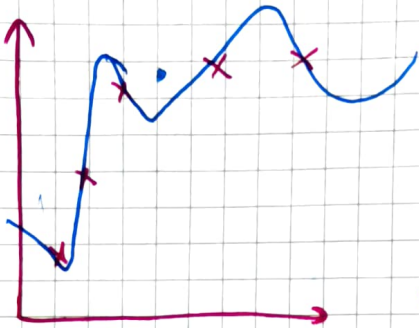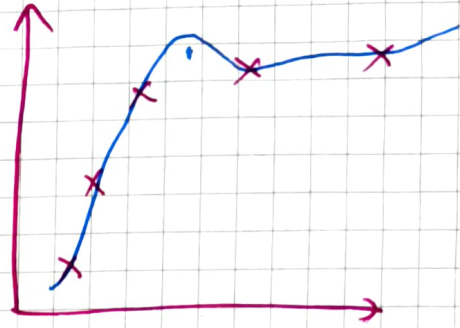**Regularization :** Gentle way to reduce the impact of the feature without being harsh i.e. eliminating the feature Outrightly.

↳ Way to Shrink the value of the parameters without necessarily demanding that parameters may be set to 0 exactly

☺ Even if you fit a higher Order polynomial, using the small value of parameters will end up in a curve which is Still better



$$\phi(x) = 28x - 385x^2 + 39x^3 - 174x^4 + 100$$

$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - 0.0001x^4 + 10$$

**Regularization :** Keeps all features but prevents the features from having an Overall large effect (which may cause Overfit)

↳ Regularized Model is Simpler

↳ Reduces Overfitting by penalizing wts & thus minimizing the Complexity of model

↳ Reduces generalization error / prediction Error

Very weights of higher degree Poly will reduce the Complexity of curve/model

---

**Ques** What is curse of dimensionality ?

If D = 1    $y = wx + b$

If D = 2    $y = w_0, x_1 x_2 + w_2 x_1^2 + w_3 x_1^2 + b$

If D = 3    $x_1 x_2$    $x_2 x_3$    $x_2^2$    $x_1^2 x_2$    $x_3^2 x_1$
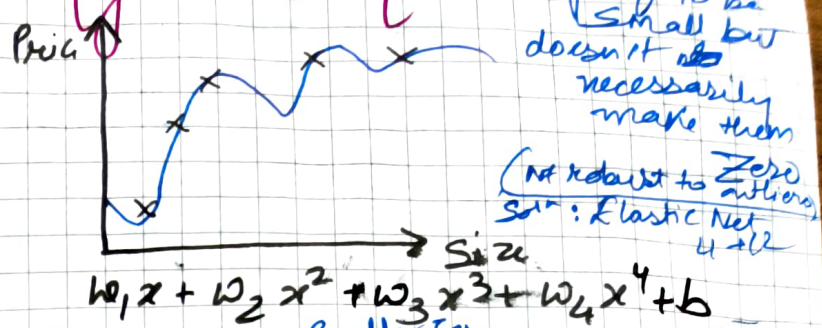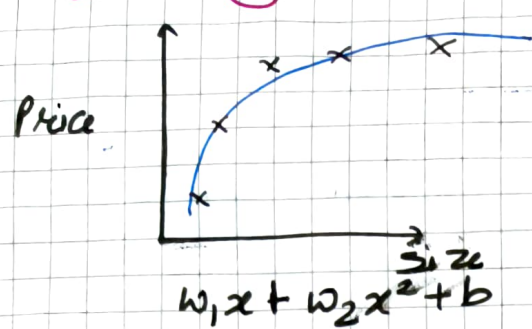                $x_1 x_3$    $x_1^2$    $x_3^2$    $x_2^2 x_3$

As D ↑, no. of Independent Coeff ↑, thus to, Capture Complex dependencies in data, We may need to use higher Order Poly^n.

# Cost Function with Regularization

*L2 reg: forces wts to be Small but doesn't necessarily make them Zero (Not robust to Zero solutions) Sol^n: Elastic Net L1 + L2



Price

$w_1 x + w_2 x^2 + b$

Price

$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

→ Penalty Term
Regularization Parameter (How many)

Make $w_3$ & $w_4$ really Small ($\approx 0$)

$$\min \frac{1}{2N} \sum \left(y^{(i)} - f(x^{(i)})\right)^2 + 1000 w_3^2 + 1000 w_4^2$$

modified Cost

→ any large no. will be ok → penalize them

⇒ The func. will penalize the model if $w_3$ & $w_4$ are large → **Penalizes Higher Values/terms**
∴ If we want to min. this Cost func., the Only way to make it Small is to Set $w_3$ & $w_4$ Small enough
Otherwise → 2 terms will be really really Big

⇒ ( $w_3$ & $w_4$ ) Will be close to ( 0 ) say 0.0006
⇒ nearly have an impact of getting rid of $x^3$ & $x^4$ which is equivalent to
⇒ will end up fitting curve quadratic found on left (∵ left fit is much better fit)

**Idea Behind Regularization** → If there are Smaller Values of parameters we will have a Simpler model

⤷ May be One with fewer features
⤷ Less prone to Overfitting

Regularize/Penalize Weights/parameters

In general, when we have d features, we don't know which features to penalize, Thus we penalize all wj para etc
⇒ Will fit a Simpler, Smoother curve that's less prone to Overfitting

L2 regulari

**Updated Cost func:**

$$J(w,b) = \frac{1}{2N} \sum_{i=1}^{N} \left(y^{(i)} - f(x^{(i)})\right)^2 + \frac{\lambda}{2N} \sum_{j=1}^{d} w_j^2 + \frac{\lambda b^2}{2N}$$

→ Regularization Parameter

So that both terms are Scaled similarly
& it turns out that Scaling both terms Similarly will help you to Chose better Value of $\lambda$

LiReg/Lasso : $\lambda$ ||w|| penalize absolute Value
L2 Reg/Ridge Reg & w2

**Objective :** To min. first term i.e mean Sq. error
 ↳ Encourages algo to fit data well by minimizing
        diff b/w predictions & actual Values

     To min. Second term
      ↳ To Keep parameters $w_i$ small, which will
   tend to Reduce Overfitting

  $\lambda$ : Controls relative Imp./relative Trade off / how you balance b/w
       2 goals

**Ques.** Why $\lambda$ should not be too small Or too high?
**If** $\lambda = 0$ ⟹) No regularization
              ⟹) Fit Overly wiggly, Overly complex curve.
              ↳ Effect of overfitting data
                 Less Traning Error, Large Generalization/
                                                      Prediction error

**If** $\lambda \to$ Very Very large ⟹ Placing heavy weight on regularization
                                term
                      ⟹ Learning algo will choose $w_1, w_2, \ldots$
                      ⟹ to be extremely close to 0.
                      ⟹ $f(x^{(i)})$ is basically equal to $b$    horizontal
       Large Traing Error         ↳ Learning algo will fit a Straight
       & Large Generalization        line & it underfits
              Prediction Error

Choose $\lambda$ that balances
first & second terms of tradeoff
   ↳ Min. mse & keep the parameters small

---

**Regularized Linear Regression** (Exactly Same for logistic regr)       $\frac{d}{dw_3} \left[ w_1^2 + w_2^2 + \ldots w_3^2 \right]$

$$\min_{w,b} J(w,b) = \min_{w,b} \left[ \frac{1}{2N} \sum_{P} \left( y^{(i)} - \underset{wx+b}{f(x^{(i)})} \right)^2 + \frac{\lambda}{2N} \sum_{j=1}^{d} w_j^2 \right]$$

**Gradient Descent**
                                        $\frac{1}{N} \sum_{i=1}^{N} \left( f(x^{(i)}) - y^{(i)} \right) * x_j^{(i)} + \frac{\lambda}{N} w_j$

repeat
  $w_j = w_j - \alpha \left( \dfrac{\partial J(w,b)}{\partial w_j} \right)$

  $b = b - \alpha \left( \dfrac{\partial J(w,b)}{\partial b} \right)$       $\frac{1}{N} \sum_{i=1}^{N} \left( f(x^{(i)}) - y^{(i)} \right)$       remember
                                                                                 we don't
                                                                                 regularize
                                                                                 $b$

Why regularization has effect of shrinking $w_j$

consider 1st & 3rd term

$$= w_j - \frac{\alpha * \lambda}{N} w_j \quad \text{—Usual update} \quad = w_j \left(1 - \frac{\alpha \lambda}{N} \xrightarrow{\;0.01\;} 10\right) \quad \text{—Usual update}$$

$$\underbrace{\qquad\qquad\qquad}_{\text{Shrink } w_j}$$

ADV $\Rightarrow$ Regularized model is simpler since it has less features
$\Rightarrow$ $\downarrow$ reduces Overfitting & thus
$\Rightarrow$ Reduces the generalization error / prediction error

$$\boxed{J = \frac{1}{2N} \sum_{i=1}^{N} \left(y^{(i)} - \hat{y}^{(i)}\right)^2 + \frac{\lambda}{2N} \sum_{j=1}^{N} w_j^2}$$

$\to w_1 x_1 + w_2 x_2 + \ldots w_d x_d + b$

**Gradient Descent for**

| Linear Regression | Regularized Linear Regr |
|---|---|
| repeat until Convergence | |
| $\{$ $b = b - \alpha \frac{1}{N} \sum_{i=1}^{N}\left(\hat{y}^{(i)} - y^{(i)}\right)$ | $\{$ $\to$ SAME |
| $w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^{N}\left(\hat{y}^{(i)} - y^{(i)}\right) x_j^{(i)}$ | $w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^{N}\left(\hat{y}^{(i)} - y^{(i)}\right) x_j^{(i)}$ |
| | $+ \frac{\lambda}{N} w_j$ |
| $\}$ | $\}$ |

$$\boxed{J = \frac{-1}{N} \sum_{i=1}^{N} \left(y^{(i)} \log\left(f(x^{(i)})\right) + (1-y^{(i)}) \log\left(1-f(x^{(i)})\right)\right) + \frac{\lambda}{2N} \sum_{j=1}^{d} w_j^2}$$

$$= \frac{1}{1 + e^{-\left(w_1 x_1 + w_2 x_2 + \ldots w_d x_d + b\right)}}$$

**Gradient Descent for**

| Logistic Regression | Regularized Logistic Regr |
|---|---|
| repeat until Convergence | |
| $\{$ $b = b - \alpha \frac{1}{N} \sum \left(f(x^{(i)}) - y^{(i)}\right)$ | $\{$ $\to$ SAME |
| $w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^{N} \left(f(x^{(i)}) - y^{(i)}\right) x_j^{(i)}$ | $w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^{N} \left(f(x^{(i)}) - y^{(i)}\right) x_j^{(i)}$ |
| | $+ \frac{\lambda}{N} w_j$ |
| $\}$ | $\}$ |

<span style="color:orange">**Orange**</span> : Optional

$$T = \frac{\lambda}{2N} \left(w_1^2 + w_2^2 + w_3^2 + \ldots w_d^2\right)$$

$$\frac{\partial T}{\partial w_2} = \frac{\lambda}{2N} * 2 w_2 = \frac{\lambda}{N} * w_2$$