

Linear Regression: Model and Algorithms (part 2)

CSE 446

Most slides by Emily Fox
Presented by Anna Karlin
April 8, 2019

- panepto
- latex
- hw1



XKCD

Linear regression: a supervised learning problem

Goal: to predict some output from some inputs using labelled examples. Example: house sales price from square footage

Supervised learning: Problem of learning a function that maps inputs to outputs based on labelled examples.

Regression: When the labels are real numbers

Linear regression: a supervised learning problem

Goal: to predict some output from some inputs/features. Example: house sales price from square footage

Step 1: Define set up and get data

- a model for how the output y depends on the inputs x .
- We assumed that y is a linear function of features + noise.

$$y = \mathbf{w}^T \mathbf{x} + \varepsilon$$

$$\mathbf{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

- A training set (labelled examples): $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$

Linear regression: a supervised learning problem

Goal: to predict some output from some inputs/features.

Step 1: Define set up ($y = \mathbf{w}^T \mathbf{x} + \epsilon$) and get data $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$

Step 2: find the parameters \mathbf{w} that minimize the “loss/cost” on the training set.

- Our **loss function** was residual sum of squares (RSS)
- Find $\hat{\mathbf{w}}$ that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^d w_j x_i[j])^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$
- Found solution by solving for gradient of $(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$

Solution: $\boxed{\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}}$ If $\mathbf{X}^T \mathbf{X}$ is invertible, could write $\hat{\mathbf{w}} = \underline{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$

5

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} & & \\ \mathbf{x}_1^T & \cdots & \mathbf{x}_n^T \end{pmatrix}$$

prediction

Linear regression: a supervised learning problem

Goal: to predict some output from some inputs/features.

Step 1: Define set up ($y = \mathbf{w}^T \mathbf{x} + \epsilon$) and get data $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$

Step 2: find the parameters $\hat{\mathbf{w}}$ that minimizes $\text{RSS} = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$

Step 3: Use $\hat{\mathbf{w}}$ to make predictions.

Given \mathbf{x} , predict output: $\hat{\mathbf{w}}^T \mathbf{x}$

Plan for today:

- Gradient descent
- Handling an intercept
- More features/more complex models
- How well does it work?

$$\vec{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \end{pmatrix}$$

Linear regression: a supervised learning problem

Goal: to predict some output from some inputs/features. Example: house sales price from square footage

Step 1: Define set up and get data

- a model for how the output y depends on the inputs x .
- We assumed that y is a linear function of features + noise.

$$y = \mathbf{w}^T \mathbf{x} + \varepsilon$$

- A training set (labelled examples): $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$

Step 2: find the parameters w that minimize the “loss/cost” on the training set.

- Our **loss function** was residual sum of squares (RSS)
- Find $\hat{\mathbf{w}}$ that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^d \mathbf{x}_i^T \mathbf{w}_j)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$
- Found solution by solving for gradient of $(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$

$$\text{Solution: } \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y} \quad \text{If } \mathbf{X}^T \mathbf{X} \text{ is invertible, could write } \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Step 3: Use $\hat{\mathbf{w}}$ to make predictions, Given \mathbf{x} , predict output: $\mathbf{w}^T \mathbf{x}$

Fitting the linear regression model

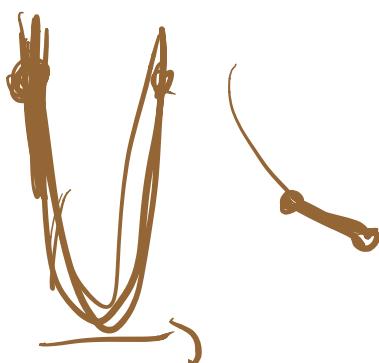
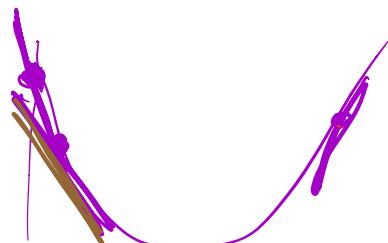
Gradient descent

$f(w)$

function convex, cont, nice

Gradient Descent – univariate case

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



$$f(w+sw) \approx f(w) + sf'(w)$$

if s very small

$$\begin{cases} f'(w) < 0 & s \text{ to be positive} \\ > 0 & s \text{ negative} \end{cases}$$

initially w_0 as long as $|f'(w_t)| > \epsilon$

$$w_{t+1} = w_t -$$

$$ts = t+1$$

$$m|f'(w)|$$

step size
step size

$$\frac{\epsilon}{t}$$

(1)

Gradient Descent – multivariate case



$$f(w_1, \dots, w_d)$$

$$f(\vec{w} + \vec{s}) \approx f(\vec{w}) + \sum_{j=1}^d s_j \frac{\partial f}{\partial w_j}(\vec{w}) = f(\vec{w}) + \vec{s}^T \nabla f(\vec{w})$$

to make the fn \downarrow as quickly as possible

want to move in direction of negative gradient

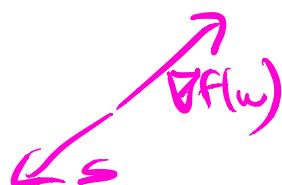
initialise \vec{w}_0

while $\|\nabla f(\vec{w}_+)\| \geq \epsilon$

$$\vec{w}_{+1} = \vec{w}_+ - \eta \nabla f(\vec{w}_+)$$

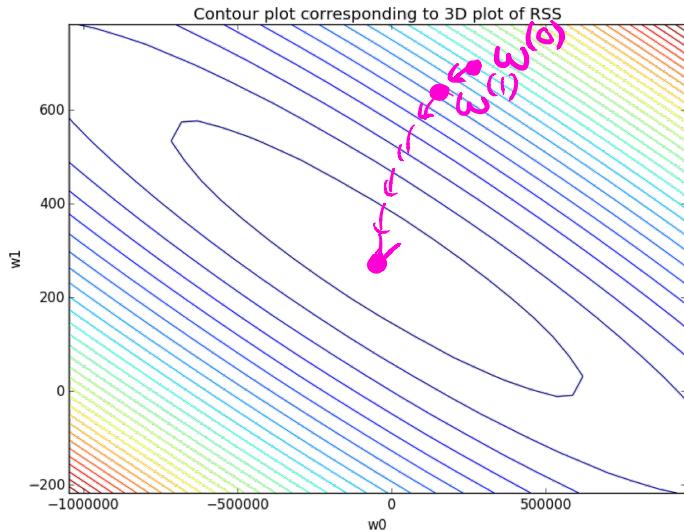
$$\|\vec{s}\| \|\nabla f(\vec{w})\| \cos \theta$$

angle between $\nabla f(\vec{w})$ and \vec{s}



$$\|\vec{x}\| = \sqrt{\sum x_i^2}$$

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



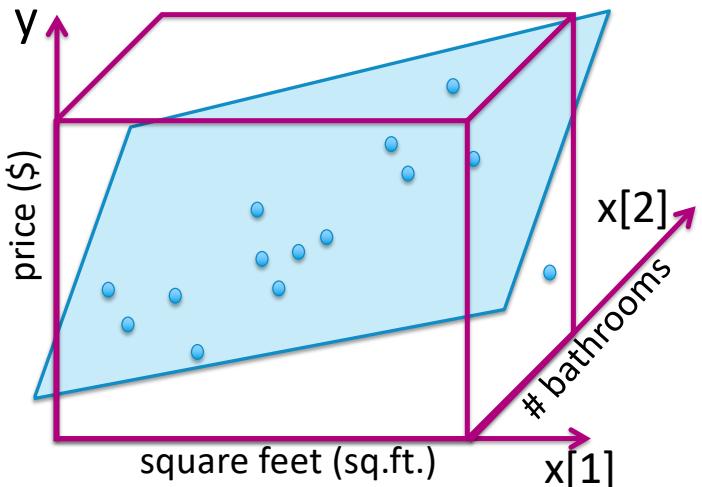
while not converged

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



$$-2X^T(y - Xw^{(t)})$$

Interpreting elementwise



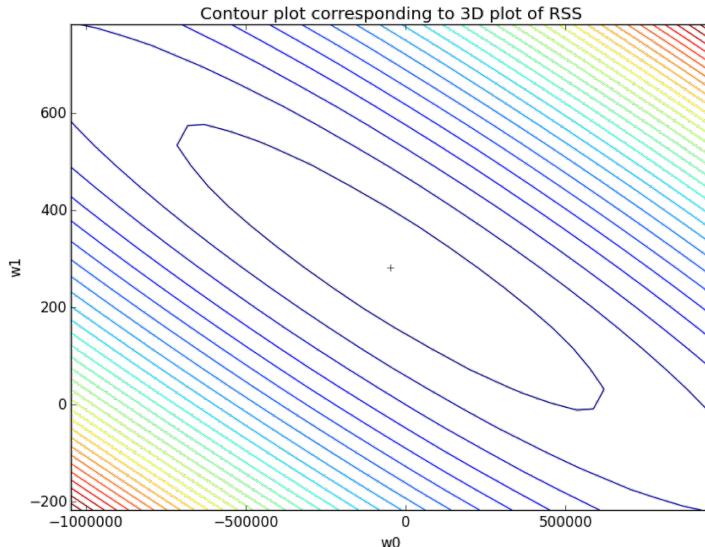
Update to j^{th} feature weight:

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

predict for x_i
Using current w
 $x_i^T w^{(t)}$

Summary of gradient descent for multiple regression

$$\frac{\partial \text{RSS}(\omega^{(t)})}{\partial w_j}$$



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

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

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

$$\sum_{j=1}^c w_j x_{f(j)} + b$$

- 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:

$$\hat{\omega}^T \mathbf{x}$$

What if we want to allow for an intercept?

Assume that $y = \mathbf{w}^T \mathbf{x} + b + \varepsilon$

Find $\hat{\mathbf{w}}, b$ that minimize

$$RSS = \sum_{i=1}^n (y_i - \sum_{j=1}^d w_j x_i[j] - b)^2$$

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

Handling an intercept (constant term)

Assume that $y = \mathbf{w}^T \mathbf{x} + b + \varepsilon$

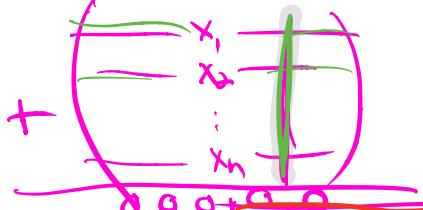
Find $\hat{\mathbf{w}}, \hat{b}$ that minimize $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^d w_j x_i[j] - b)^2$

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

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.



$$\begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \\ 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

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

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

Find \hat{w}, \hat{b} that minimize

$$\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^d w_j x_i[j] - b)^2$$

$$= (y - Xw - b\mathbf{1})^T (y - Xw - b\mathbf{1})$$

$$\nabla_w = 0$$

coeff of b is
 $2 \sum_{i=1}^n x_i[1]$

$$\sum_i x_i[1] y_i - \sum_i x_i[1] x_i^T w - \sum_i x_i[1] b = 0$$

$$\frac{\partial \text{RSS}}{\partial w_j} = \text{partial}[w_j] = -2 \sum_{i=1}^n x_i[j] (y_i - x_i^T w - b) = 0$$

$$\text{partial}[b] = -2 \sum_{i=1}^n (y_i - x_i^T w - b) = 0$$

$$\sum_i y_i - \cancel{\sum_i x_i^T w} - nb = 0$$

$$(\sum_i x_i^T) w = \cancel{nb}$$

$$b = \frac{1}{n} \sum_{i=1}^n y_i$$

CSE 446: Machine Learning

$$\nabla_w \text{RSS} = -2X^T(y - Xw)$$

$$b = \frac{1}{n} \sum_{i=1}^n y_i$$

$$X^T X \hat{w} = X^T y$$

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

$$(\vec{x}_i, y_i)$$



$$(\vec{x}_i - \bar{\mu}, y_i)_{i=1}^n$$

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i$$

$$\tilde{X} = \begin{pmatrix} \vec{x}_1 - \bar{\mu} \\ \vec{x}_2 - \bar{\mu} \\ \vdots \\ \vec{x}_n - \bar{\mu} \end{pmatrix}$$

$$\begin{aligned} \sum_{i=1}^n \vec{x}_i &= \sum_{i=1}^n (\vec{x}_i - \bar{\mu}) + n\bar{\mu} \\ &= \sum_{i=1}^n (\vec{x}_i - \bar{\mu}) = \sum \vec{x}_i - n\bar{\mu} \\ &= \vec{0} \end{aligned}$$

$$\vec{x}$$

predict

$$(\vec{x}^T - \vec{\mu}^T) \hat{w} + \hat{b}$$

18

Soln:

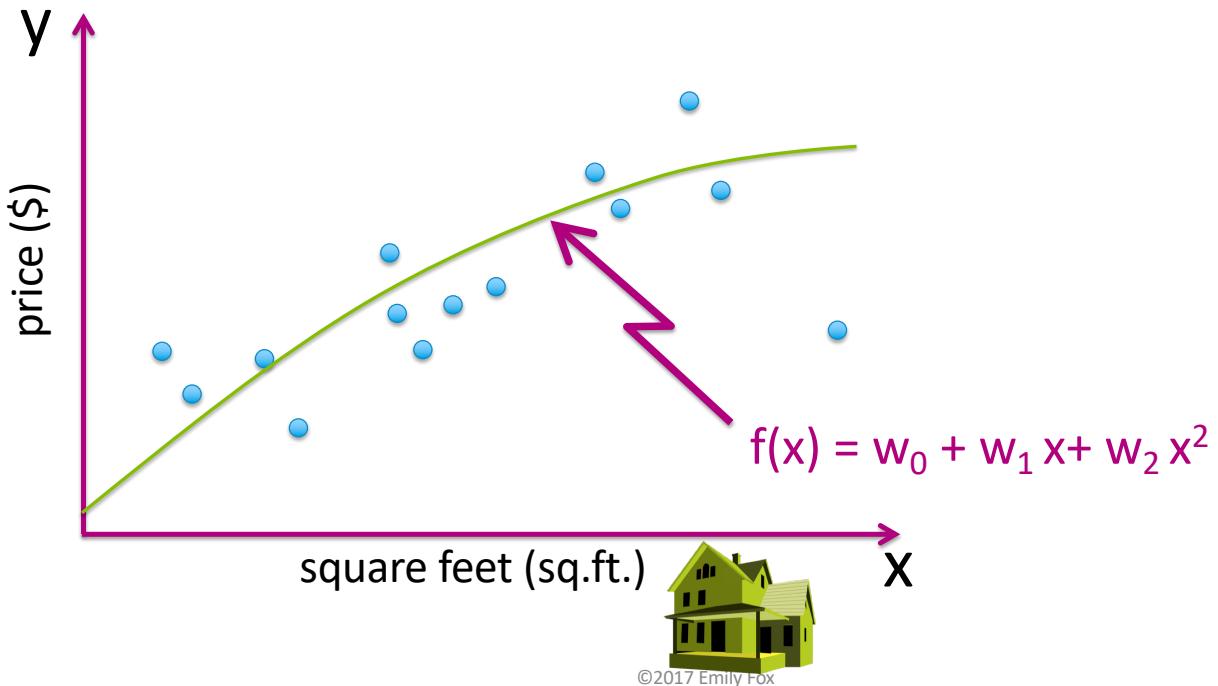
$$\hat{w}, \hat{b}$$

=

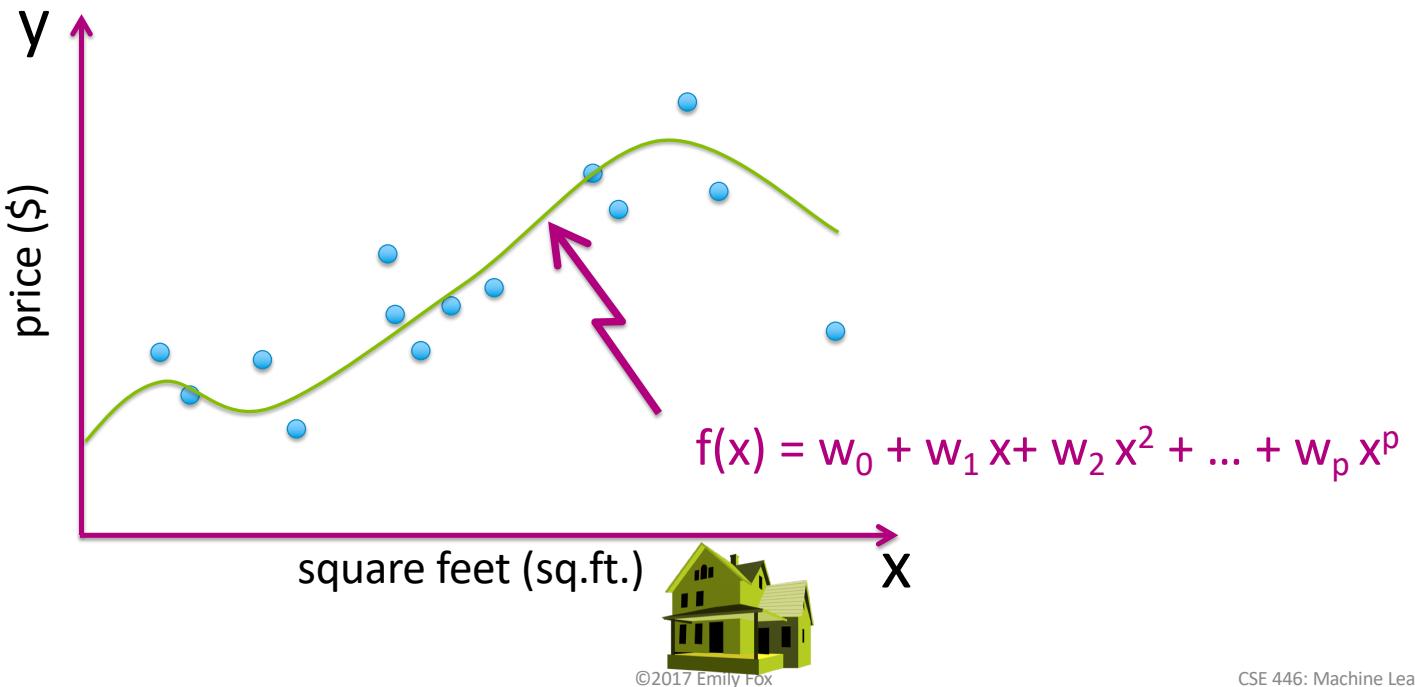
$$\begin{cases} \tilde{X}^T \tilde{X} \hat{w} = \tilde{X}^T y \\ \hat{b} = \frac{1}{n} \sum_{i=1}^n y_i \end{cases}$$

More features, more complex models

What about a quadratic function?



Even higher order polynomial



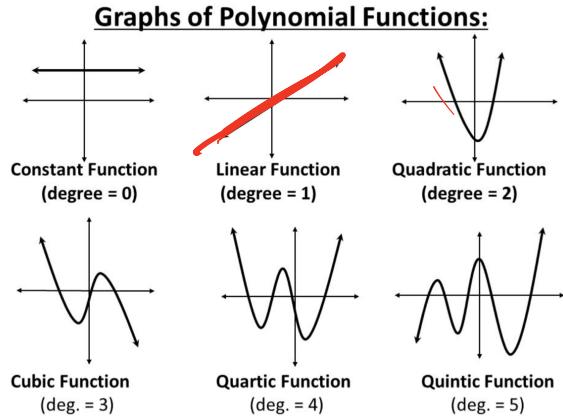
Polynomial regression (single input)

Goal: to predict some output from some inputs/features.

Step 1: Assume that y (sales price) is a polynomial function of feature (square footage)+ noise.

$$y_i = \sum_{j=0}^p w_j x_i^j + \epsilon \quad \text{A training set (labelled examples):}$$

(x_i^j)



Polynomial regression

Goal: to predict some output from some inputs/features.

Step 1: $y_i = \sum_{j=0}^p w_j x_i^j + \epsilon$ A training set (labelled examples):

Step 2: find params w that minimize the “loss/cost” on training set $\{(x_i, y_i)\}_{i=1..n}$

- **Loss function** is residual sum of squares (RSS)
- Find \hat{w} that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=0}^{p-1} w_j x_i^j)^2$

Polynomial regression

Goal: to predict some output from some inputs/features.

Step 1: $y_i = \sum_{j=0}^p w_j x_i^j + \epsilon$ A training set (labelled examples):

Step 2: find params w that minimize the “loss/cost” on training set $\{(x_i, y_i)\}_{i=1..n}$

- Find \hat{w} that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=0}^{p-1} w_j x_i^j)^2$ $\sum w_j x_i^j$
- Just as easy to solve! Just think of x_i^j as one of p features associated with the i^{th} observation.
- Instead of single input x_i , define features $h(x) = (1, x, x^2, \dots, x^p)$

$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$
$$h(x) = \begin{pmatrix} 1 & x^2 & x^3 & \dots \end{pmatrix}$$
$$H = \begin{pmatrix} h(x_1) \\ h(x_2) \end{pmatrix}$$

$$h(x_i) = (h_0(x_i), h_1(x_i), h_2(x_i), h_3(x_i), h_4(x_i), h_5(x_i))$$
$$= (1, x_i, x_i^2, x_i^3, x_i^4, x_i^5)$$

CSE 446: Machine Learning

find w to min

$$RSS(\omega) = \underset{x}{\uparrow} (y - H\omega)^T (y - H\omega)$$

Polynomial regression

Step 1: $y_i = \sum_{j=0}^p w_j x_i^j + \epsilon$

Step 2: find the parameters w that minimize the “loss/cost” on the training set.

- Find \hat{w} that minimizes $RSS = \sum_{i=1}^n (y_i - \sum_{j=0}^p w_j x_i^j)^2 = (y - Hw)^T (y - Hw)$
- Find solution by solving for gradient of $(y - Hw)^T (y - Hw) = 0$

Solution: $H^T H \hat{w} = H^T y$

Polynomial regression

Step 1: $y_i = \sum_{j=0}^p w_j x_i^j + \epsilon$

Step 2: find the parameters w that minimize the “loss/cost” on the training set.

- Find \hat{w} that minimizes $RSS = \sum_{i=1}^n (y_i - \sum_{j=0}^p w_j x_i^j)^2 = (y - Hw)^T(y - Hw)$
- Find solution by solving for gradient of $(y - Hw)^T(y - Hw) = 0$

Solution: $H^T H \hat{w} = H^T y$

Step 3: Use $\hat{w} = (\hat{w}_0, \hat{w}_1, \dots, \hat{w}_p)$ to make predictions. Given x , let $h(x) = (1, x, x^2, \dots, x^p)$ and predict output:

$$h(x)^T \hat{w} = \hat{w}_0 + \hat{w}_1 x + \dots + \hat{w}_p x^p$$

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)

feature 2 = x

feature 3 = x²

...

feature p+1 = x^p

parameter 1 = w₀

parameter 2 = w₁

parameter 3 = w₂

...

parameter p+1 = w_p

Why might we want to use polynomial regression?

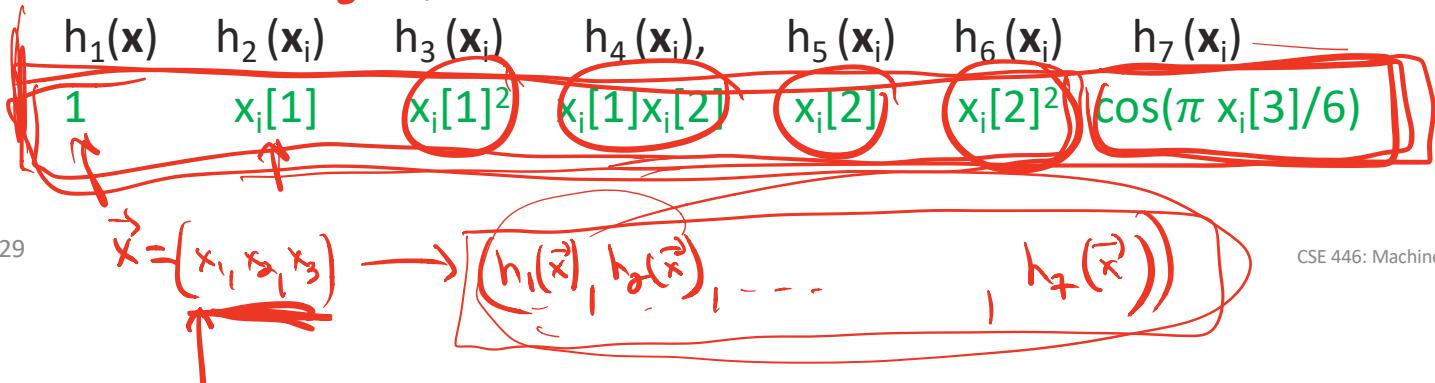
- Taylor Series!

$$f(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \dots$$

More generally

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

Example: $x_i[1] \quad x_i[2] \quad x_i[3]$
sq. feet \rightarrow ~~days~~ \rightarrow ~~the day when the house sold.~~



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*)
- $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_p(\mathbf{x}))$ 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)

- H = n by p matrix whose i^{th} row is $\mathbf{h}(\mathbf{x}_i)$

To fit these more general functions

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

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