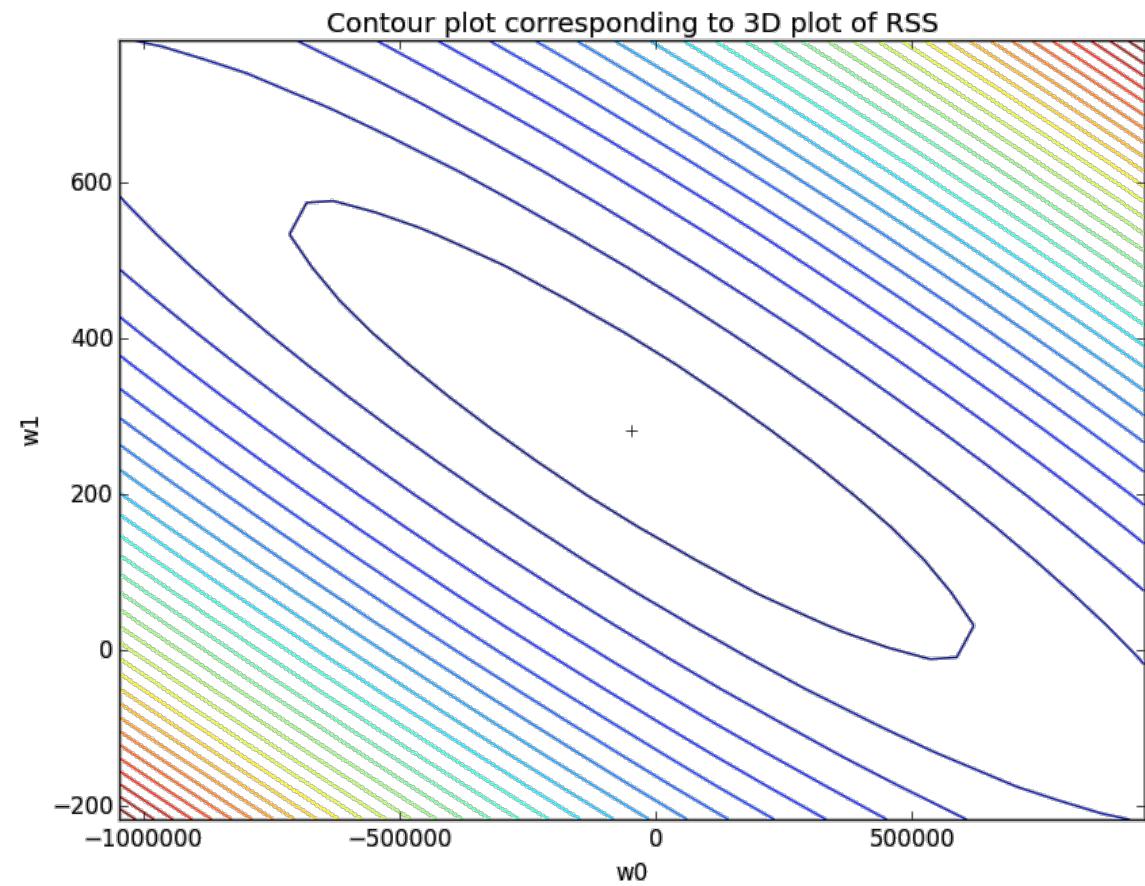


Update to j^{th} feature weight:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + 2\eta \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$$

If underestimating impact of #bath ($\hat{w}_j^{(t)}$ is too small)
then $(y_i - \hat{y}_i(w^{(t)}))$ on average
weighted by #bath will be positive
 $\Rightarrow w_j^{(t+1)} > w_j^{(t)}$ (increase)

Summary of gradient descent for multiple regression



```
init  $\mathbf{w}^{(1)} = \mathbf{0}$  (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  $\leftarrow t + 1$ 
```

An extremely useful algorithm

Summary for multiple linear regression

What you can do now...

- Describe polynomial regression
- Detrend a time series using trend and seasonal components
- Write a regression model using multiple inputs or features thereof
- Cast both polynomial regression and regression with multiple inputs as regression with multiple features
- Calculate a goodness-of-fit metric (e.g., RSS)
- Estimate model parameters of a general multiple regression model to minimize RSS:
 - In closed form
 - Using an iterative gradient descent algorithm
- Interpret the coefficients of a non-featurized multiple regression fit
- Exploit the estimated model to form predictions
- Explain applications of multiple regression beyond house price modeling