

Linear Regression: Model and Algorithms

CSE 446

Slides by Emily Fox (with minor modifications)

Presented by Anna Karlin

April 3, 2019

Linear regression: The model

How much is this house worth?



How much is this house worth?



Data



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



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



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



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



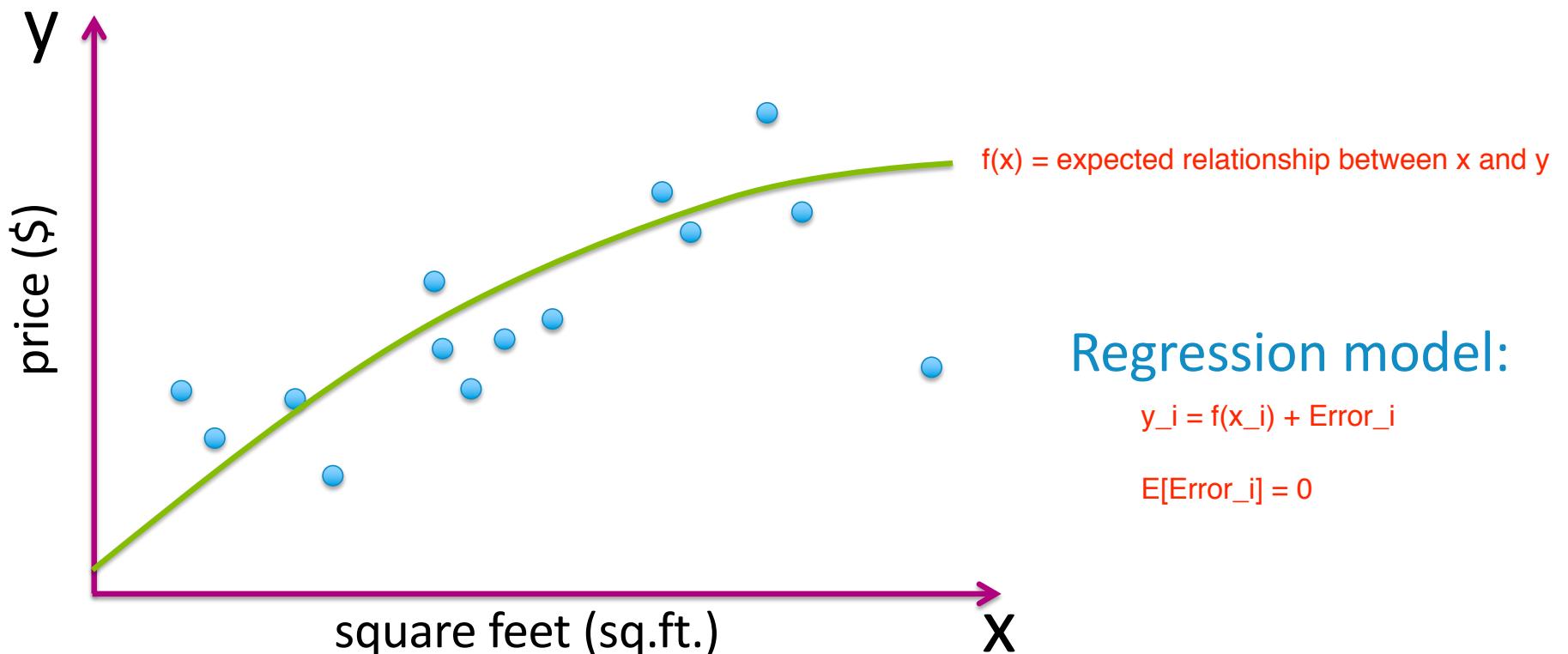
$(x_5 = \text{sq.ft.}, y_5 = \$)$

⋮

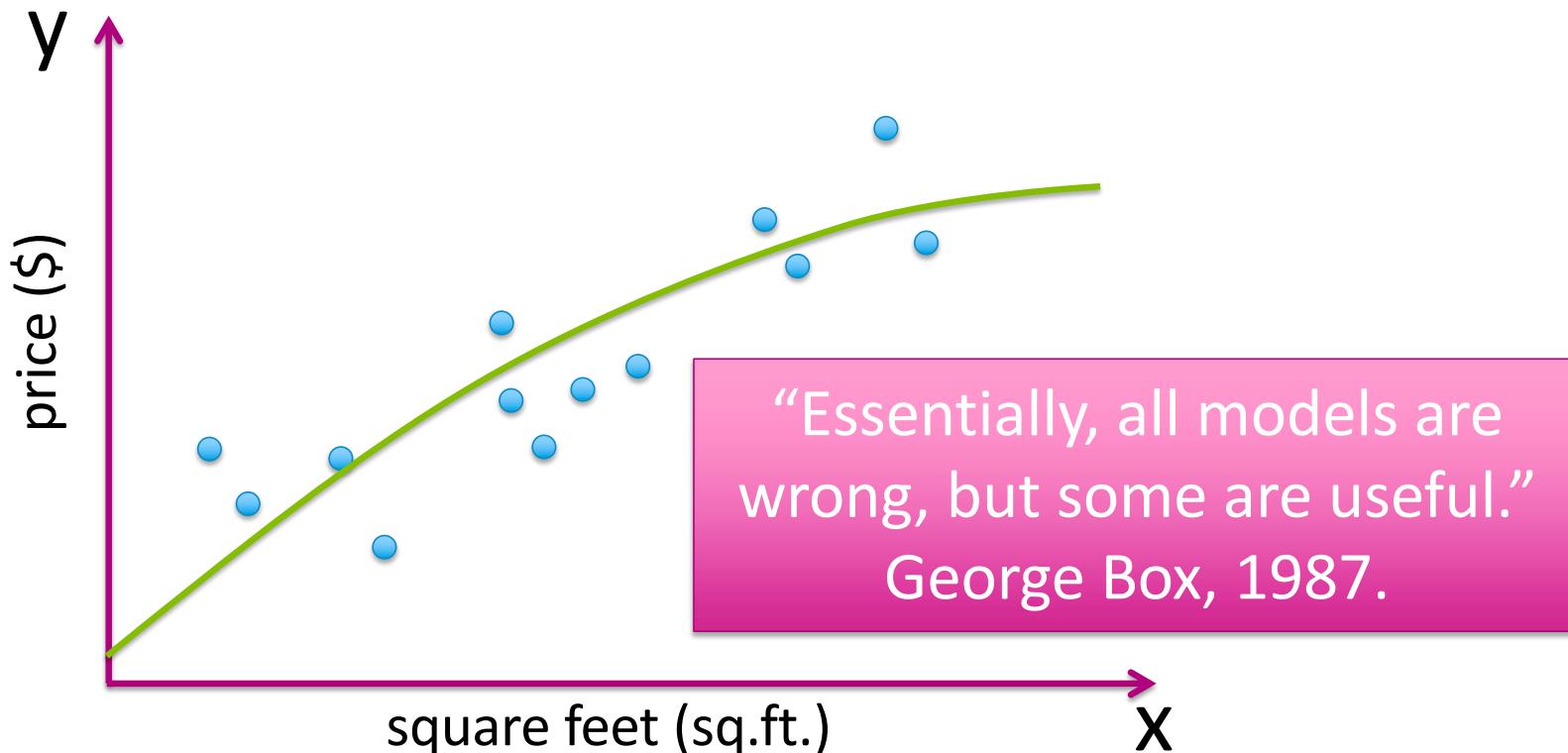
Input vs. Output:

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

Model – How we *assume* the world works



Model – How we *assume* the world works



Process

- Chose a model

a family of functions that we feel best captures the dynamics of the system.

$\{ f_w(w) \mid (w_1, w_2, \dots, w_n) \in R^n \} \rightarrow$ a function parameterized by parameters w_1, w_2, \dots, w_n in R

- Find the best fit to our data set.

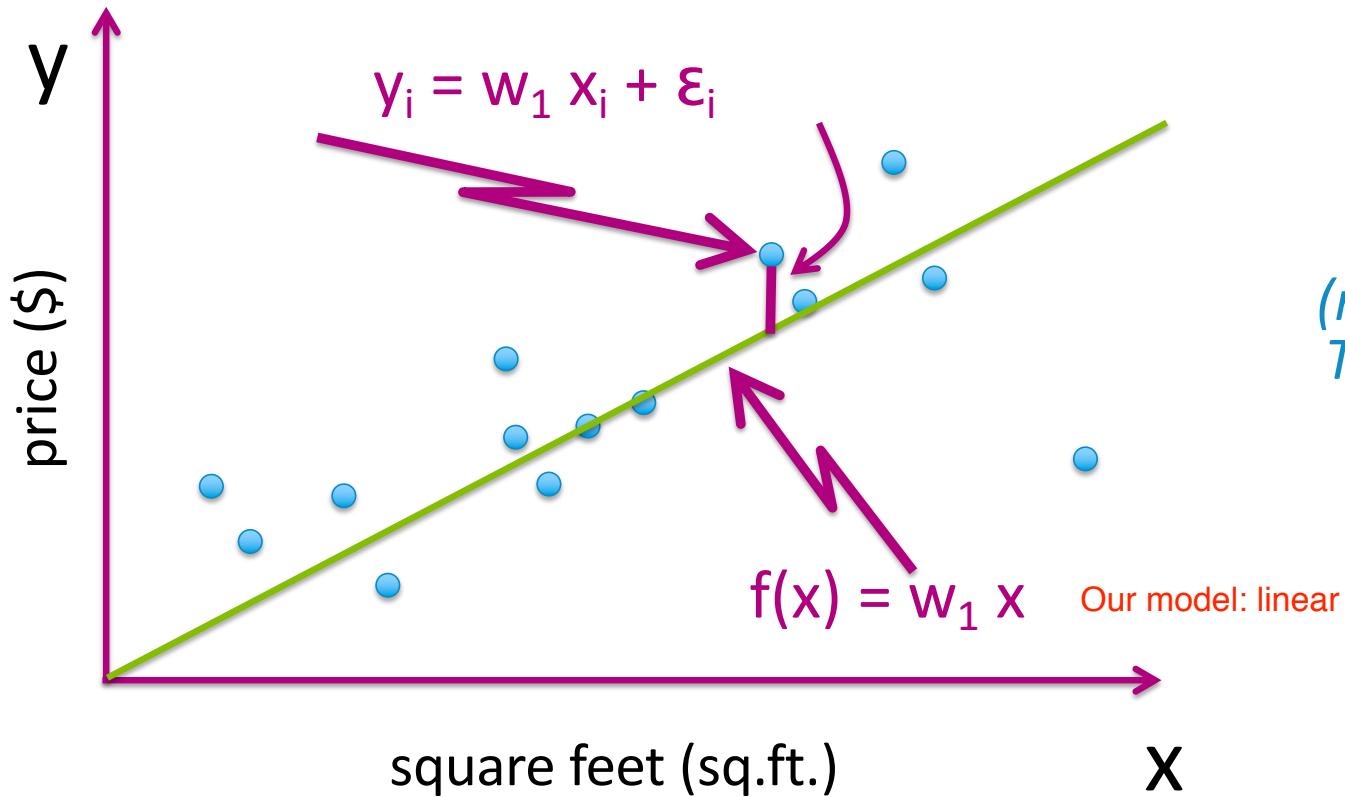
Which function in family is best? What are the best choice of parameters?

Find w^* that meets some sort of condition we find ideal in family of function; i.e. minimizes or maximizes some quantity.

- Use the fitted function to make future predictions.

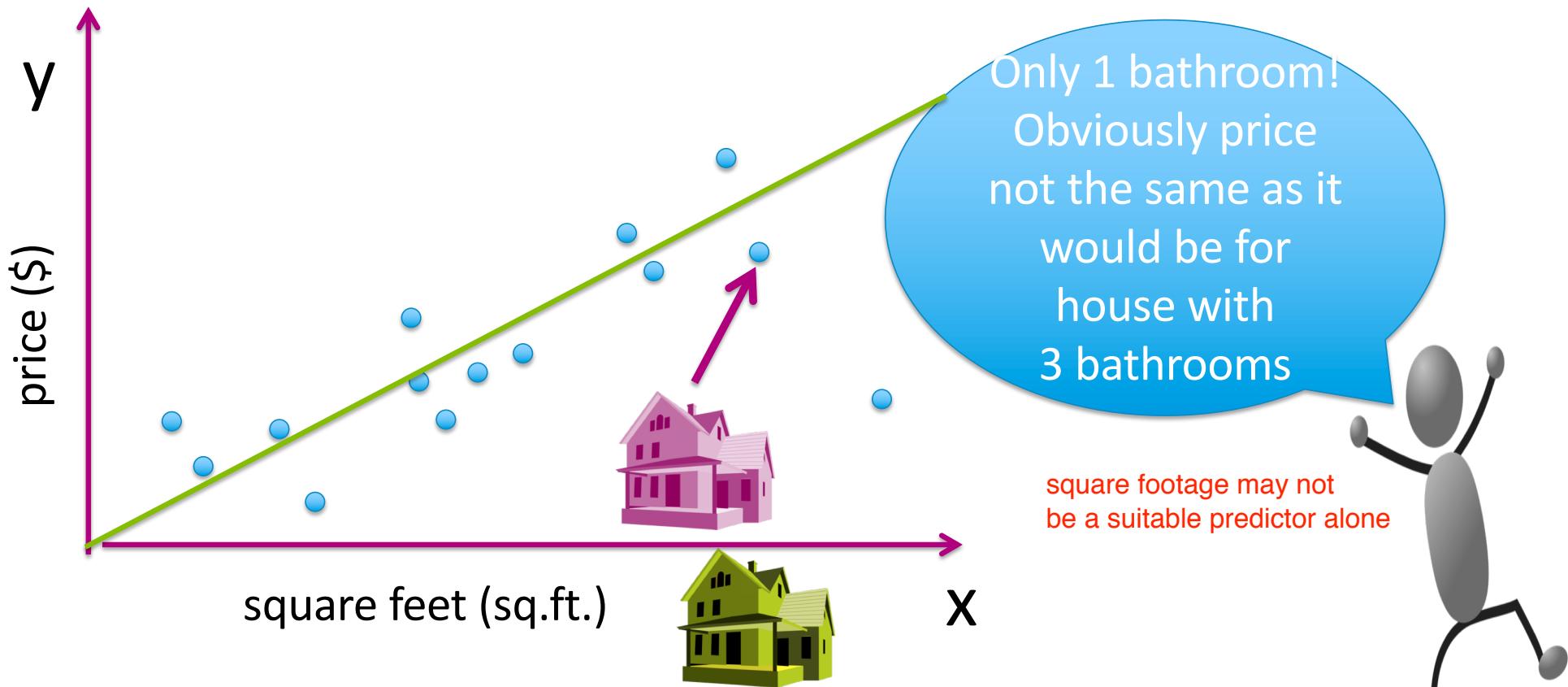
`y_future = f(x_future) + Error`

Simple linear regression model



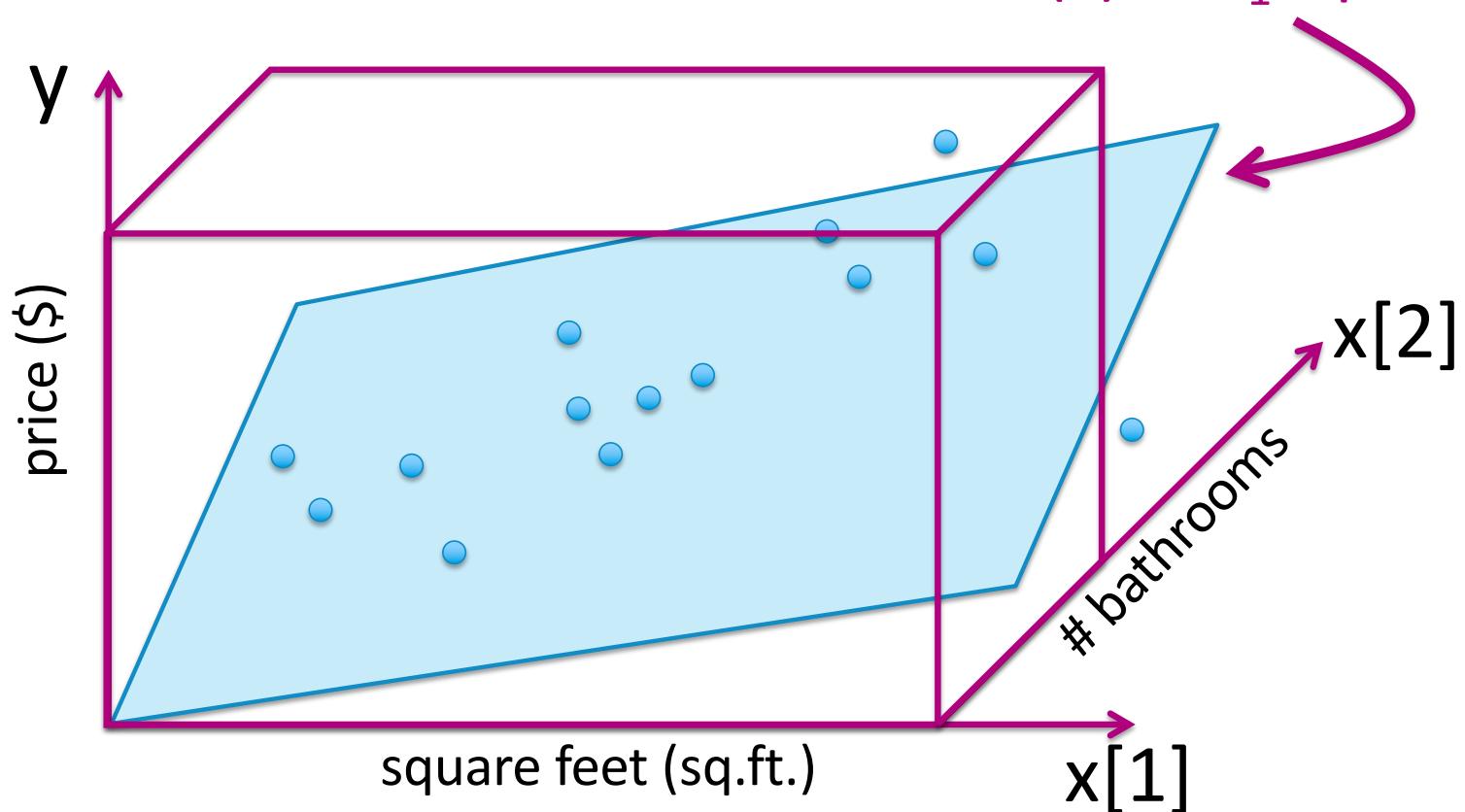
w_1 is the parameter
(regression coefficient)
That we want to learn
from our data set.

Predictions just based on house size



Add more inputs

$$f(x) = w_1 \text{sq.ft} + w_2 \#\text{bath}$$



Many possible inputs

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

General notation

Output: y scalar

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

e.g., $x[1] = \text{sq. ft}$, $x[2] = \#\text{baths}$ and so on.

Notational conventions:

training set: $\{(x_i, y_i)\}_{i=1..n}$

x_i = input of i^{th} data point/observation (*vector*); y_i is output

$x_i[j]$ = j^{th} input of i^{th} data point (*scalar*)

n = number of observations; d = number of input features

Generic linear regression model

Model: Given feature vector $\mathbf{x}_i = (x_i[1], x_i[2], \dots, x_i[d])$

$$y_i = w_1 x_i[1] + w_2 x_i[2] + \dots + w_d x_i[d] + \varepsilon_i$$

$$= \sum_{j=1}^d w_j x_i[j] + \varepsilon_i = \mathbf{w}^\top \mathbf{x}_i + \varepsilon_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i$$

feature 1 = $x[1]$ = sq. ft.

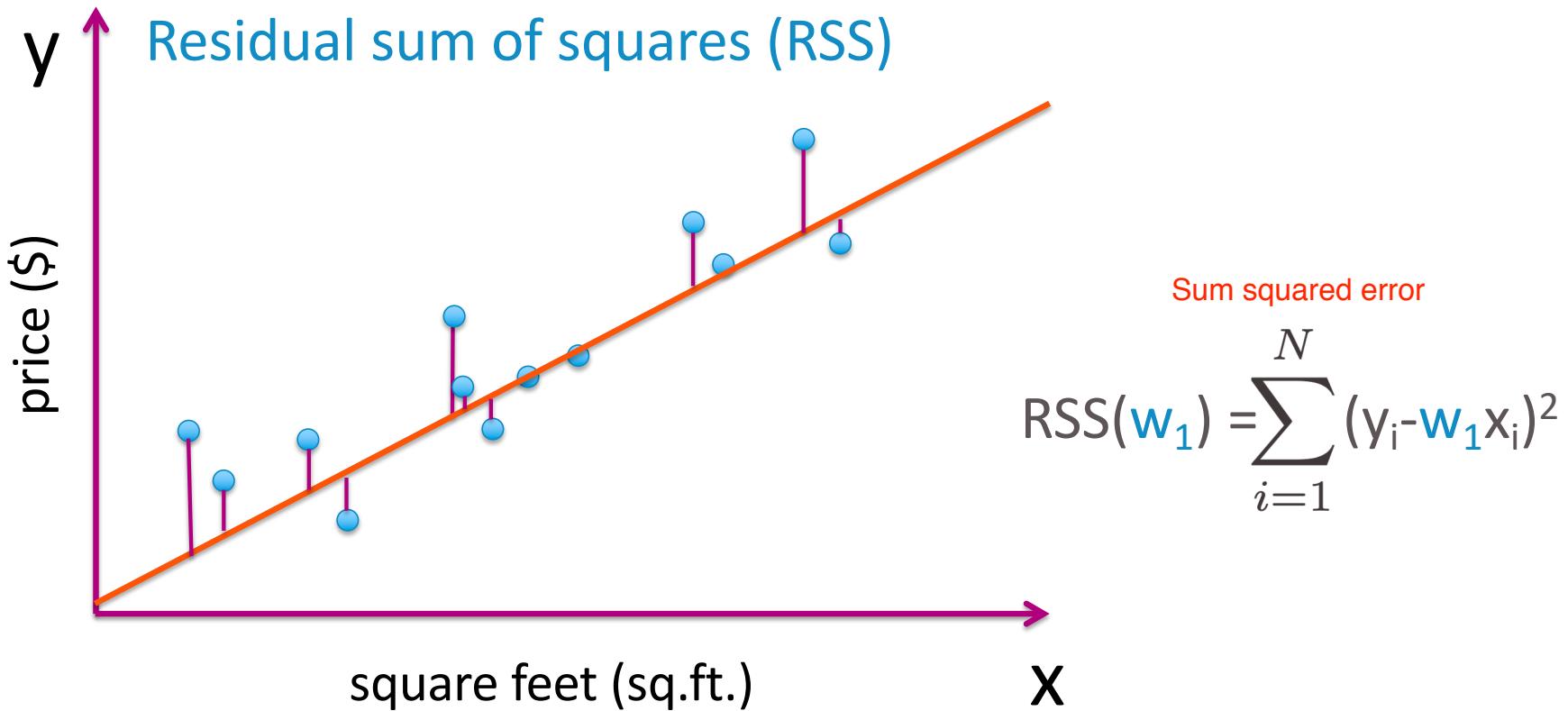
feature 2 = $x[2]$ = #bath

...

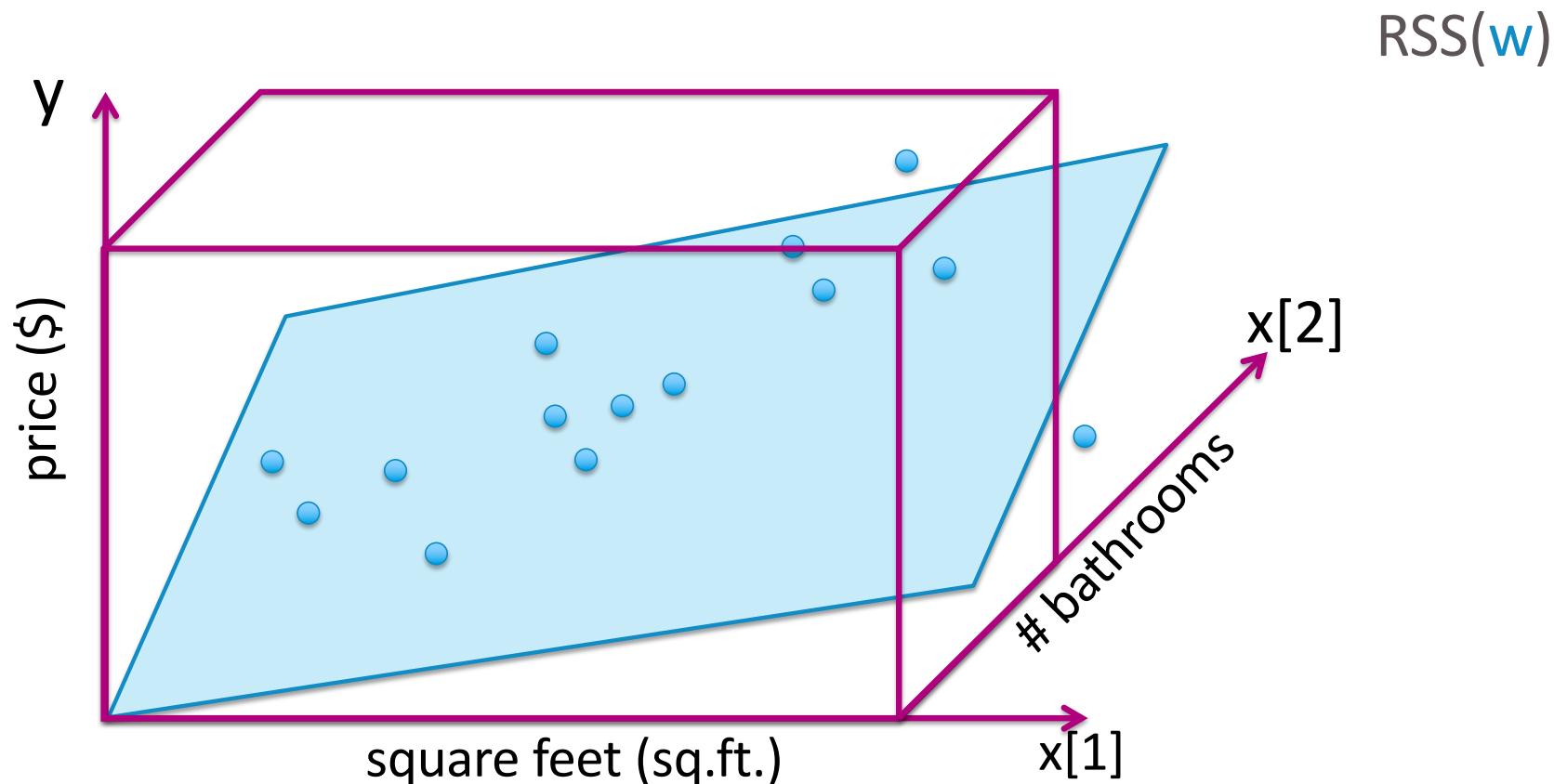
feature d = $x[d]$ = lot size

Fitting the linear regression model

“Cost” of using a given line



RSS for multiple regression



Rewrite in matrix notation

For all observations together

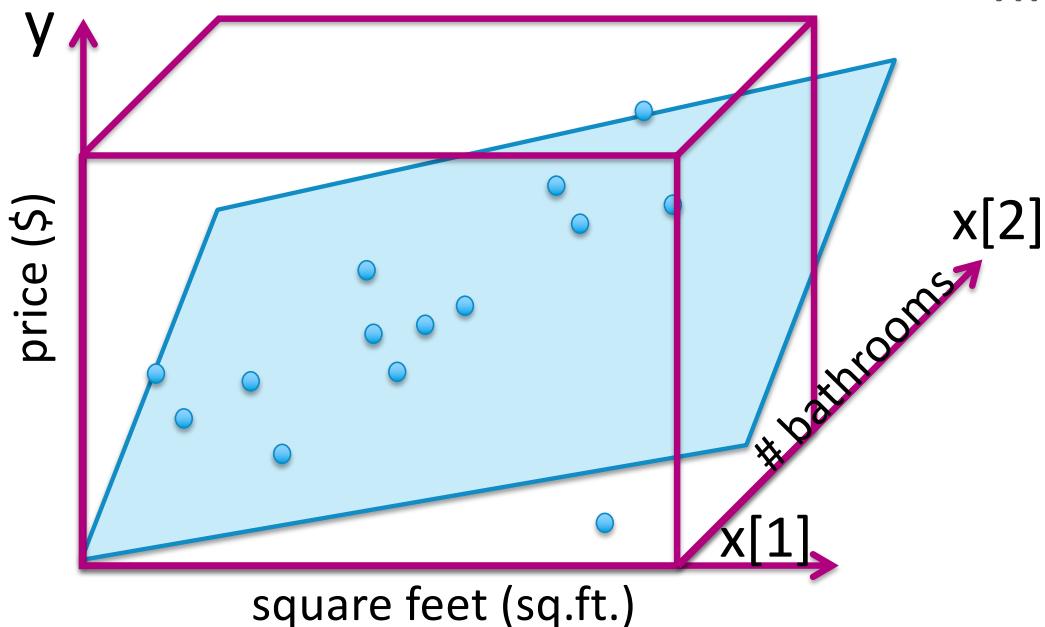
$$\text{RSS}(w) = \sum_{i=1}^n (y_i - x_i^T w)^2 = (y - Xw)^T (y - Xw)$$

Diagram illustrating the components of the equation:

- y**: A vertical vector containing observations y_1, y_2, \dots, y_n .
- x_i**: A horizontal vector representing the features for observation i , labeled $[1] [2]$ at the top.
- X**: A matrix where each row is a feature vector x_i , labeled $f(x_i)$ below the matrix.
- w**: A vertical vector containing weights w_1, w_2, \dots, w_d .

The diagram shows the subtraction operation $y - Xw$ as a horizontal line with a minus sign between the vectors y and Xw .

RSS for multiple regression



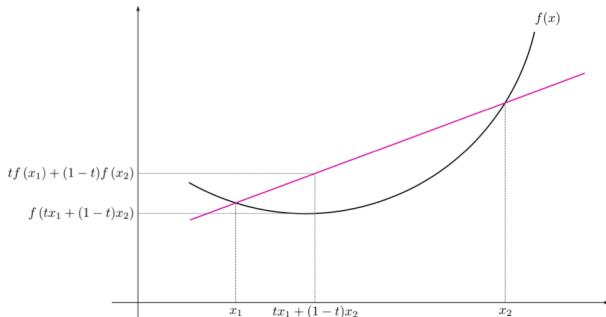
Objective: Find the best fit, i.e., find the \mathbf{w} that minimizes

$$\text{RSS}(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w})$$

Interlude: Optimization (and convex functions)

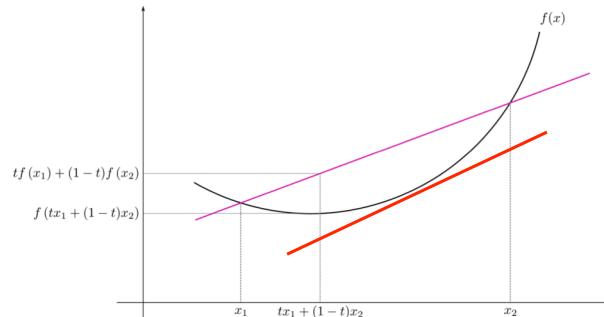
Minimizing a univariate function

Convex functions $f''(x) > 0$ for all x ,



21

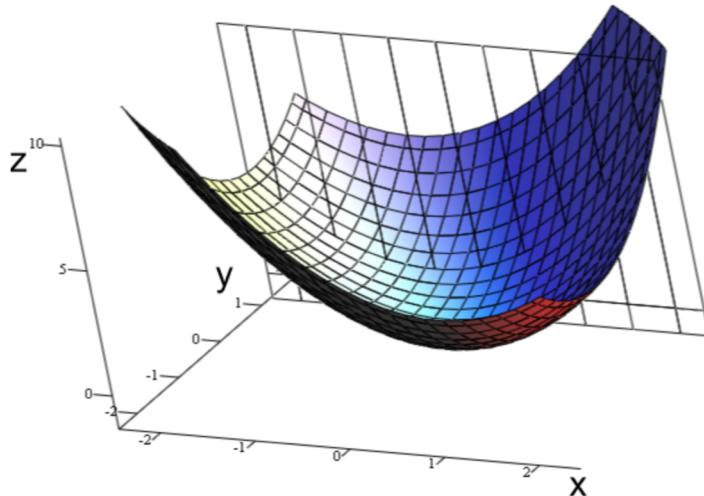
line between any two points always overestimates $f(x)$ between those two points



at any point the tangent plane will always underestimate the function values at other points

Minimizing a multivariate function

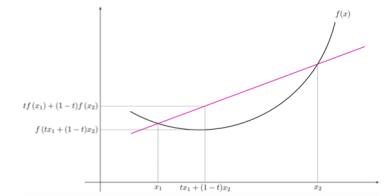
Convex functions



gradient $f = 0 \rightarrow$ global minimum

Any two points on surface the line between them overestimates height of surface.

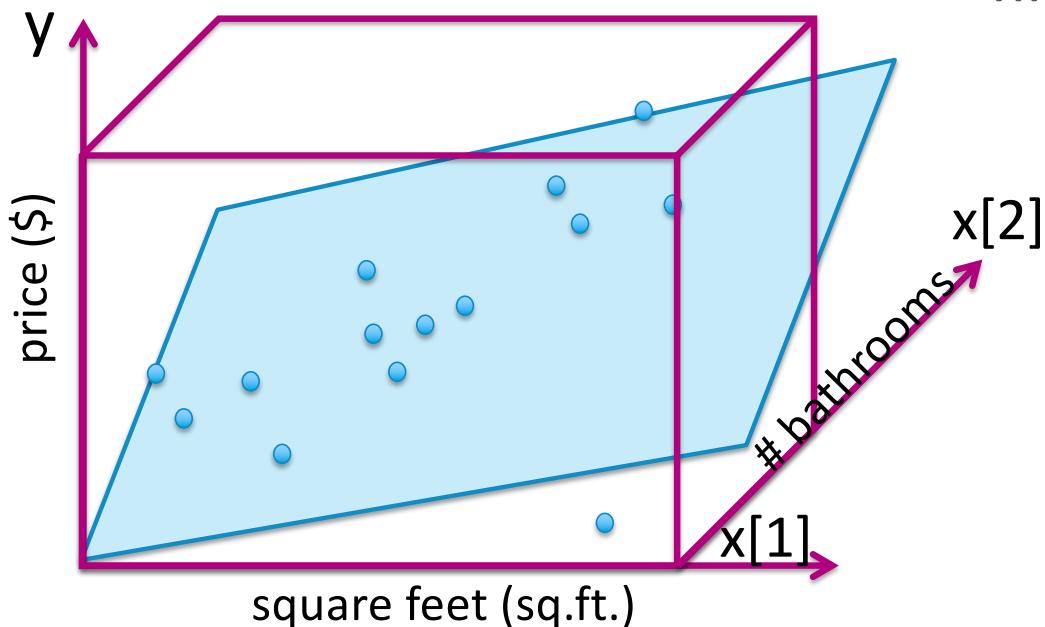
Tangent plane always underestimates surface and only touches at tangent point



$f(w_0) + (w-w_0)^T \text{gradient } f(w_0) \leq f(w)$ for all w
tangent plane at w_0

Back to our specific optimization problem

Two ways to solve our optimization problem



Objective: Find the best fit, i.e., find the \mathbf{w} that minimizes

$$\text{RSS}(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$\text{RSS}(\mathbf{w})$ is a convex function of \mathbf{w}

Take the second partial derivative with respect to w_i and note that it is always > 0

1. Solve for $\nabla \text{RSS}(\mathbf{w}) = 0$

Gradient of RSS

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

Closed form solution

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

We want solution to $-2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) = 0$

Solution: "normal equations" $\mathbf{X}^\top \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}$

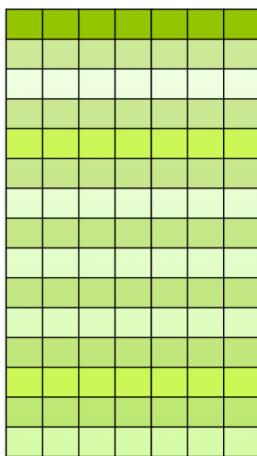
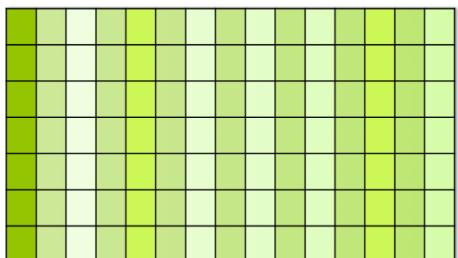
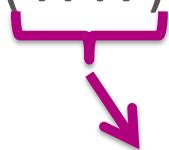
could solve this with an inverse, but it might exist, doing so is slow, and there are better ways of going about doing this. We know there is a solution

Geometric intuition for solution

Again, back to getting closed form solution.

Solution to normal equations $X^T X \hat{w} = X^T y$

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



Invertible if:

Complexity of inverse:

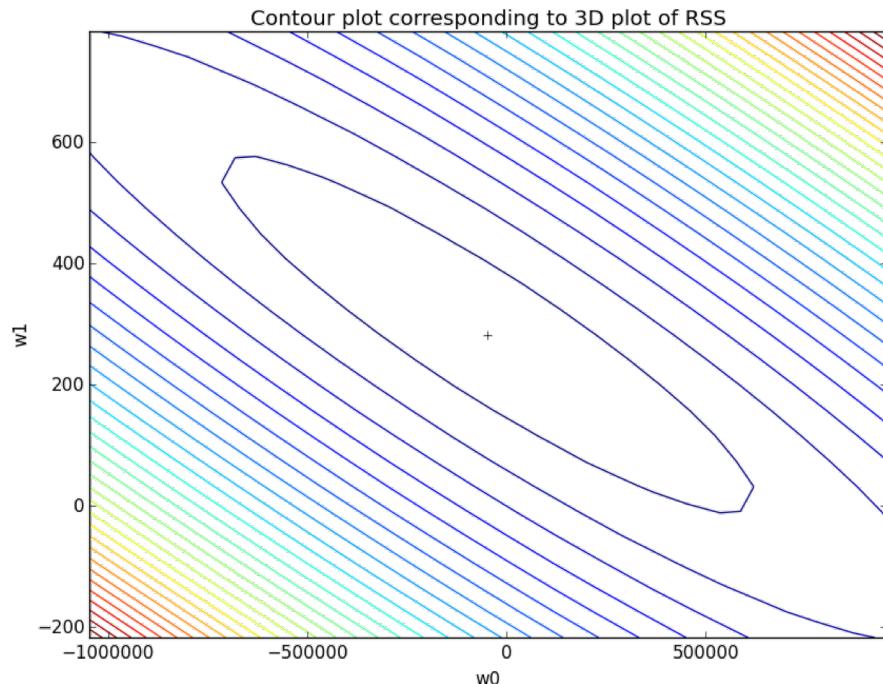
2. Gradient descent

Gradient Descent – univariate case

- Repeatedly move in direction that reduces the value of the function.

Gradient Descent – multivariate case

Gradient descent for linear regression: repeatedly move in direction of negative gradient

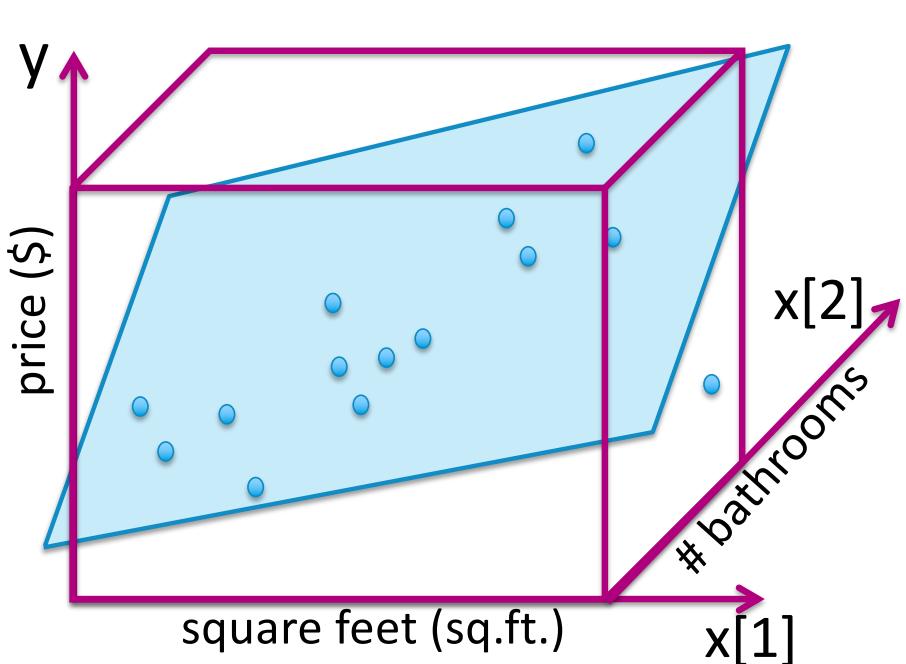


while not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}^{(t)})$$

$\underbrace{-2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}^{(t)})}$

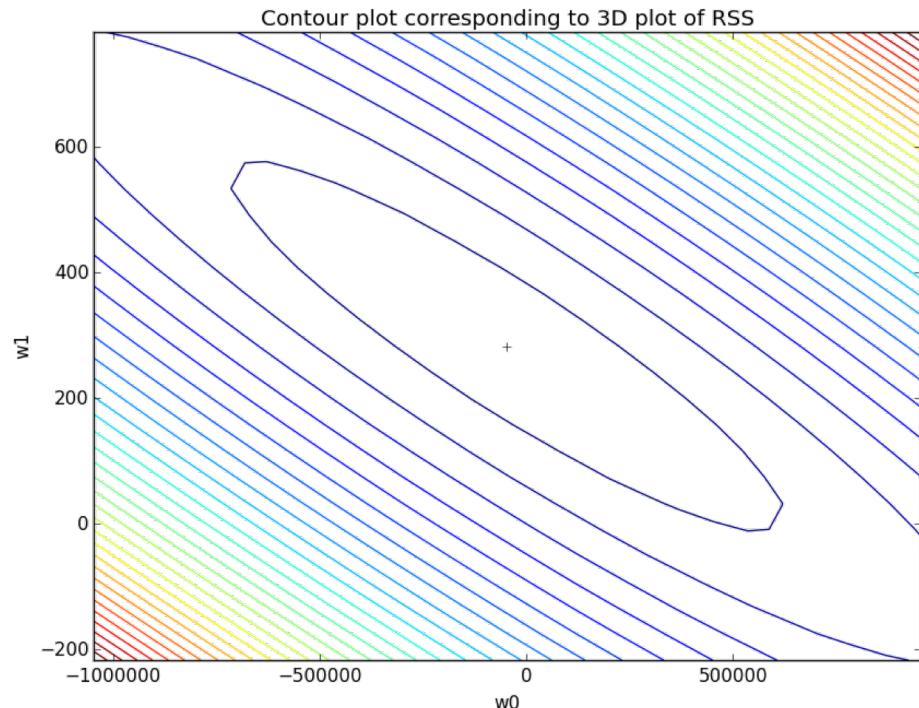
Interpreting elementwise



Update to j^{th} feature weight:

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

Summary of gradient descent for multiple regression



```
init  $w^{(1)}=0$  (or randomly, or smartly),  $t=1$ 
while  $\|\nabla \text{RSS}(w^{(t)})\| > \epsilon$ 
    for  $j=1,\dots,d$ 
         $\text{partial}[j] = -2 \sum_{i=1}^n x_i[j](y_i - \hat{y}_i(w^{(t)}))$ 
         $w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \text{ partial}[j]$ 
     $t \leftarrow t + 1$ 
```

Adding an intercept – “demeaning”

Once we have a fitted function

- We use it to predict the sales price for new houses, by plugging in square footage, number of bathrooms, etc for the new house \mathbf{x} whose sales price we want to predict.
- Prediction is:
- What if we want to allow for an intercept?

Handling an intercept (constant term)

Two step approach:

1. Show that if $\frac{1}{n} \sum_i \mathbf{x}_i = \mathbf{0}$ (*) then solution is simple.
2. Show how to transform, aka ``demean'' any linear regression problem so that (*) holds.

1. Show that if $\frac{1}{n} \sum_i \mathbf{x}_i = \mathbf{0}$ (*) then solution is simple.

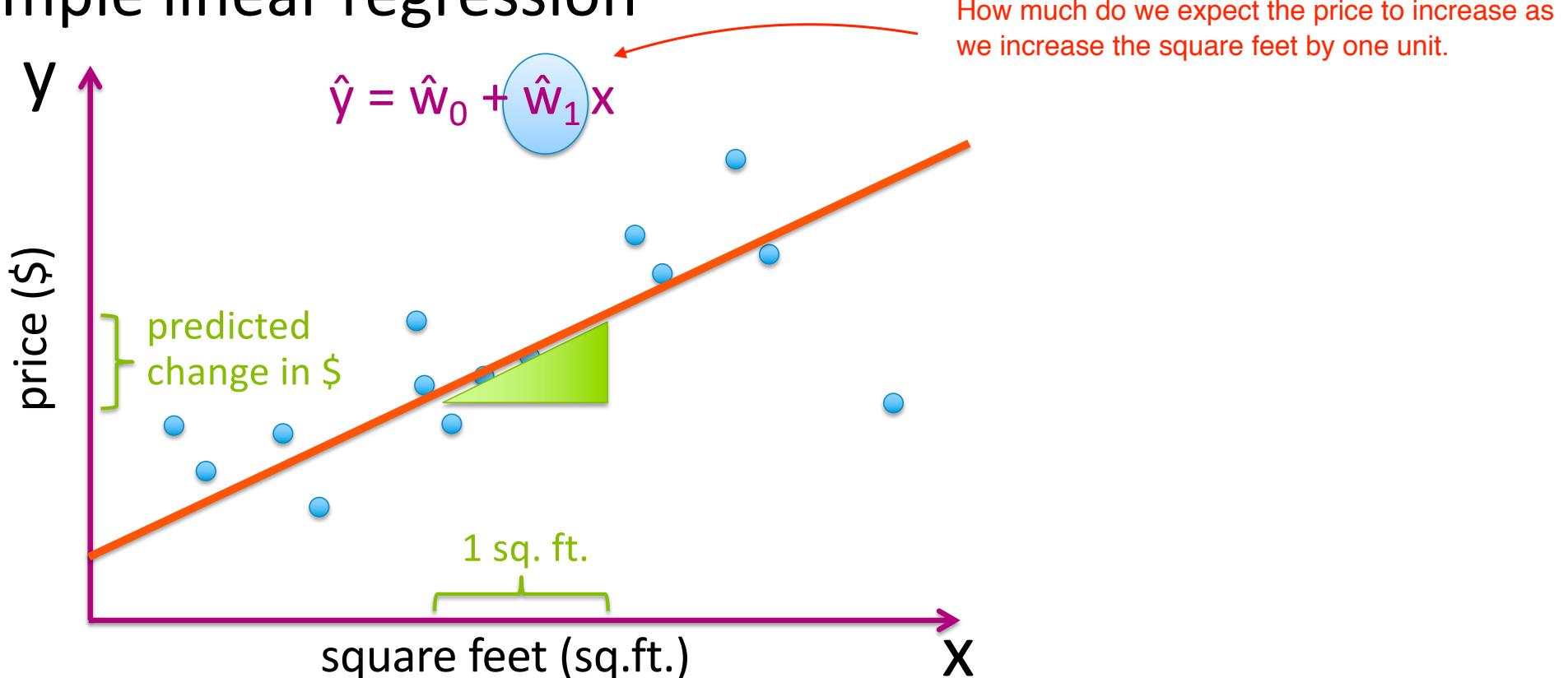
Same as saying that $\mathbf{X}^T \mathbf{1} = \mathbf{0}$.

2. Show how to transform, aka ``demean'' any linear regression problem so that (*) holds.

$$\frac{1}{n} \sum_i x_i = \mathbf{0} \quad (*)$$

Interpreting the fitted function

Interpreting the coefficients – Simple linear regression

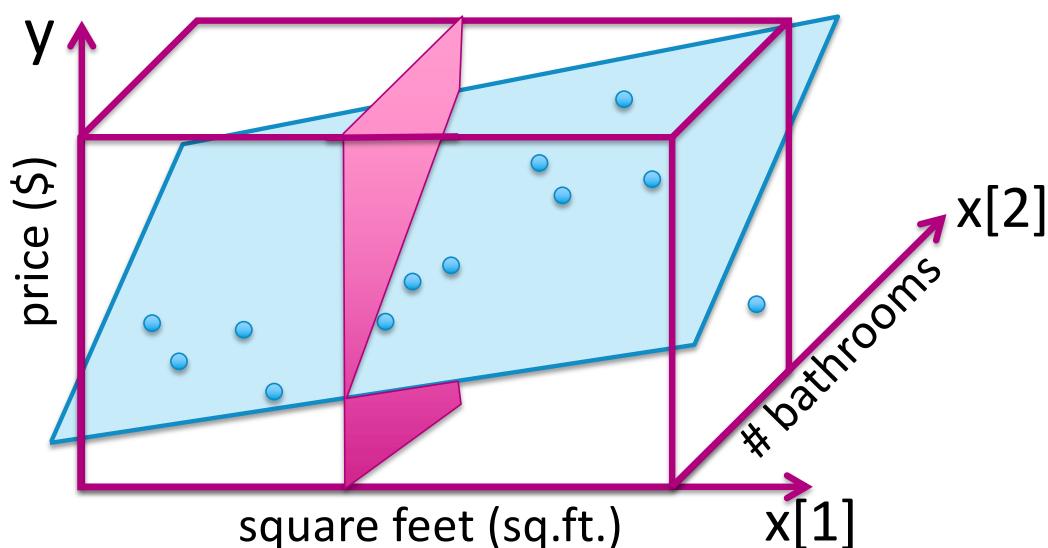


Interpreting the coefficients – Two linear features

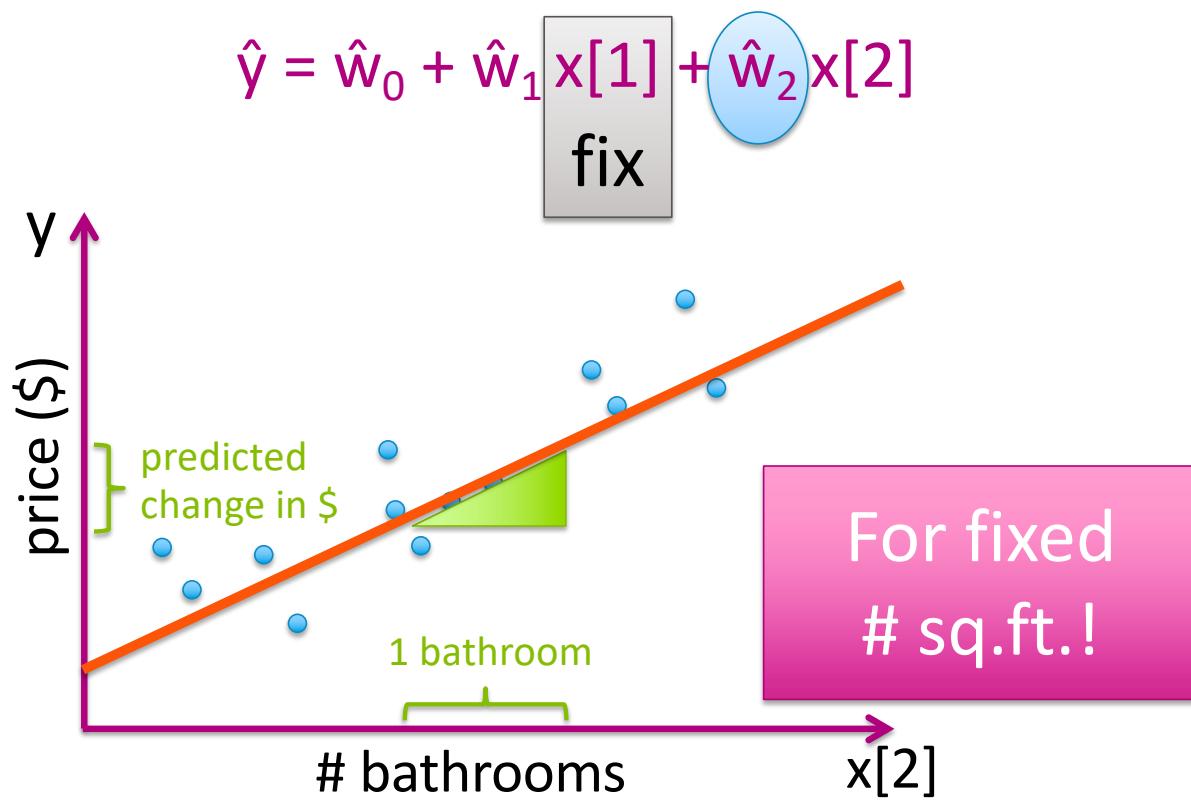
$$\hat{y} = \hat{w}_0 + \hat{w}_1 x[1] + \hat{w}_2 x[2]$$

fix

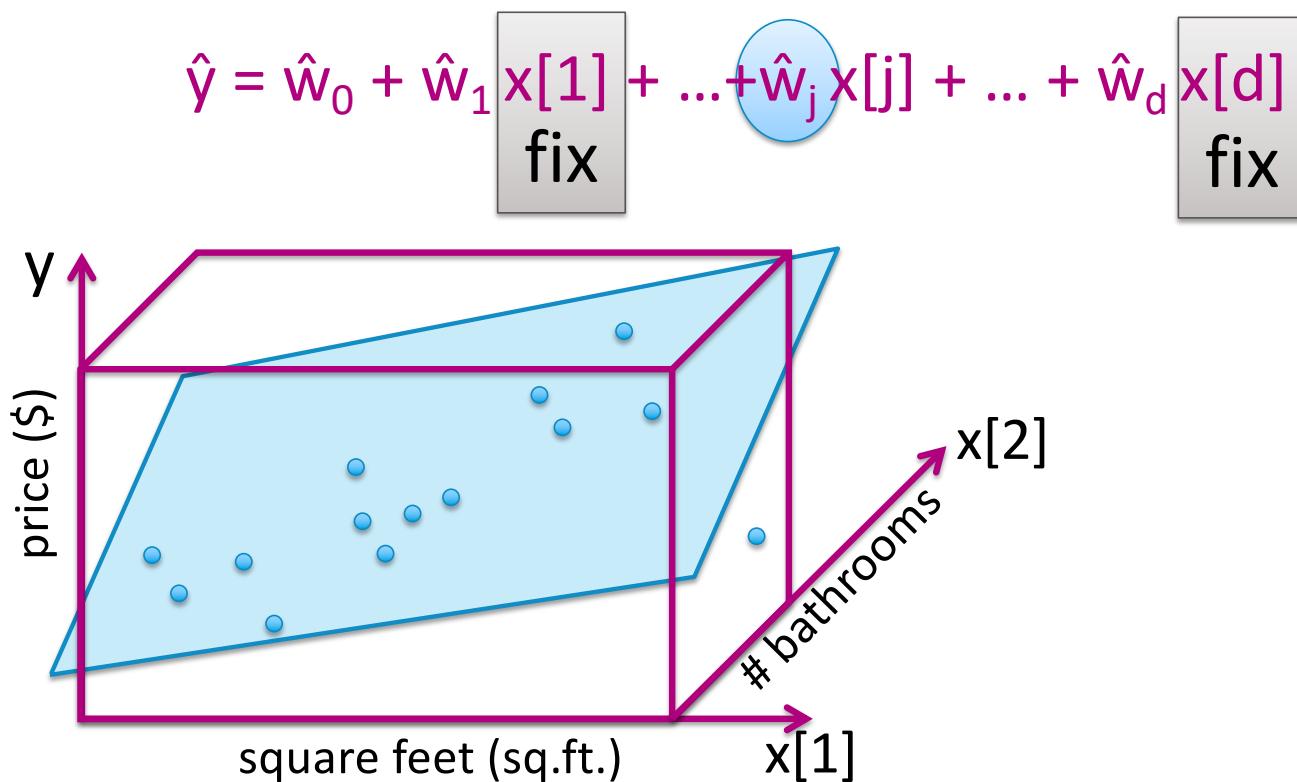
w2 = Fixing x[1] , i.e the square footage, how much does increasing x[2], number of bathrooms, by one unit increase the price.



Interpreting the coefficients – Two linear features



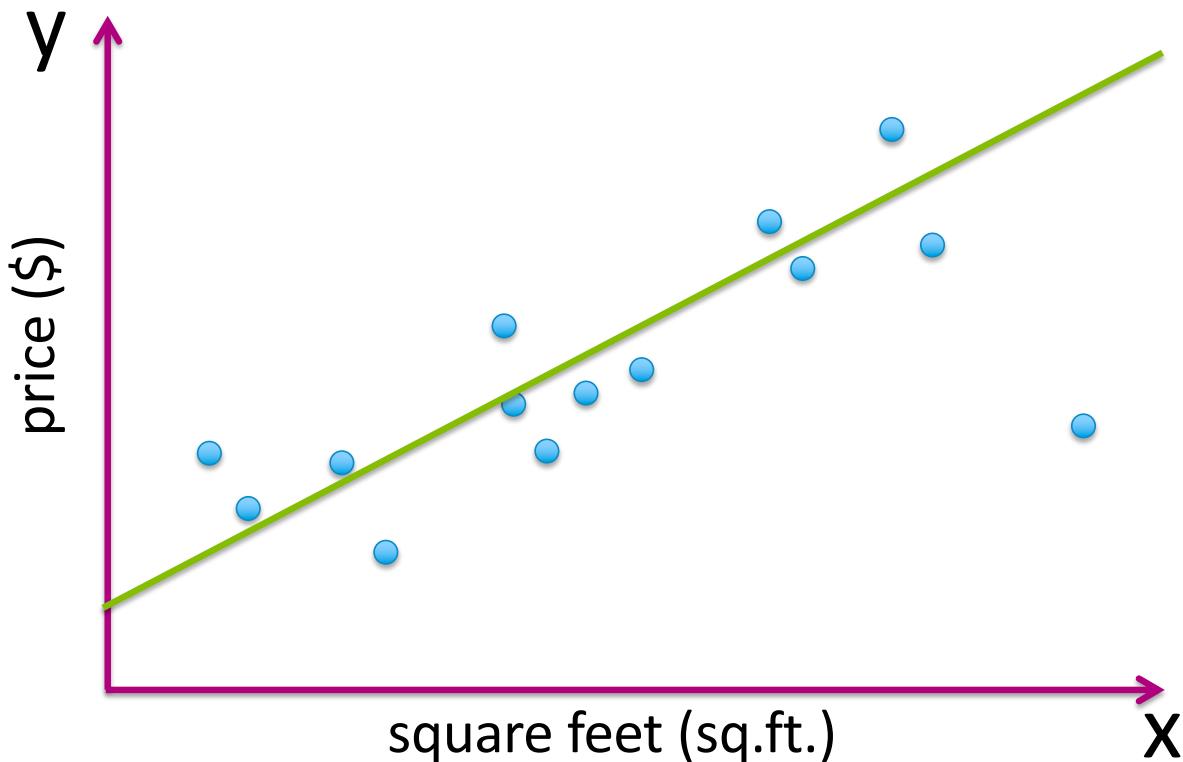
Interpreting the coefficients – Multiple linear features



Why min RSS?

Makes sense for gaussian noise...

Assuming Gaussian noise



Model for ϵ_i :
 $N(0, \sigma^2)$

Implied distribution
on y_i :
 $y_i \sim N(w^T x_i, \sigma^2)$

Maximum likelihood estimate of params

If $y_i = \mathbf{x}_i^T \mathbf{w} + \varepsilon_i$ is $N(\mathbf{x}^T \mathbf{w}, \sigma^2)$ Maximum Likelihood estimate

And we see data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, what is the MLE of \mathbf{w} ?

$P(y | x, w, \sigma^2) = \prod_{i=1}^n \phi(y_i; \mathbf{x}_i^T \mathbf{w}, \sigma^2)$
where ϕ is the probability density function of the normal distribution

find the maximum of this function; note is not strictly a probability so this equality is not quite true
but ϕ is still a relative measure of probability so it performs adequately.

taking the log we see that the likelihood

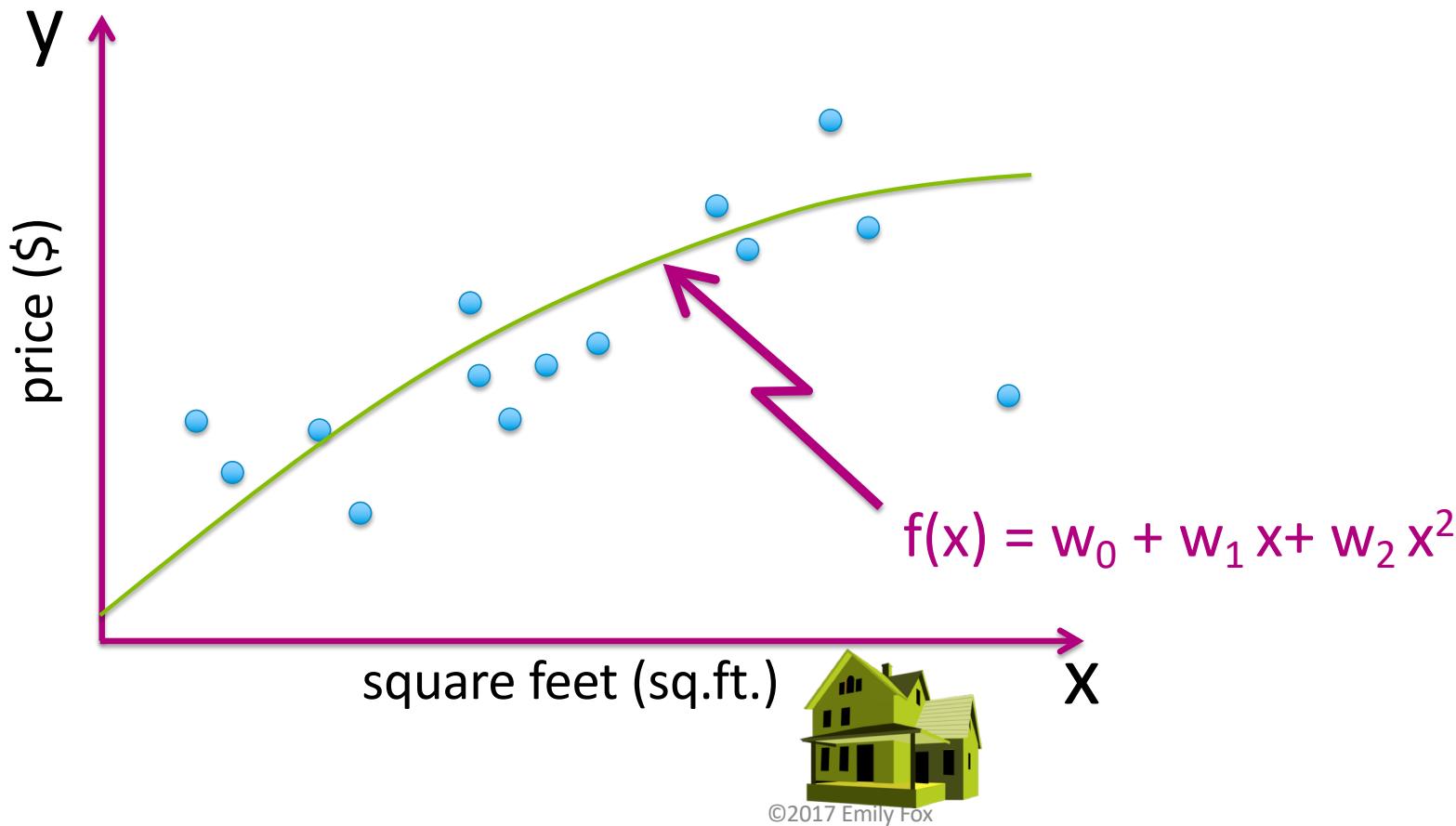
Conclusion

If $y_i = x_i^T \mathbf{w} + \varepsilon_i$ is $N(x^T \mathbf{w}, \sigma^2)$

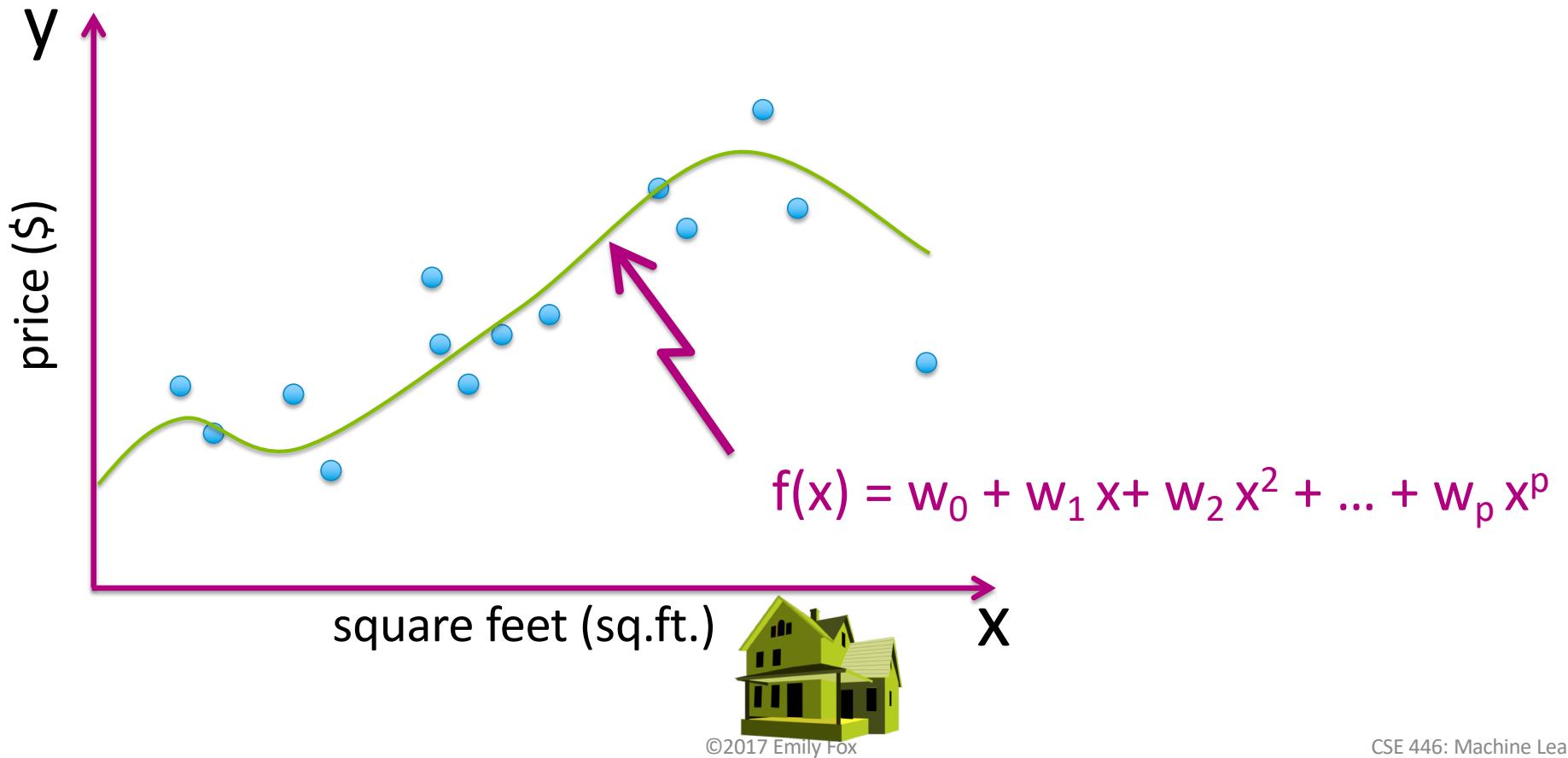
Then choosing \mathbf{w} to maximize log-likelihood
is same as choosing \mathbf{w} to minimize RSS!

Feature maps, polynomial regression and basis expansion

What about a quadratic function?



Even higher order polynomial



Polynomial Regression

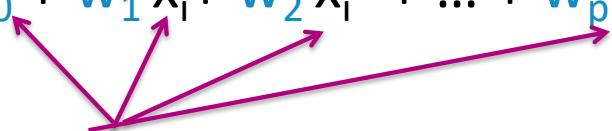
- Start with single input feature x (e.g. square footage of house) and training set: $\{(x_i, y_i)\}_{i=1..n}$
- Define feature map that transforms each x_i to higher dimensional feature vector $h(x_i)$.

Example: x_i

$$\begin{array}{ccccccc} h_0(x_i) & h_1(x_i) & h_2(x_i) & h_3(x_i) & h_4(x_i) & h_5(x_i) \\ 1 & x_i & x_i^2 & x_i^3 & x_i^4 & x_i^5 \end{array}$$

Polynomial regression

Model:

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


treat transformed inputs as different features

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 $p+1 = w_p$

Why might we want to use polynomial regression?

- Taylor Series!

More generally

- Start with input features $\mathbf{x} = (x[1], x[2], \dots, x[d])$ and training set: $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$
- Define feature map that transforms each \mathbf{x}_i to higher dimensional feature vector $h(\mathbf{x}_i)$.

Example: $x_i[1] \ x_i[2] \ x_i[3]$

$$\begin{array}{ccccccc} h_1(\mathbf{x}) & h_2(\mathbf{x}_i) & h_3(\mathbf{x}_i) & h_4(\mathbf{x}_i), & h_5(\mathbf{x}_i) & h_6(\mathbf{x}_i) & h_7(\mathbf{x}_i) \\ 1 & x_i[1] & x_i[1]^2 & x_i[1]x_i[2] & x_i[2] & x_i[2]^2 & \cos(\pi x_i[3]/6) \end{array}$$

General notation

Output: y  scalar

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

Notational conventions:

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

$x_i[j]$ = j^{th} input of i^{th} data point (*scalar*)

$\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_d(\mathbf{x}))$ feature map applied to input \mathbf{x} (*vector*)

$h_j(\mathbf{x})$ = j^{th} feature associated with input \mathbf{x} (*scalar*) (j^{th} basis function)

To fit these more general functions

- Start with input features $\mathbf{x} = (x[1], x[2], \dots, x[d])$ and training set: $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$
- Define feature map that transforms each \mathbf{x}_i to higher dimensional feature vector $h(\mathbf{x}_i)$.
- Model: $y_i = \sum_{j=1}^p w_j h_j(\mathbf{x}_i) + \varepsilon_i$
- Find $\hat{\mathbf{w}}$ that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j h_j(\mathbf{x}_i))^2$

Recap of concepts

What you can do now...

- Describe linear regression (and feature maps)
- Write a regression model using multiple inputs or features thereof.
- 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