

Regression = Predicting continuous output

Linear Regression (parametric)

gets the parameter from learning process

hyperparameter \Rightarrow you provide the value not from learning

$$\hat{y} = w_1 x_1 + w_2 x_2 + b$$

$$\hat{y} = \theta_1 x_1 + \theta_2 x_2 + \theta_0$$

* of predictors = * of features + 1

* to evaluate a line \Rightarrow predict on the training set and calculate

the Loss. (How far our prediction from the correct)

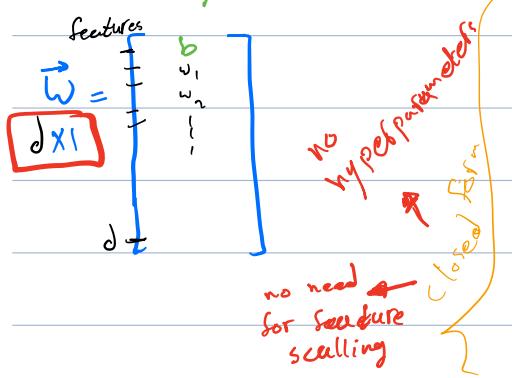
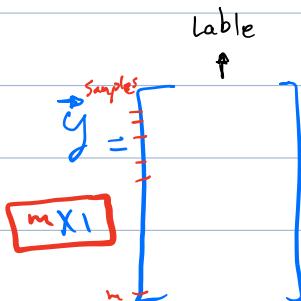
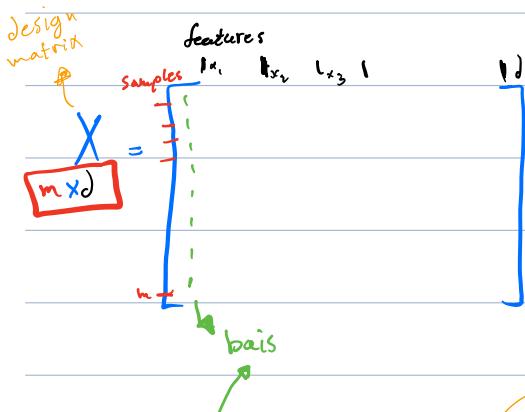
$$\text{loss} = |y - \hat{y}| \quad (L^1)$$

$$\text{cost}(J) = \frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$$

$$\text{RSS loss} = (y - \hat{y})^2 \quad (L^2) \text{ (better)}$$

$$\text{cost}(J) = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

A residual sum of squares (RSS) is a statistical technique used to measure the amount of **variance** in a data set that is not explained by a regression model.



$$\begin{aligned} J &= (\vec{w} - y)^T (\vec{w} - y) \\ &= \vec{w}^T \vec{w} - \underbrace{\vec{w}^T y}_{\text{closed form}} - \underbrace{y^T \vec{w}}_{\text{closed form}} + y^T y \\ \frac{\partial J}{\partial \vec{w}} &= 2 \vec{w}^T \vec{x} - 2 \vec{x}^T y \\ 2 \vec{w}^T \vec{x} - 2 \vec{x}^T y &= 0 \\ \vec{w}^T \vec{x} &= \vec{x}^T y \\ \vec{w} &= (\vec{x}^T \vec{x})^{-1} \vec{x}^T y \end{aligned}$$

very rare way to get \vec{w} ??

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

$m < d$

- Sometimes there is no inverse
- inverse is costly (d^3)

$$A \cdot A^T = A^2$$

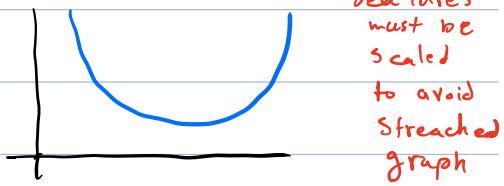
no problems with
big # of d

Local Search (Iterative solution) :-

Gradient descent :-

better than closed-form because the graph is convex (one minimum)

* Step size = learning rate



* big step = overshoot, may fail to converge

* small step = more time

more time to reach the solution

Step 1: Initial w, b randomly

However, the solution doesn't change

Step 2: calculate J (High value initial)

Step 3: calculate $\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

Step 4: $w \leftarrow w - \eta \frac{\partial J}{\partial w}$

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) \cdot x_j$$

$b \leftarrow b - \eta \frac{\partial J}{\partial b}$

Step 5: repeat 2 to 4

$$\text{until } J^K - J^{K+1} \leq \epsilon^{10^{-3}}$$

$$\frac{\partial J}{\partial b} \approx \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)$$

to

$$\frac{\delta}{\delta} =$$

Initially

| Sleep | Study | Score | \hat{y} | Loss | $\frac{\partial J}{\partial w_1}$ | $\frac{\partial J}{\partial w_0}$ | $\frac{\partial J}{\partial b}$ |
|-------|-------|-------|-----------|-------------------|-----------------------------------|-----------------------------------|---------------------------------|
| 6 | 6 | 8 | 0 | 64 | -48 | -48 | -8 |
| 5 | 6 | 8 | 0 | 64 | -48 | -48 | -8 |
| 7 | 8 | 9 | 0 | 81 | -63 | -72 | -9 |
| 7 | 7 | 7 | 0 | 49 | -49 | -49 | -7 |
| | | | | $\bar{J} = 32.75$ | $\frac{-50}{-50}$ | $\frac{-49}{-49}$ | $\frac{-7}{-7}$ |

$$w_1 \leftarrow 0 - 0.01 \times (-50)$$

$$w_1 = +0.5$$

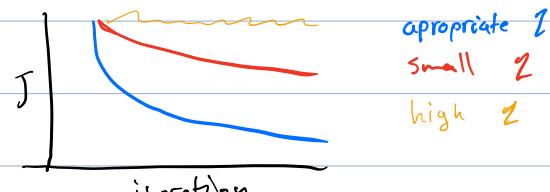
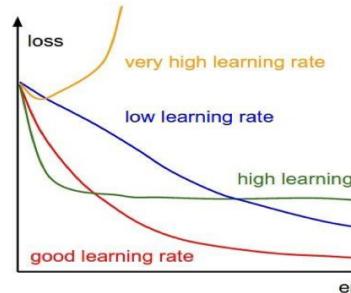
$$w_0 \leftarrow 0 - 0.01 \times (-44.75)$$

$$w_0 = +0.4475$$

$$b \leftarrow 0 - 0.01 \times (-8)$$

$$b = 0.08$$

then update \hat{y}
and repeat



* Polynomial regression

in Linear Regression \Rightarrow we assumed that features have linear relation with predictor

if you have complex relation:- Features expansion

expand before scaling

* Study \rightarrow (* study)²



adding new feature
which is squared (or something else)

then do linear Regression \rightarrow
but plot in linear plot

| | |
|-------|-------|
| x_1 | x_2 |
| → | |
| x_1 | x_2 |

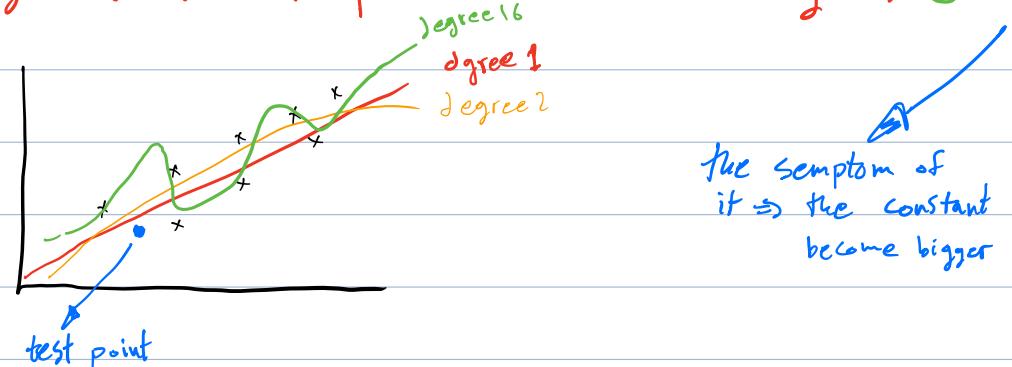
in Real life, no need to expand more than square:-

- not that complicated
- it's hard to know the relationship

✗ overfitting

if you used expanded features
you may have Overfitting \Rightarrow same as when you have lots of features

you will start with line until you get very complicated curve that goes through each point (each label in training set) overfitting

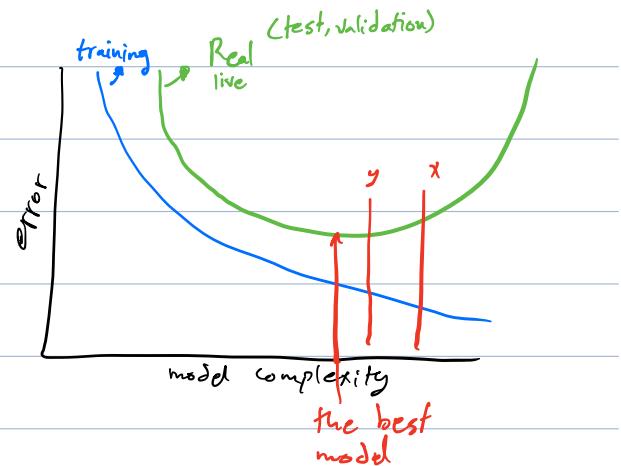


the symptom of
it \Rightarrow the constant with coefficients
become bigger

What is Overfitting ??

why model x is over fitting?
because there is a better model (y) in testing set

The criteria is to reduce error in generalization set



How to deal with Overfitting?? \Rightarrow ① Regularization ② add more training data

Regularization (techniques to deal with overfitting): [Penalize for large value of w]

① Ridge Regularization

→ Ridge Regression

$$\hat{y} = b_0 + b_1 x_1 - 2b_2 x_2$$

→ L^2

for training

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \lambda \|w\|_2^2$$

Regulation hyper parameter
don't include bias

$$\|w\| = \sqrt{b_0^2 + b_1^2}$$

$$\|w\|_2 = b_0^2 + b_1^2$$

we get it using trial and error

if λ is high \Rightarrow underfitting

if λ is 0 \Rightarrow no Regularization (overfitting)

② closed-form

$$J = (Xw - y)^T (Xw - y) + \lambda w^T w$$

Identity matrix

fields of ML

{ define cost
optimize cost
Regularization }

$$w = (X^T X + \lambda I)^{-1} X^T y \Rightarrow \text{by adding } \lambda I \text{ always you can get the inverse}$$

→ but still you have the computation problem when λ is Large

$$w_j \leftarrow w_j - \lambda \frac{\partial J}{\partial w_j} \rightarrow \frac{1}{m} \sum (\hat{y} - y) - x_j + \lambda w_j$$

$$b \leftarrow b - \lambda \frac{\partial J}{\partial b} \rightarrow \frac{1}{m} \sum (\hat{y} - y)$$

$$w_j \leftarrow w_j (1 - \lambda) - \lambda \frac{\partial J}{\partial w_j}$$

Steps

- ① start with simple model
 - ② go more complex ^{→ polynomial} until get the best model
 - ③ go 2 more degree then regularize it
- why? to reduce the bias
and regularization will reduce the variance

The Sources of error:-

- ① bias² (simple model → high bias)

How far the avg model from actual data (complex model → low bias)

- ② variance how much each model different from each other
How far each model from actual

Simple model → low variance
Complex model → high variance

- ③ irreducible error E

- noisy in tables
- lack of information
- Inherence between data

* Ideal \Rightarrow low bias + low variance

Feature Selection:-

* more features with less data set leads to overfitting

ways:-

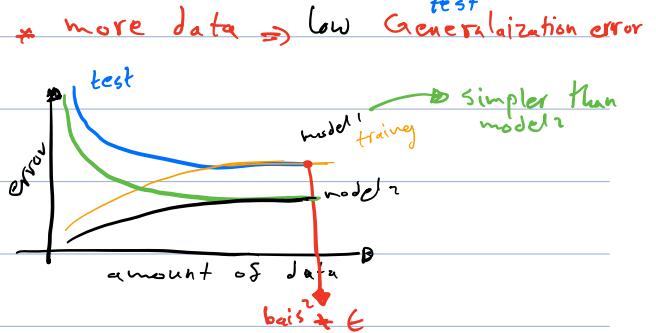
- ① heuristic approach:- (doesnt give the optimal)
Randomly

- ② forward Features selection

Start with 1 feature then add the best feature

- ③ backward Features deduction

Start with all features then delete the worst feature

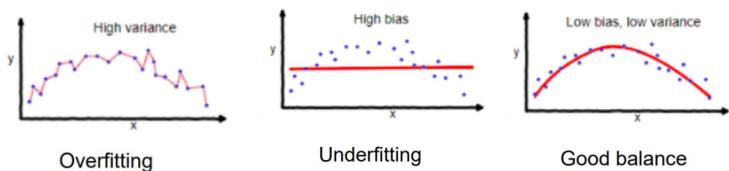


* theoretically \Rightarrow you can have 0 variance but you cant have 0 bias

All of them are slow

The better way is

Lasso



② Lasso Regularization

$$J = \frac{1}{2m} \sum (\hat{y} - y)^2 + \lambda \|w\|_1$$

more $\lambda \rightarrow$ more features become zero

→ don't have closed-form solution or gradient

→ we solve it by coordinate descent

→ it will do feature selection automatically

→ so it's deal with overfitting

□ Lasso tends to select amongst highly correlated features arbitrarily randomly (without any reason)

□ Empirically, ridge has better predictive performance

to get the best solution:-

then
Lasso \Rightarrow Ridge or Elastic net

combin Lasso and Ridge)

$$J = \frac{1}{2m} \sum (\hat{y}^i - y^i)^2 + \lambda (\alpha \|w\|_2^2 + (1-\alpha) \|w\|_1)$$

□ Perform training and check the error:

□ If error is high: (underfitting) [when $d > m$]

➤ Add more features

➤ More complex model by Polynomial

□ If overfitting:

➤ Add more data

➤ Perform regularization

Bias and variance

