

CS 66: Machine Learning

Prof. Sara Mathieson

Spring 2019



- Lab 2 due **THURSDAY** at midnight
- Lab 3 released today, due next Thursday
- Scribe notes for **extra credit!**
- Let me know of any partner issues (anytime)
- Reading posted (logistic regression & naïve bayes)

Research Talk

Thursday, February 14, 2019

11:30-12:30 p.m.

Science Center 256

Anne Cocos

Department of Computer and Information Science

University of Pennsylvania

“Semantic Structure from Paraphrase Pairs”

Outline for February 13

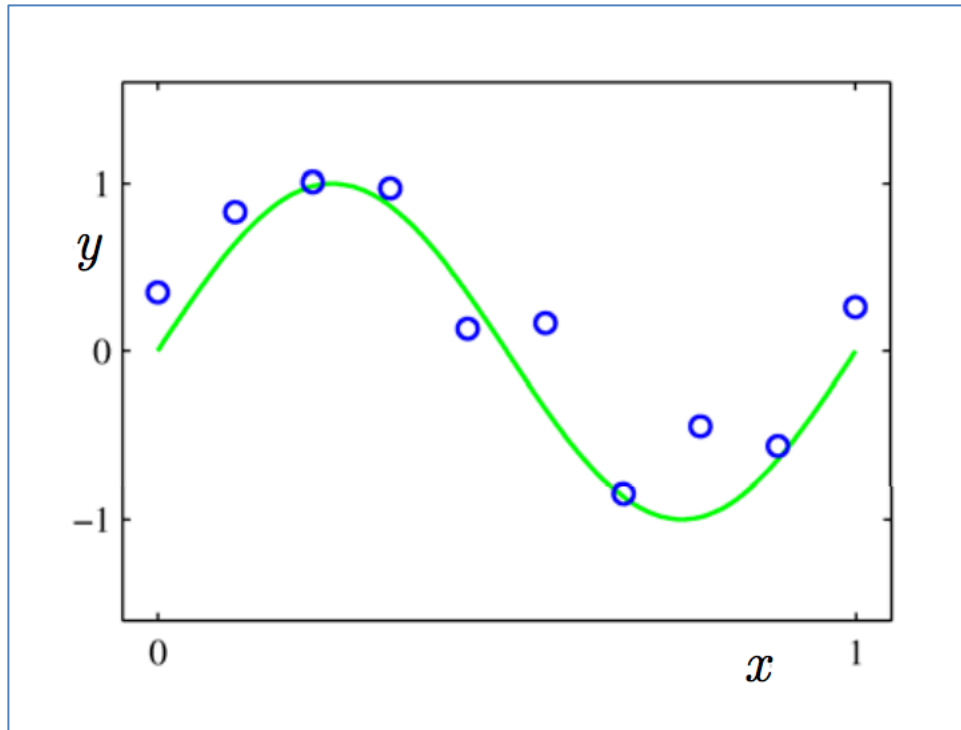
- Polynomial regression
- Regularization
- Linear regression for classification
- Begin: logistic regression

Outline for February 13

- Polynomial regression
- Regularization
- Linear regression for classification
- Begin: logistic regression

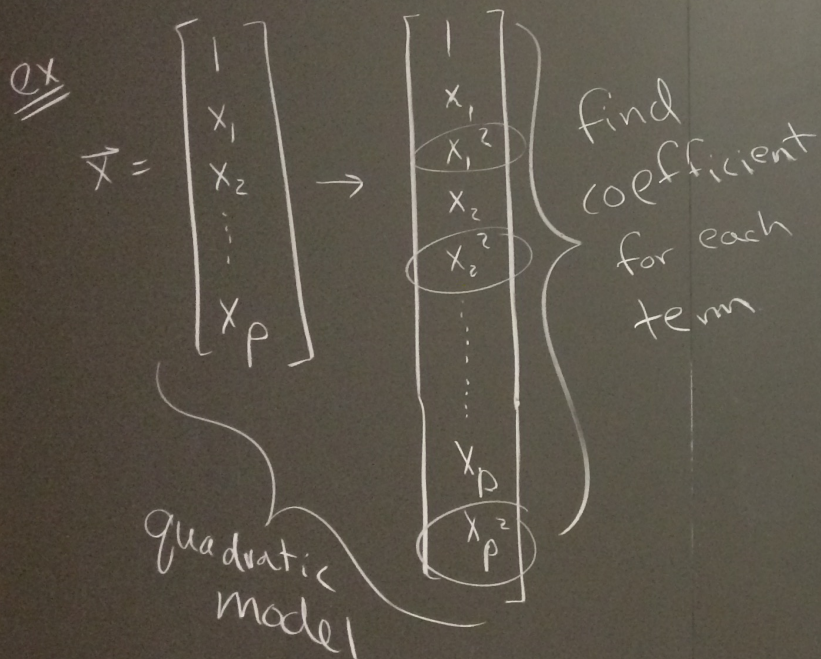
Polynomial Regression

- Can be thought of as regular linear regression with a change of basis



Polynomial Regression

before: $h_{\vec{b}}(\vec{x}) = \vec{b}^T \vec{x}$



Generalized linear model

$$h_{\vec{b}}(\vec{x}) = \vec{b}^T \underbrace{\phi(\vec{x})}_{\text{some transformation}}$$

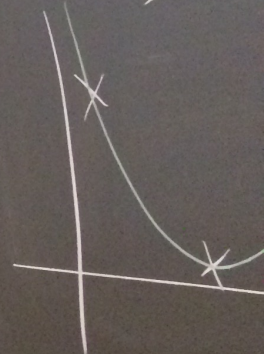
$$\vec{b} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\Phi = \begin{bmatrix} - & \phi(\vec{x}_1) & - \\ - & \phi(\vec{x}_2) & - \\ & \vdots & \\ - & \phi(\vec{x}_n) & - \end{bmatrix}$$

replace \vec{x} with Φ

n data points
=> fit

$n=3$, $d=2$



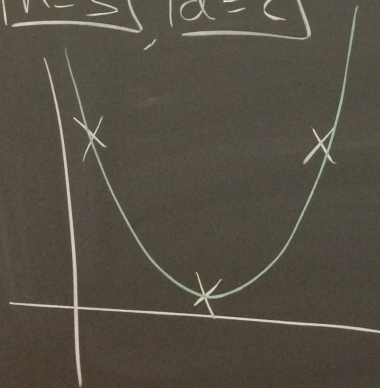
Analytic Solution

Transformation

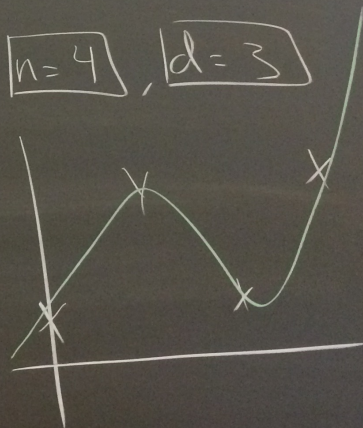
$$\vec{b} = (\Phi^T \Phi)^{-1} (\Phi^T \vec{y})$$

n data points, degree = $d = n-1$
 \Rightarrow fit perfectly!

$n=3$, $d=2$



$n=4$, $d=3$



Reg

J

$\lambda >$

- Sm

- la

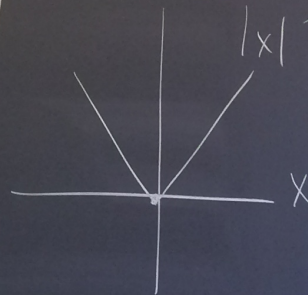
$$\boxed{p=1}, d = \text{degree}$$

$$\Rightarrow \boxed{h_{\vec{b}}(x) = b_0 + b_1 x + b_2 x^2 + \dots + b_d x^d}$$

in general. p features, degree $d \Rightarrow \boxed{pd+1}$ coeffs

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x_p \end{bmatrix}$$

$$\boxed{\vec{b}^T \phi(x)}$$



SGD

$b_0 \leftarrow$

$b_j \leftarrow$

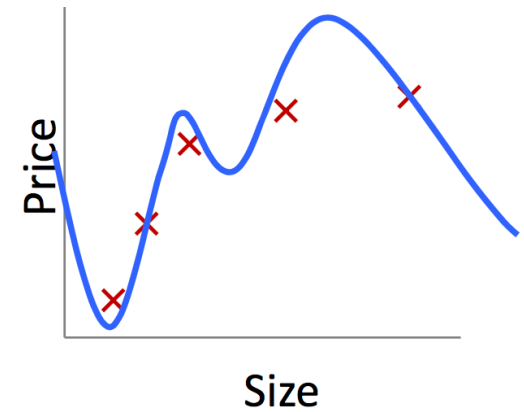
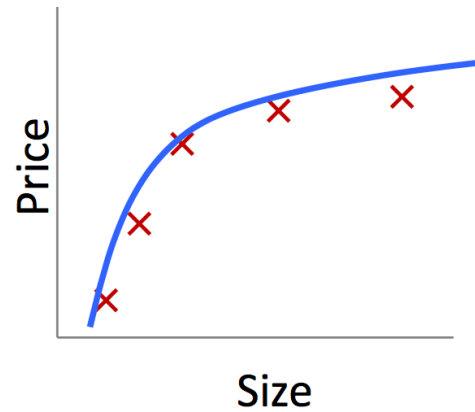
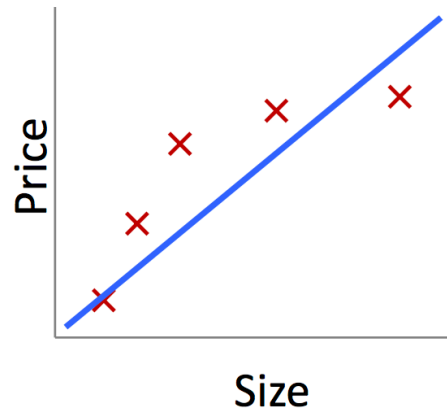
hyperparameter

Outline for February 13

- Polynomial regression
- **Regularization**
- Linear regression for classification
- Begin: logistic regression

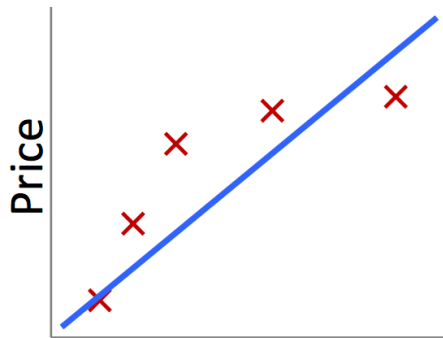
Generalization error

- Example: price vs. size (i.e. of a house or car)



Generalization error

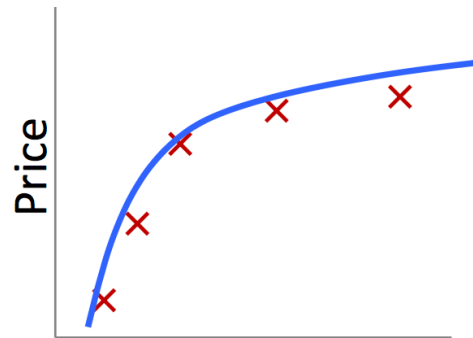
- Example: price vs. size (i.e. of a house or car)



Size

$$b_0 + b_1x$$

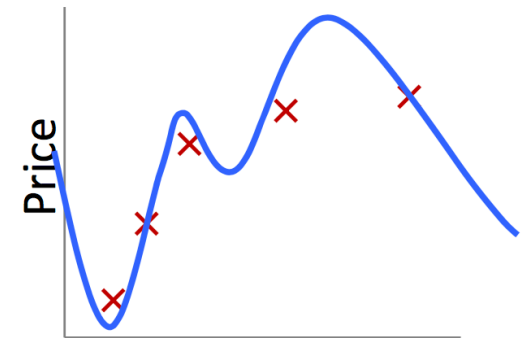
underfitting
(high bias)



Size

$$b_0 + b_1x + b_2x^2$$

correct fit



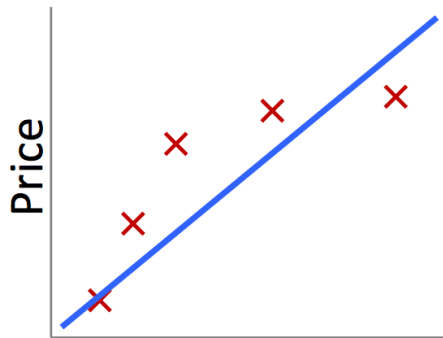
Size

$$b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4$$

overfitting
(high variance)

Generalization error

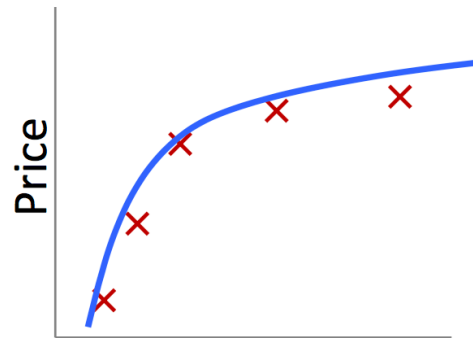
- Example: price vs. size (i.e. of a house or car)



Size

$$b_0 + b_1x$$

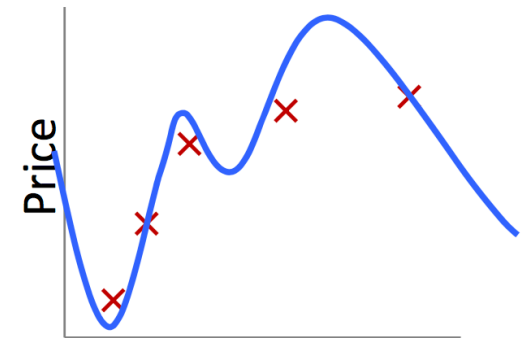
underfitting
(high bias)



Size

$$b_0 + b_1x + b_2x^2$$

correct fit



Size

$$b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4$$

overfitting
(high variance)

Structural error:

Hypothesis space cannot model true relationship

- ⇒ More data doesn't help
- ⇒ Need a more flexible model

Estimation (approximation) error:

Hypothesis space *can* model true relationship, BUT hard to identify correct model due to large hypothesis space, small n , or noise

- ⇒ Reduce hypothesis space
- ⇒ Add more data



Regularization

What if ...

- we have a limited # of training examples ($n < p$), or
- we want to automatically control the complexity of the learned hypothesis?

Regularization

What if ...

- we have a limited # of training examples ($n < p$), or
- we want to automatically control the complexity of the learned hypothesis?

Idea: penalize large values of b_j

Why prefer small weights?

- if large weights, small change in feature can result in large change in prediction
- prevent giving too much weight to any one feature
- might prefer zero weight for useless features

Common Regularizers

$$\|\mathbf{b}\|_0 = \sum_{j:b_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

Common Regularizers

$$||\mathbf{b}||_0 = \sum_{j:b_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

$$||\mathbf{b}||_1 = \sum_{j=1}^p |b_j|$$

L_1 norm

- Sum of magnitude of weights
- Not differentiable

Common Regularizers

$$||\mathbf{b}||_0 = \sum_{j:b_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

$$||\mathbf{b}||_1 = \sum_{j=1}^p |b_j|$$

L_1 norm

- Sum of magnitude of weights
- Not differentiable

$$||\mathbf{b}||_2 = \sqrt{\sum_{j=1}^p b_j^2}$$

L_2 norm

- Sum of squared weights
- Differentiable

Common Regularizers

$$||\mathbf{b}||_0 = \sum_{j:b_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

$$||\mathbf{b}||_1 = \sum_{j=1}^p |b_j|$$

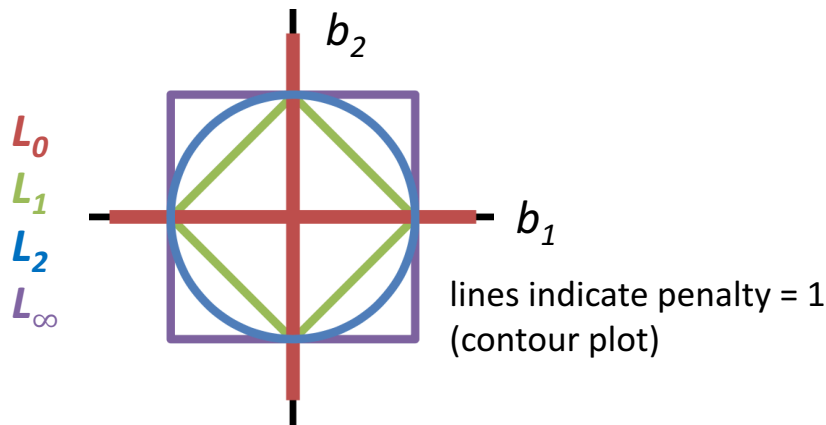
L_1 norm

- Sum of magnitude of weights
- Not differentiable

$$||\mathbf{b}||_2 = \sqrt{\sum_{j=1}^p b_j^2}$$

L_2 norm

- Sum of squared weights
- Differentiable

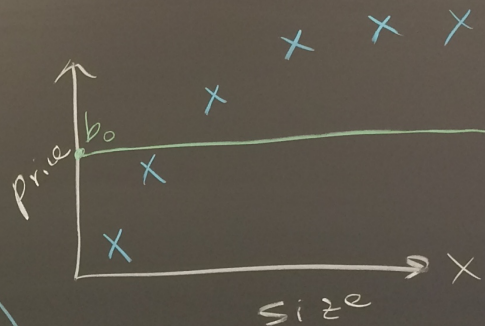


Regularization

$$J(\vec{b}) = \underbrace{\frac{1}{2} \sum_{i=1}^n (b^T \vec{x}_i - y_i)^2}_{\text{fit to training data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^p (b_j)^2}_{\text{regularization}}$$

$\lambda \geq 0$ is the regularization parameter.

- Small λ : more fit to training data.
- large λ : keeps weights small more generalization.



$$b_0 + \overset{0}{\cancel{b_1}}x + \overset{0}{\cancel{b_2}}x^2 + \overset{0}{\cancel{b_3}}x^3 + \overset{0}{\cancel{b_4}}x^4$$

Choose $\lambda = 10^{10}$ (very large)

do not
regularize
 (b_0)

SGD *

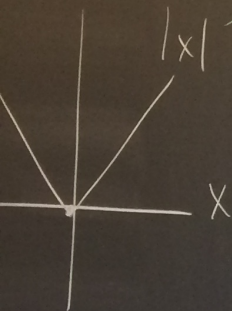
$$b_0 \leftarrow b_0 - \alpha (\vec{b}^T \vec{x}_i - y_i)$$

$$b_j \leftarrow \textcircled{b_j} - \alpha \left[(\vec{b}^T \vec{x}_i - y_i) x_{ij} + \textcircled{\lambda} b_j \right]$$

d+1 coeffs

hyperparameters

$$= \underbrace{(1 - \textcircled{\alpha \lambda})}_{\text{pulling } b \text{ terms toward } 0 \text{ each iteration.}} b_j - \alpha \left[(\vec{b}^T \vec{x}_i - y_i) x_{ij} \right]$$



Analytic

$$J(\vec{b}) = (\vec{X} \vec{b} - \vec{y})^T (\vec{X} \vec{b} - \vec{y}) + \underbrace{\lambda \vec{b}^T \vec{b}}$$



$$\vec{b} = (\vec{X}^T \vec{X} + \underbrace{\lambda \mathbf{I}})^{-1} \vec{X}^T \vec{y}$$

b_0 →

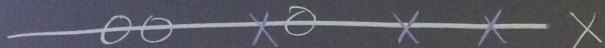
$$\mathbf{I} = \begin{bmatrix} \textcircled{1} & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\underbrace{[\overset{(b_1)}{\vec{b}}]}_{\text{row vector}} \underbrace{\begin{bmatrix} \overset{(b_1)}{\vec{b}} \\ 1 \end{bmatrix}}_{\text{column vector}} = \sum_{j=1}^p b_j^2$$

$$\begin{aligned} & \vec{b} \cdot \vec{b} \quad \text{dot product} \\ &= \vec{b}^T \vec{b} \end{aligned}$$

Outline for February 13

- Polynomial regression
- Regularization
- Linear regression for classification
- Begin: logistic regression



x	likes toy
8	Y
6	N
2	N
5	Y
10	Y
3	N

$$\hat{y} = \begin{cases} 1 & \text{if } \vec{b}^T \vec{x} \geq 0.5 \\ 0 & \text{o.w.} \end{cases}$$

