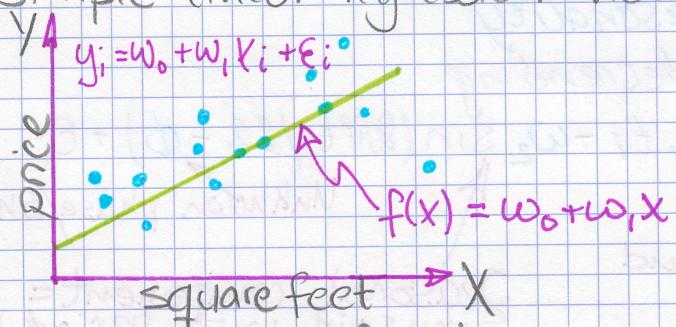


Multiple Regression

Linear regression with multiple features

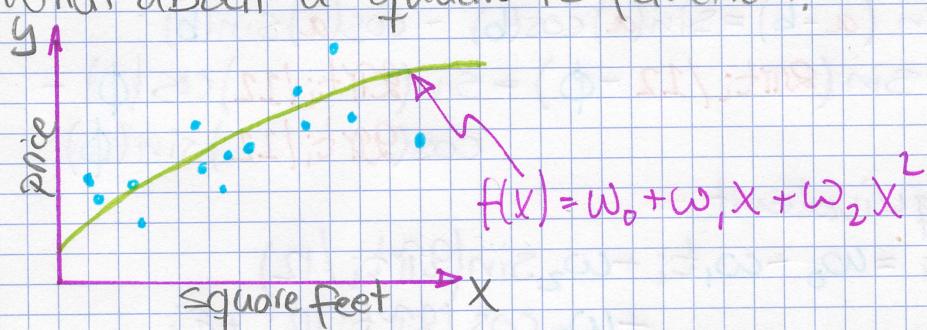
Polynomial regression

Simple linear regression model

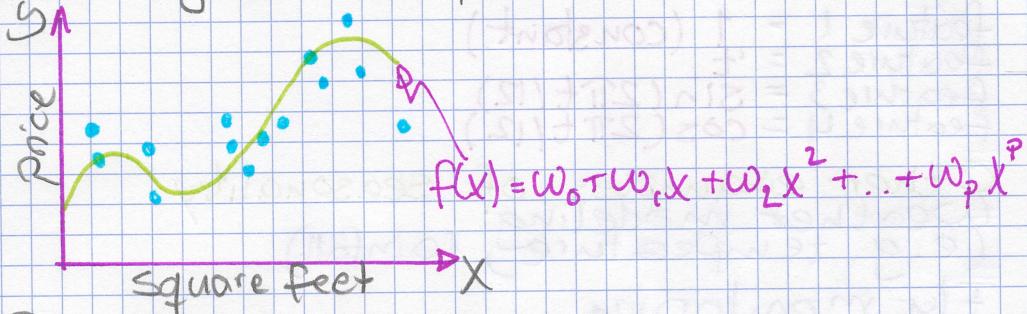


More complex functions
of a single input

What about a quadratic function?



Ever higher order polynomial



Polynomial regression
Model:

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

treat as different features

feature 1 = 1 (constant)
feature 2 = x
feature 3 = x^2
...

parameter 1 = w_0
parameter 2 = w_1
parameter 3 = w_2
...

feature $p+1 = x^p$

parameter $p+1 = w_p$

Modeling Seasonality

An example detrending

$$y_i = w_0 + w_1 t_i + w_2 \sin(2\pi t_i / 12 - \phi) + \epsilon_i$$

↑
Linear trend
Seasonal component =
Sinusoid with period 12
(resets annually)

Unknown phase/shift

Trigonometric identity:

$$\sin(a - b) = \sin(a)\cos(b) - \cos(a)\sin(b)$$
$$\rightarrow \sin(2\pi t_i / 12 - \phi) = \sin(2\pi t_i / 12)\cos(\phi) - \cos(2\pi t_i / 12)\sin(\phi)$$

Equivalent

$$y_i = w_0 + w_1 t_i + w_2 \sin(2\pi t_i / 12) + w_3 \cos(2\pi t_i / 12) + \epsilon_i$$

feature 1 = 1 (constant)
feature 2 = t
feature 3 = $\sin(2\pi t / 12)$
feature 4 = $\cos(2\pi t / 12)$

Other examples of seasonality
(Weather modeling (e.g. temperature, rainfall))

Flu monitoring

Demand forecasting
(e.g. jacket purchases)

Motion capture data

Regression with general features of y input

Generic basis expansion

Model:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i$$

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

\uparrow \uparrow
 j^{th} feature
 j^{th} regression coefficient or weight

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

Motivating the use of multiple inputs

Adding more inputs

$$f(x) = w_0 + w_1 \text{sq.ft} + w_2 \# \text{bath}$$



Many possible inputs

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

General notation

Output: y scalar

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

\uparrow
 d -dim vector

Notational conventions:

$x[j]$ = j^{th} input (Scalar)

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

x_i = input of i^{th} data point (vector)

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

Regression with features of multiple inputs

Simple hyperplane

Model:

$$y_i = w_0 + w_1 x_{i[1]} + \dots + w_d x_{i[d]} + \epsilon_i$$

feature 1 = 1

feature 2 = $x_{i[1]}$ sq.ft

feature 3 = $x_{i[2]}$ #bath

...
feature $d+1 = x_{i[d]}$ lot size

More generically D-dimensional curve

Model:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i$$

feature 1 = $h_0(x)$... e.g. 1

feature 2 = $h_1(x)$... e.g. $x_{i[1]} = 5 \text{ q.ft}$

feature 3 = $h_2(x)$... e.g. $x_{i[2]} = \# \text{bath}$
or $\log(x_{i[1]}) x_{i[2]} =$
 $= \log(\#\text{bed}) \times \#\text{bath}$

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

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

More on notation

observations (x_i, y_i) : N

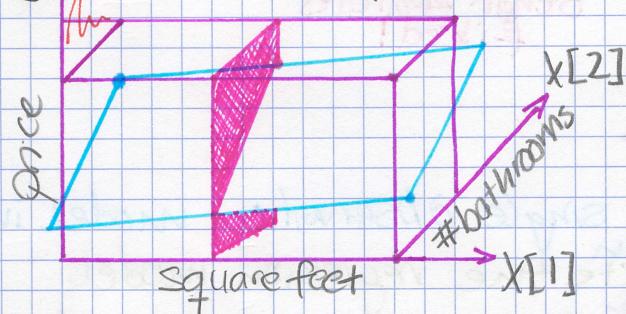
inputs $x_{i[j]}$: d

features $h_j(x)$: D

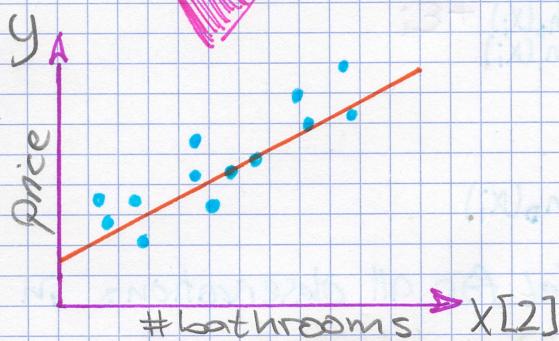
Interpreting the multiple regression fit

Interpreting the coefficients - Two linear features

$$\hat{y} = \hat{\omega}_0 + \hat{\omega}_1 \underbrace{x[1]}_{\text{fix}} + \hat{\omega}_2 x[2]$$



For fixed
Sq. ft



Interpreting the coefficients -
Multiple linear features

$$\hat{y} = \hat{\omega}_0 + \hat{\omega}_1 \underbrace{x[1]}_{\text{fix}} + \dots + \underbrace{\hat{\omega}_j x[j]}_{\text{fix}} + \dots + \hat{\omega}_d \underbrace{x[d]}_{\text{fix}}$$

Interpreting the coefficients - Polynomial regression

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x + \dots + \hat{w}_j x^j + \dots + \hat{w}_p x^p$$



Can't hold
other features
fixed!

Recognizing the single observation model in vector notation

Step 1: Rewrite the regression model

Matrix notation
For observation i

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

$$y_i = \underbrace{\begin{bmatrix} w_0 & w_1 & \dots & w_p \end{bmatrix}}_{\omega^T} \underbrace{\begin{bmatrix} h_0(x_i) \\ h_1(x_i) \\ \vdots \\ h_p(x_i) \end{bmatrix}}_{h(x_i)} + \epsilon_i$$

$$= \omega^T h(x_i) + \epsilon_i$$

Recognizing the model for all observations in matrix notation

Rewrite in matrix notation
For all observations together

$$\begin{array}{c|ccccc} Y & h_0(x_1) & h_1(x_1) & \dots & h_p(x_1) & \omega \\ \hline y_1 & & & & & w_0 \\ y_2 & & & & & w_1 \\ y_3 & & & & & w_2 \\ \vdots & & & & & \vdots \\ y_n & & & & & w_p \end{array} + \begin{array}{c|c} \epsilon & \epsilon_1 \\ \hline \vdots & \vdots \\ \epsilon_n & \epsilon_p \end{array} \Rightarrow \boxed{Y = H\omega + \epsilon}$$

Computing the cost of a D-dimensional curve

Step 2: Compute the cost

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

RSS for multiple regression

$$RSS(w) = \sum_{i=1}^N (y_i - h^T(x_i)w)^2$$

$$\begin{matrix} \square = & \begin{matrix} h^T(x_i) & w \end{matrix} \\ \hat{y}_i & h_0(x_i) & h_1(x_i) & \begin{matrix} w_0 \\ w_1 \end{matrix} \end{matrix}$$

$$RSS \text{ in matrix notation } RSS(w) = \sum_{i=1}^N (y_i - h(x_i)^T w)^2$$

$$= (y - Hw)^T (y - Hw)$$

$$\hat{y} = Hw$$

$$y - Hw = (y - \hat{y}) =$$

$$\begin{bmatrix} \text{residual}_1 \\ \text{residual}_2 \\ \vdots \\ \text{residual}_N \end{bmatrix}$$

$$[\text{residual}_1, \text{residual}_2, \dots, \text{residual}_N]$$

$$\begin{bmatrix} \text{residual}_1 \\ \text{residual}_2 \\ \vdots \\ \text{residual}_N \end{bmatrix}$$

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

Computing the gradient of RSS

Step 3: Take the gradient

$$\nabla \text{RSS}(\omega) = \nabla [(\mathbf{y} - \mathbf{H}\omega)^T (\mathbf{y} - \mathbf{H}\omega)] \quad \left\{ \begin{array}{l} \text{Gradient of} \\ \text{RSS} \end{array} \right.$$
$$= -2\mathbf{H}^T (\mathbf{y} - \mathbf{H}\omega)$$

Why? By analogy to 1D case:

$$\frac{\partial}{\partial \omega} (\mathbf{y} - \mathbf{H}\omega)(\mathbf{y} - \mathbf{H}\omega) = \frac{\partial}{\partial \omega} (\mathbf{y} - \mathbf{H}\omega)^2 = 2 \cdot (\mathbf{y} - \mathbf{H}\omega)^T (-\mathbf{H})$$

↑ ↑
scalar

$$= -2\mathbf{H}(\mathbf{y} - \mathbf{H}\omega)$$

Approach 1: closed-form solution

Step 4: Approach 1: Set the gradient = 0

Closed-form solution

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

Solve for ω :

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

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

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

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

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

$$D \quad \underbrace{\begin{bmatrix} & & \\ & & \end{bmatrix}}_N \quad \underbrace{\begin{bmatrix} & & \\ & & \end{bmatrix}}_{\substack{\# \text{features} = D \\ \# \text{obs} = N}} = \boxed{D \times D} \quad \# \text{features}$$

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

really
of linearly
Ind: observations

Complexity of inverse:
 $O(D^3)$

Approach 2: Gradient descent
Step 4 approach 2

$$\omega^{(t+1)} \leftarrow \omega^{(t)} - \eta \nabla_{\omega} \text{RSS}(\omega^{(t)})$$
$$-2H^T(y - H\omega^{(t)})$$

$$\omega^{(t)} + 2\eta H^T(y - H\omega^{(t)})$$
$$\underbrace{\quad}_{\hat{y}(\omega^{(t)})}$$

Feature-by-Feature update

$$\text{RSS}(\omega) = \sum_{i=1}^N (y_i - h(x_i)^T \omega)^2$$
$$= \sum_{i=1}^N (y_i - \omega_0 h_0(x_i) - \omega_1 h_1(x_i) - \dots - \omega_D h_D(x_i))^2$$

Partial with respect to ω_j :

$$\sum_{i=1}^N 2(y_i - \omega_0 h_0(x_i) - \omega_1 h_1(x_i) - \dots - \omega_D h_D(x_i)) \cdot (-h_j(x_i))$$
$$= -2 \sum_{i=1}^N h_j(x_i) (y_i - h(x_i)^T \omega)$$

update to jth feature weight:

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

bathrooms
jth feature

bath for houseⁱ

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

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 $w^{(1)} = 0$ (or randomly, or smartly), $t=1$

while $\|\nabla \text{RSS}(w^{(t)})\| > \epsilon$ tolerance

for $j=0, \dots, D$

$$\text{partial}[j] = -2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \text{partial}[j]$$

$t \leftarrow t+1$

$$\sqrt{\text{partial}[0]^2 + \dots + \text{partial}[D]^2}$$