

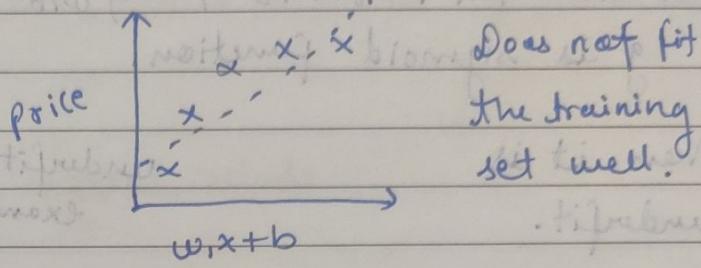
11. The Problem of Overfitting. (underfit vs overfitted)

a. Overfitting and Underfitting.

Regression example:

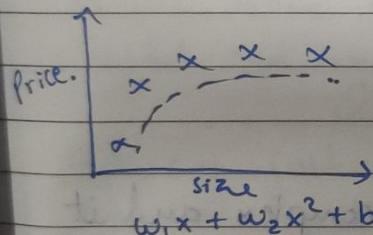
The model is underfit.

the model has high bias.



Underfit - Not able to fit the data.

High bias - Pre conception that the data is linear.

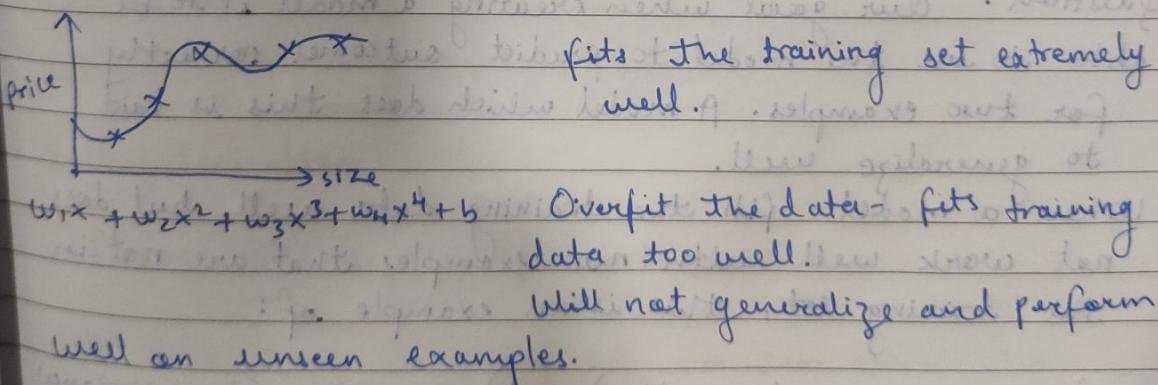


- fits training set pretty well.
Just right.

generalization - We want our model to be generalized well.

It should perform well on examples.

that it has never seen before.



It has high variance. (Same as overfit).

Trying very hard to fit every training example.

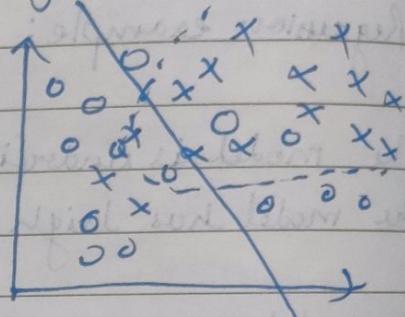
A well generalized model (just right) is a good model.

Classification Problem

for example, if we have our model as $z = w_1x_1 + w_2x_2 + b$, $f_{\text{WIB}}(\vec{x}) = g(z)$.
 g is sigmoid function.

Then it is underfit.

underfit example.



$$\text{gf, } z = w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + b$$

Then, it is just right model fit.

$$\text{if } z = w_1x_1 + w_2x_2 + w_3x_1^2 + \dots + b$$

Then, the decision boundary will be complex and it will overfit.

C. Question. Our goal when creating a model is to be able to use the model to predict outcomes correctly for two examples. A model which does this is said to generalize well.

When a model fits the training data well but does not work well with new examples that are not in the training set, this is an example of:

- a) Underfitting (high bias)
- b) Overfitting (high variance)
- c) A model that generalizes well.

Answer - b.

II. Addressing Overfitting

- a. Collect more training examples.
If we get more data, the model will fit a f^n that is less wiggly.
- b. Select features to include/exclude.
e.g. size, bedrooms, floors, age, avg. income, distance to coffee shop, price (y).
we all \rightarrow features.
if we have all features and insufficient data, our model can again overfit.

We can just select few features, like size, bedrooms, and age to get it just right.
This is called feature selection.

This advantage:

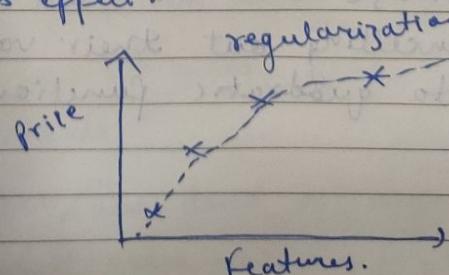
We are throwing away information.

c. Regularization

Reduce the size of parameters w_j .

We had large values of w_j , we can eliminate a parameter to reduce its effect.

But let's us keep features.



small values of w_j .

$$f(x) = 13x - 0.23x^2 + 0.0000041x^3 - 0.00001x^4 + 10.$$

d. Summary: Addressing Overfitting.

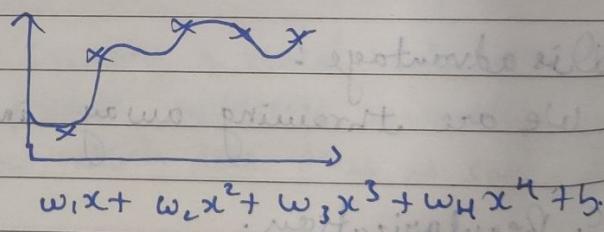
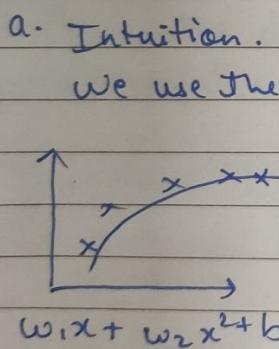
- Options - Collect more data (e.g. prior knowledge)
- Select features or set threshold to avoid overfitting
- Reduce size of parameters using ~~Regularization~~
Regularization.

e. Question: Applying regularization, increasing the number of training examples, or selecting a subset of the most relevant features are methods for:

- a) Addressing underfitting.
- b) -- Overfitting.

Answer - b. Underfitting has smooth line and overfitting has wavy line.

III. Cost function with Regularization.



make w_3, w_4 really small (≈ 0)

$$\text{choose } \min_{w_1, b} \frac{1}{2m} \sum_{i=1}^m (f_{w, b}(x^{(i)}) - y^{(i)})^2 + 1000w_3^2 + 1000w_4^2$$

if we make them small, we minimize their cost. We are reducing the values. We make them closer to quadratic function.

b. Regularization.

smaller values of w_1, w_2, \dots, w_n, b , we have a simpler model which is less likely to ~~overfit~~ overfit.

More generally, if we have a lot of features, we ~~can't~~ penalize all of our features! say n features, we penalize w_1, w_2, \dots, w_n .

Our new cost function should look like this.

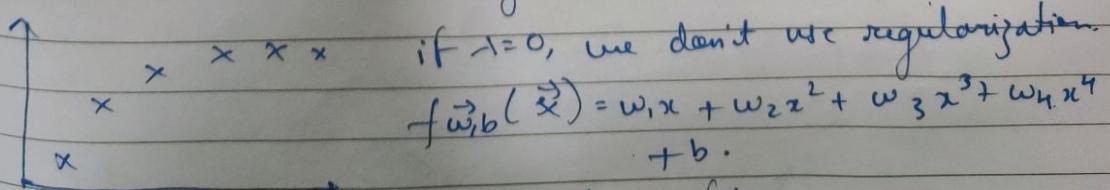
$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\lambda = \text{lambda}$ = is called regularization parameter
 $\lambda > 0$.

By scaling both terms by $2m$, we have better way to calculate value of λ . We don't penalize b . We can include or exclude b .

c. Regularized cost function.

$$\min_{w, b} J(\vec{w}, b) = \min_{w, b} \left[\underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2}_{\text{fit data.}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularized term.}} \right]$$

we balance these two goals.

 if $\lambda = 0$, we don't use regularization.
 $f_{\vec{w}, b}(\vec{x}) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$.
 This will overfit.

if $\lambda = 10^{10}$, we choose $w \approx 0$, and $f(x) = b$.
This will underfit.

We need to choose just right λ .

d-Quiz : For a model that includes the regularization parameter λ (lambda), increasing λ will tend to:

- a) Decrease w_1, w_2, \dots, w_n
- b) Increase w_1, w_2, \dots, w_n
- c) Increase b
- d) Decrease b

Answer: a).

IV. Regularized Linear Regression.

a. Regularized Linear Regression.

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^n w_j^2.$$

Gradient Descent : repeat { $w_j = w_j - \alpha \frac{\partial J(\vec{w}, b)}{\partial w_j}$ } $j = 1, \dots, n$.

$$b = b - \alpha \frac{\partial J(\vec{w}, b)}{\partial b}$$

} simultaneous update.

New, $\frac{\partial}{\partial w_j}$ was $\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

and $\frac{\partial}{\partial b}$ was $\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})$

With added regularized term, we will have $\frac{\lambda}{m} w_j$
for w_j . b remains the same.

b. Implementing gradient descent:

$$\text{repeat } \left\{ \begin{array}{l} w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right] \\ \quad + \frac{\lambda}{m} w_j \end{array} \right. \\ b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}). \quad \left. \right\} \text{ simultaneous update.}$$

Now, $w_j = 1 \cdot w_j - \alpha \frac{\lambda}{m} w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

$\boxed{w_j \left(1 - \alpha \cdot \frac{\lambda}{m} \right)}.$ usual update.

Now, α is a small number, λ is a small number, $\alpha(1 + \lambda)$, $\alpha(0.1)$, and m is 50.

$$\text{Now, } \frac{\alpha \lambda \alpha}{m} = \frac{0.01 \times 1}{50} = 0.0002. 1 - = (1 - 0.0002)^{50}$$

$$\text{Now } (1 - 0.0002) = 0.9998.$$

We are multiplying it by a small number. We are restricting the value of w_j in every iteration.

c. How we get the derivative term ($-w = jw$)

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{\partial}{\partial w_j} \left[\frac{1}{m} \sum_{i=1}^m (f(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^n w_j^2 \right].$$

$$= \frac{1}{m} \sum_{i=1}^m \left[(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \times \cancel{\vec{x}} x_j^{(i)} \right] + \frac{\lambda}{m} \cancel{\vec{x}} x_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m \left[(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \cdot \cancel{x_j^{(i)}} \right] + \frac{\lambda}{m} w_j$$

$$= \frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \cdot \cancel{x_j^{(i)}} \right] + \frac{\lambda}{m} w_j$$

d. Question: Assuming α , the learning rate, is a small number like 0.001, t is 1, and $m=50$, what is the effect of the 'new part' on updating w_j ?
 a) decreases value of w_j b) increases value of w_j

Answer - a.

IV. Regularized Logistic Regression.

Now, we have $Z = w_1x_1 + w_2x_2 + w_3x_1^2x_3 + \dots + b$.

$$f_{w,b}(\vec{x}) = \frac{1}{1+e^{-Z}}$$

a. New cost function:

$$\begin{aligned} J(\vec{w}, b) = & -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{w,b}(\vec{x}^{(i)})) + (1-y^{(i)}) \times \right. \\ & \left. \log(1-f_{w,b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \quad \min_{\vec{w}, b} J(\vec{w}, b). \end{aligned}$$

Now, using gradient descent, repeat

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$j=1, \dots, n, \quad b = b - \frac{\partial}{\partial b} J(\vec{w}, b)$$

$$\text{Now, } \frac{\partial}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m [(f_{w,b}(\vec{x}^{(i)}) - y^{(i)}) \cdot x_j^{(i)}] + \frac{\lambda}{m} w_j$$

$$\frac{\partial}{\partial b} = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)})$$

$f_{w,b}(\vec{x}^{(i)})$ is the logistic function.