

Feature scaling and normalization continued ...

We prepare 程序代写代做 CS 编程辅导 features to have dynamic all

range [-1,



Step1: "Mean normalization": find the maximum likelihood

WeChat: cstutorcs

$\underset{k \text{th training instances}}{\overbrace{x_i^{(k)}}}$

ML mean for each feature and

subtract it for all training instances ... Except the Bias!

Assignment Project Exam Help

$$x_{i,\text{normal}}^{(k)} = x_i^{(k)} - \hat{\mu}_{ML,i} \quad \text{for instances } k=1, \dots, N$$

actual
training
feature (i)

Email: tutorcs@163.com
For feature $i=1, \dots, d$
(not bias!)

QQ: 749389476

Note: This makes the average value of all features equal to zero:

$$E_{x_i|w_i} \left[x_{i,\text{normal}}^{(k)} \right] = E_{x_i|w_i} \left[x_i^{(k)} - \hat{\mu}_{ML,i} \right]$$

$$\begin{aligned} &= E_{x_i|w_i} \left[x_i^{(k)} \right] - \underbrace{\hat{\mu}_{ML,i}}_{=} \\ &= E_{x_i|w_i} \left[x_i^{(k)} \right] \end{aligned}$$

since ML mean is unbiased

= 0

Step2: "Feature Scaling": find the ML variance, take its square root to get the standard deviation and divide the features.

$$x_{i,scaled,n}^{(k)} = \frac{x_{i,normal}^{(k)} - \hat{\mu}_{ML,i}}{\sqrt{\hat{\sigma}_{ML,i}^2}}$$

for instances $k=1, \dots, N$
for features $i=1, \dots, d$

WeChat: cstutorcs
(NO Bias!)

new Scaled normalized features
all fall more or less in interval [-1,1]

NOTE: If you have enough training data, then what should happen is that weights for useless attribute (like pupil diameter) will be zero or near zero.

↳ we can ditch useless features if their weights are zero

What if N , my number of training data, is small?

And if d is small?

↳ then $\frac{\partial J}{\partial w}$, set that equal to zero, then solve for \vec{w}_{MLE}

$$\Rightarrow \vec{w}_{\text{MSE}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{P}$$

You get

程序代写代做 CS 编程辅导

↳ vector of training labels

$\vec{X}^T \vec{X}$ is big matrix
if many feature

$\vec{X}^T \vec{X}$ is big matrix

\Rightarrow inverting it is expensive



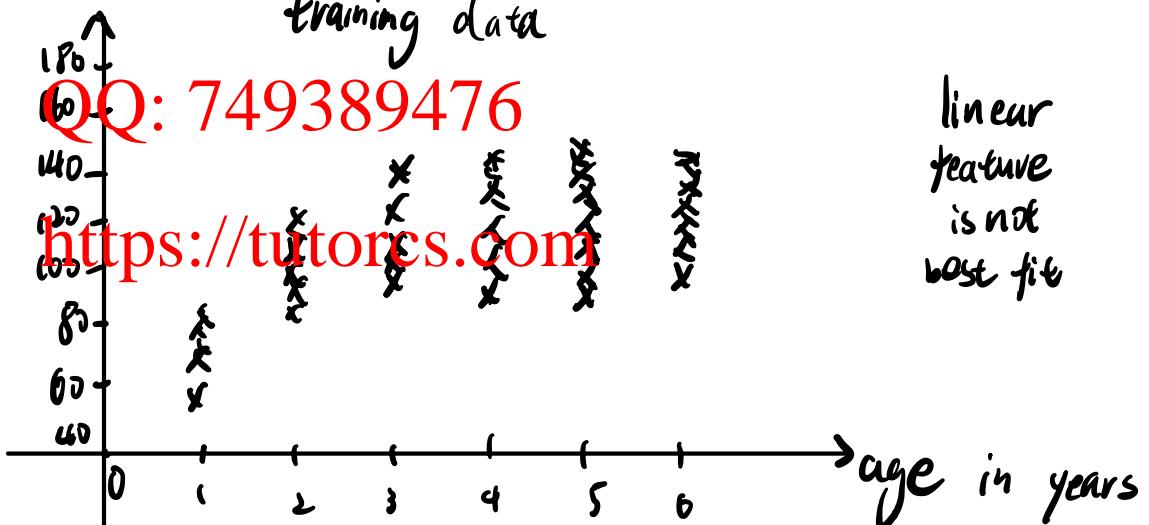
Non-Linear Polynomial Regression

Ex] Say we have WeChat: cstutorcs (age of the fish in years)
to predict its length in mm.

Assignment Project Exam Help

(Not classification, but gradient) training data

Email: tutorcs@163.com



we could try

$$\vec{x} = \begin{bmatrix} 1 \\ x_1 \\ x_1^2 \end{bmatrix} \Rightarrow \text{prediction} = \vec{w}^T \vec{x} = w_0 + w_1 x_1 + w_2 x_1^2$$

or not polynomial: $\vec{x} = \begin{bmatrix} 1 \\ \sqrt{x_1} \end{bmatrix}$

NOTE: when you use higher order polynomial, the dynamic range gets more pronounced

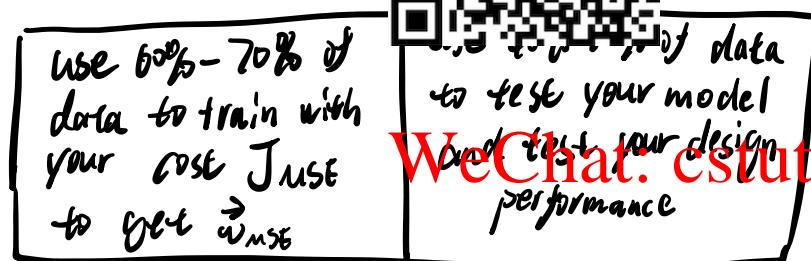
Ex if $x_1 = 50 \Rightarrow x_1^2 = 50^2$ bigger

if $x_1 = 0.1 \Rightarrow x_1^2 = 0.1^2$ smaller

\Rightarrow normalizing 程序代写代做 CS 编程辅导

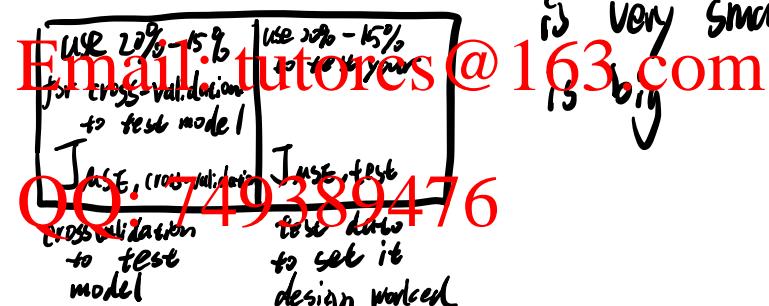
How do we choose model to use?

To test your model's performance:



why? Because in cross-validation you may overfit your model to that cross-validation data.

Assignment Project Exam Help

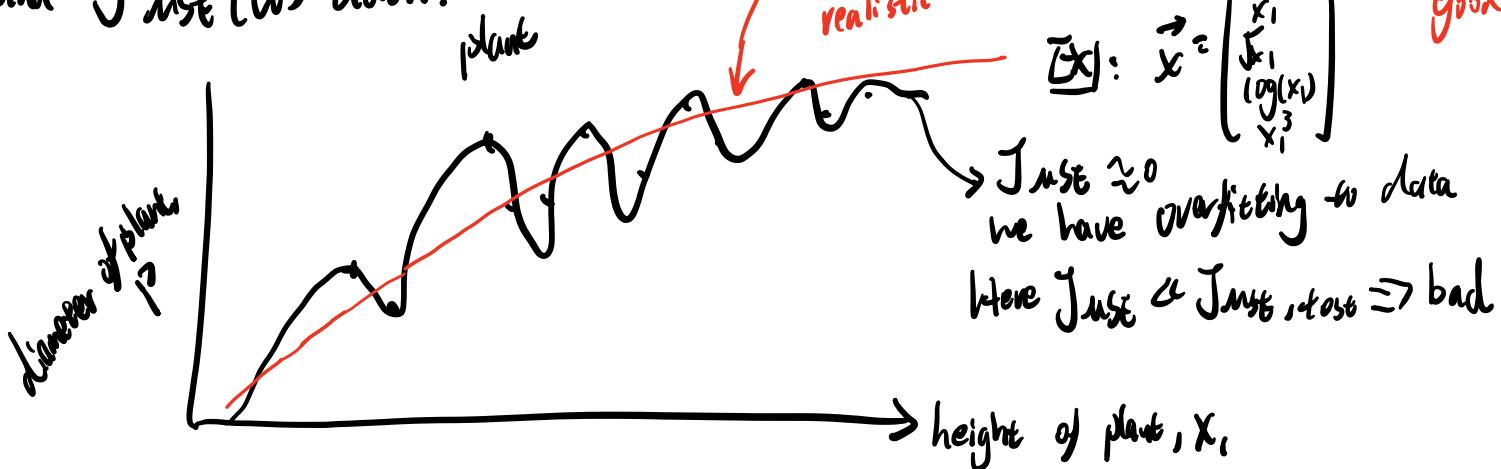


you will see that J_{USE}, cross-validation is very small, but J_{USE, test} is big
this is the one we most care about

<https://tutorcs.com>

Avoid overfitting: Regularization

If your $J_{USE}(\vec{w})$ is high, you can invest a lot of time experimenting with high order polynomials that bring J_{USE} , cross-validation and $J_{USE}(\vec{w})$ down.



When you have many features, it is difficult to know how to simplify the model to reduce overfitting.

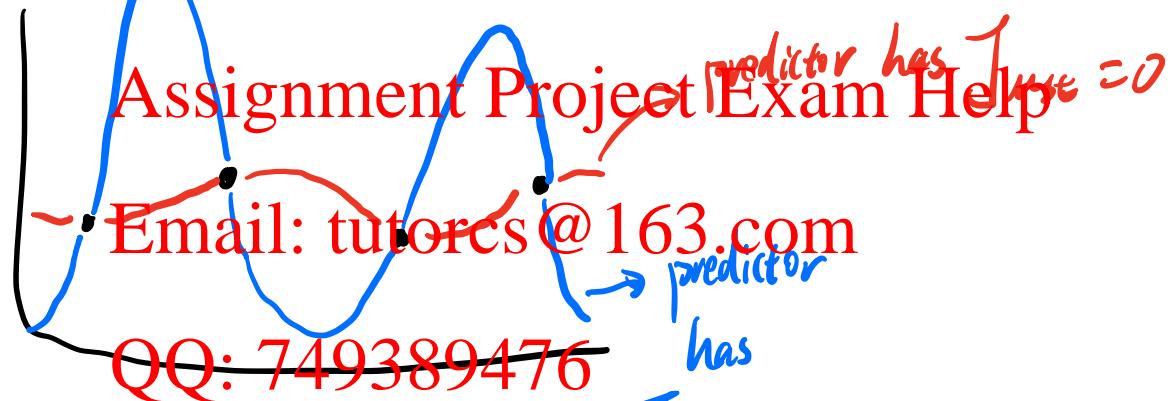
程序代写代做CS编程辅导

Step 1: discard feature that are redundant.

Step 2: "REGULARIZE the weight vector", \vec{w}



fit example & points
WeChat: cstutorcs



Regularization means

"punishing" choice of \vec{w} that has large weight values.

We do this by including it in the cost function.

$$J_{\text{MSE-W-regularization}}(\vec{w}) \stackrel{\text{def}}{=} \frac{1}{N} \left\{ \sum_{k=1}^N \underbrace{(h_{\vec{w}}(\vec{x}^{(k)}) - p_k)^2}_{\text{prediction}} + \lambda \sum_{j=1}^d w_j^2 \right\}$$

$$= J_{MSE} + \frac{\lambda}{N} \overbrace{\| \text{weight without bias} \|}^2$$

Def $\lambda \geq " \text{程序代写代做CS编程辅导}$

To find best \vec{w}  get case:

$$\underbrace{\partial J_{MSE - w - \text{reg}}}_{\partial w} = \frac{2}{N} \sum_{k=1}^K (\vec{w}^T \vec{x}^{(k)} - p_k) \cdot \vec{x}_i^{(k)} + \frac{2\lambda}{N} \sum_{i=1}^d w_i$$

∂w_i :

WeChat: cstutorcs

so \sum that minimizes $J_{MSE - w - \text{regularization}}$ (with regularization) is

$$\vec{w}_{k+1} = \vec{w}_k - \eta \left[\frac{2}{N} \sum_{j=1}^N (\vec{w}_k^T \vec{x}^{(j)} - p_j) + \frac{2\lambda}{N} w_i \right]$$

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

$$= \left[\begin{array}{c} w_0 (1-\eta) \\ w_1 (1 - \frac{2\lambda\eta}{N}) \\ \vdots \\ w_d (1 - \frac{2\lambda\eta}{N}) \end{array} \right] - \eta \left[\begin{array}{c} \frac{2}{N} \sum_{j=1}^N (\vec{w}_k^T \vec{x}^{(j)} - p_j) \\ \frac{2}{N} \sum_{j=1}^N (\vec{w}_k^T \vec{x}^{(j)} - p_j) \vec{x}_1^{(j)} \\ \frac{2}{N} \sum_{j=1}^N (\vec{w}_k^T \vec{x}^{(j)} - p_j) \vec{x}_d^{(j)} \end{array} \right]$$

If we wanted a closed form solution

with Regularization, (for small d), then 

because:

$$\vec{w}_{MSE - w - \text{regularization}} = (\vec{X}^T \vec{X} + \begin{bmatrix} 0 & \lambda & & \\ 0 & 0 & \ddots & \\ & & \ddots & \lambda \end{bmatrix})^{-1} \vec{X}^T \vec{P}$$

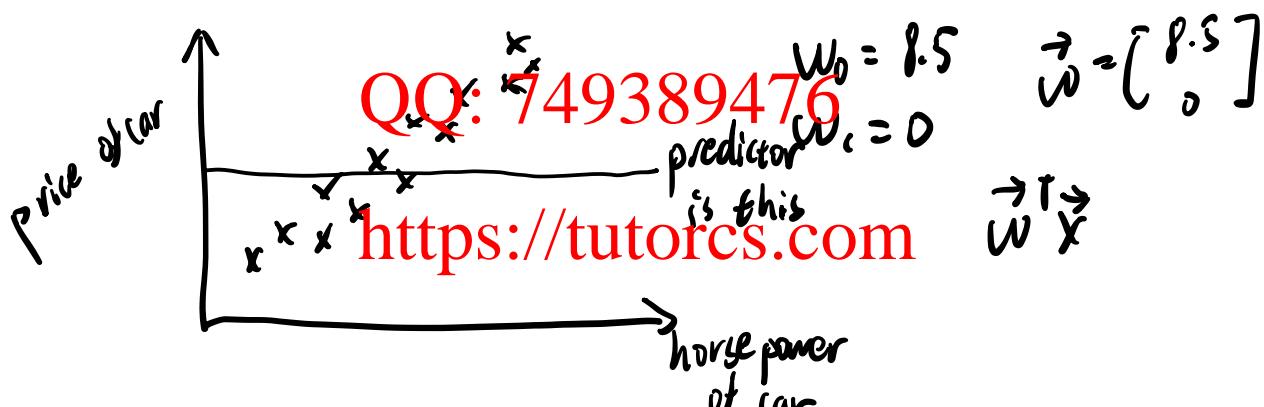
HIGH BIAS VS HIGH VARIANCE TRADE-OFF

Ex] car price prediction → say after mean normalization, scaling, and tuning λ (regularization parameter), your J_{reg} remains high. our gradient descent coverage

Typically:  is that you are either underfitting or overfitting
WeChat: cstutorcs

Def] A model suffers from "HIGH Bias" if it underfits the data.

Email: tutorcs@163.com



This model underfits yielding bad J_{reg, train} and J_{reg, test}

This example ignores horsepower

Note: collecting more data will not solve a high bias problem

How do we solve high bias (overfitting)?

- add more features ($x_i^2, x_i^3, \sqrt{x_i}, x_i \cdot x_2, x_i \cdot x_2 \cdot x_3$, etc...)

- see if you add more attributes
- if you are underfitting, then you may be regularizing too much, 程序优与代做CS编程辅导



→ to select λ

rep1: select a λ that scales over orders of magnitude

$$\lambda \in \{0, 0.001, 0.1, 1\}$$

Def A model suffers from "High Variance" if it overfits the data

Assignment Project Exam Help



To detect high variance, we'll see $J_{MSE, train}$ is small

and $J_{MSE, test}$ is big

How to solve high Variance (overfitting)?

- collect more training data (make N bigger).
(Do NOT use your test data).
- ditch some non-linear attributes (x_1^2, x_1^3 , etc)

- if you are overfitting, then you are not regularizing enough.
so make λ bigger

BIAS - VARIANCE TRADE-OFF

Def "general



Tutor CS
Assignment Project Exam Help

"expected test error"

measure of how accurately an algorithm is able to predict outcomes for the ~~test~~ data (not used in training)

WeChat: cstutorcs

$$\cdot D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

= set Email: tutorcs@163.com
training instances with features and labels

Note: for any feature ~~749389476~~ may have more than one label

<https://tutorcs.com>

Ex x_i has label $P_i = 7$ } there is an underline
 x_j has label $P_j = 7$ } conditional PDF $f(y|x)$

Ex two plants may have exactly the same height but different widths.

The most likely or expected width is $E_{y|x}[y]$

Hence: there is a joint distribution $f(x, y)$

- $h_0(x)$ = our imperfect prediction for y
 = prediction hypothesis (say for linear regression
 or for logistic regression), as a function
 of vector x , and h was trained from
 all x 's in our limited set D .
 $\Rightarrow h$'s must also be random

- $\bar{h}(x) \stackrel{\text{WeChat: tutorcs}}{=} \int h_0(x) f(D) dD$
 = expected classifier or predictor using
 Assignment Project Exam Help
 same imperfect algorithm

but with Email: tutorcs@163.com

Def we can define the generalized error of our
 imperfect algorithm <https://tutorcs.com>

$$E_{x,y,D-\text{trained}}[(h_p(x) - y)^2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (h_p(x) - y)^2 f(x, y, D) dy dx dD$$

- $\vec{y}(x) = \text{MAP estimate of label for } x$, which is the very
 best estimate in the world, but it requires us
 to know the underlying pdf of y given x : $f(y|x)$

$$= E_{Y|X}[y] = \int_{-\infty}^{\infty} y f(y|x) dy = \text{MAP estimator} = \text{Expected label}$$

FACT: You can show that the generalized error has three components that shows trade-off between bias and variance

Track #1 and subtract expected classifier

$$E_{x,y,D-\text{trained}}[(h_0(x) - \bar{h}(x))^2] = E_{x,y,D-\text{trained}}[(\{h_0(n) - \bar{h}(x)\})^2 + \{\bar{h}(x) - y\}^2]$$

WeChat: cstutorcs
= algebra, probabilities, cross-term is zero

$$= E_{x,D-\text{trained}}[(h_0(x) - \bar{h}(x))^2] + E_{x,y}[(\bar{h}(x) - y)^2]$$

Email: tutorcs@163.com

QQ: 749389476 "variance of $h_0(x)$ "
exactly definition

The classifier <https://tutorcs.sovance.com> shows how sensitive your classifier/predictor is to changing from one set D to one other set D'. How much are you overfitting your h to your predictor D?

How far off are we from the average classifier that uses all the data in the world and not our limited D?

Next: Trick #2: add and subtract the MAP estimate

$\bar{y}(x)$ of label for feature x :

程序代写代做 CS 编程辅导

$$\left[\sum_{x,y,D-\text{train}} \left[(\hat{h}_0(x) - y)^2 \right] \right] = \text{Var of } h_0(x)$$



$$\left[\left\{ \hat{h}(x) - \bar{y}(x) \right\} + \left\{ \bar{y}(x) - y \right\} \right]^2$$

... algebra, probability, cross term goes to zero

WeChat: cstutorcs

$$= \underbrace{\text{Var of assignment}}_{\text{Assignment}} + \underbrace{\left[\sum_x \left[(\hat{h}_0(x) - \bar{y}(x))^2 \right] + \sum_{x,y} \left[(\bar{y}(x) - y)^2 \right] \right]}_{\text{Project Exam Help}}$$

Email: tutorcs@163.com

"Estimator"

"noise" of estimator"

QQ: 749389476

<https://tutorcs.com>