



LINEAR REGRESSION

Closed Form
Equation

A direct mathematical formula for computing best fit parameters

It minimize the cost func
(MSE)

↓ by
Solving a matrix equation

Gradient Decent

An iterative optimization algorithm that starts with random parameters and gradually improves them.

↓
It minimize the cost func
(MSE)

↓ by
Calculating the gradient of cost func

↓
Update parameters in opp dir

NORMAL EQN

Linear Regression Equation:

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

\hat{y} = predicted Value

n = number of features

x_i = i th feature

θ_j = j^{th} model parameter

θ_0 = bias term / intercept term

\therefore Linear Regression more specifically
predicts weighted sum of inputs

$$\hat{y} = h_{\theta}(x) = \theta \cdot x = \theta^T x$$

Vectorized form of Linear Regression

θ : model's parameter vector
having bias term θ_0 and feature weights
 θ_1 to θ_n

x : instance's feature vector
containing x_0 to x_n , with $x_0 = 1$.

h_{θ} : hypothesis function

To find parameters we will minimize RMSE or
say MSE as it leads to same result

$$MSE(x_i, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

to find the value of θ which minimize the cost function

Normal Equation

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

$\hat{\theta}$ value that minimize the cost fun^c
 y is the vector of target values containing
 $y^{(1)}$ to $y^{(m)}$

Implementation of linear Reg using Normal Eqn:

```
import numpy as np

# Sample Data
# X: Feature matrix (m x n), y: Target vector (m x 1)
X = np.array([[1], [2], [3], [4], [5]])
y = np.array([2, 4, 6, 8, 10])

# Step 1: Add a bias (intercept) term -> add a column of ones to X
m = X.shape[0]
X_b = np.c_[np.ones((m, 1)), X] # Shape: (m x n+1)

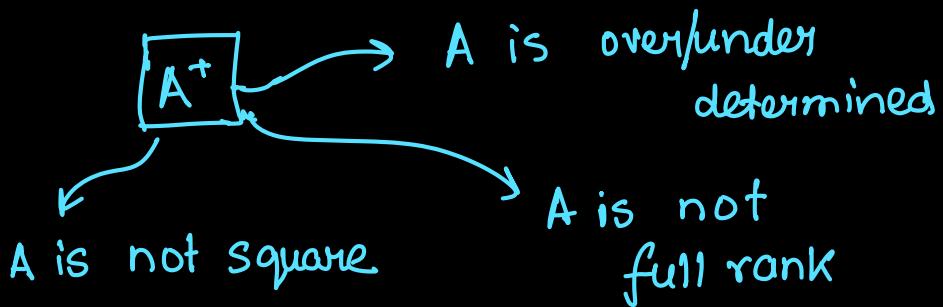
# Step 2: Apply the Normal Equation
#  $\theta = (X^T X)^{-1} X^T y$ 
theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)

print("Parameters (θ):", theta_best)

# Step 3: Predict for new data
X_new = np.array([[0], [6]])
X_new_b = np.c_[np.ones((X_new.shape[0], 1)), X_new] # add bias
y_predict = X_new_b.dot(theta_best)

print("Predictions:", y_predict)
```

Pseudo Inverse or Moore-Penrose Inverse:



$$\begin{aligned} \theta &= (X^\top X)^{-1} X \cdot y \\ &= X^+ y \end{aligned}$$

How computed ??

Singular Value Decomposition (SVD):

$$A = U \Sigma V^\top \Rightarrow A^+ = V \Sigma^+ U^\top$$

• Gives best approximate solution

$$X = U \Sigma V^T$$

↓ ↓
 orthogonal matrix diagonal matrix
 of singular values

(left singular
 vector)

$$X^+ = V \Sigma^+ U^T$$

To compute Σ^+ :

1. Take Σ - the diagonal matrix of singular values
2. Set to zero any singular value smaller than a small threshold (to avoid dividing by near zero)
3. Take the reciprocal ($1/\sigma_i$) for each non zero value
4. Transpose the result

Complexity in Computation:

Normal Equation computes

$X^\top X$ which is a

$(n+1) \times (n+1)$ matrix

$O(n^{2.4})$ to $O(n^3)$

If no. of features double :

5.3 to 8 times more time

But , SVD approach takes $O(n^2)$

\therefore Normal or SVD approach is long enough
to train but fast enough to predict...

For number of instances m :

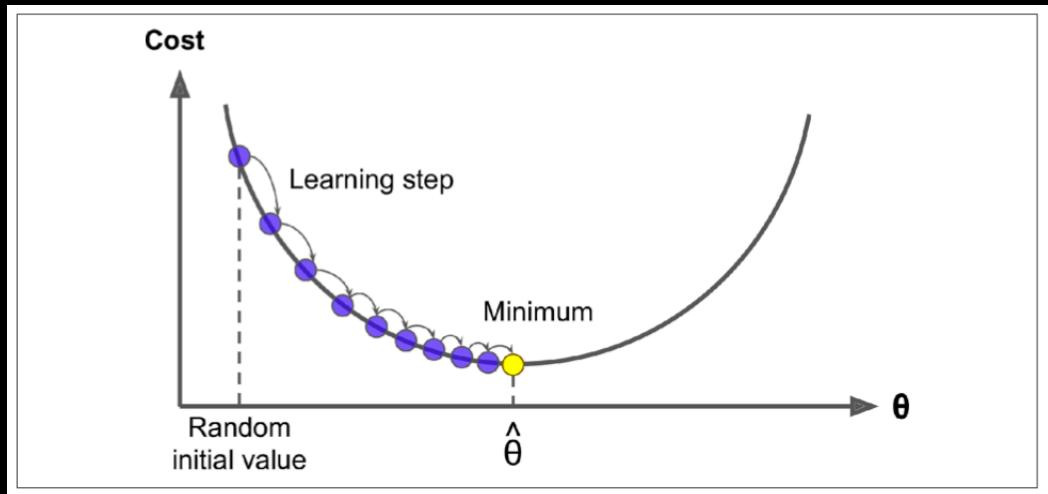
$O(m) \rightarrow$ linear

If (instances $\times 2$) : (prediction time $\times 2$)

GRADIENT DESCENT

generic optimization algorithm

↳ tweak parameters iteratively
in order to minimize c_f



Start by feeding Θ with random values
called random initialization

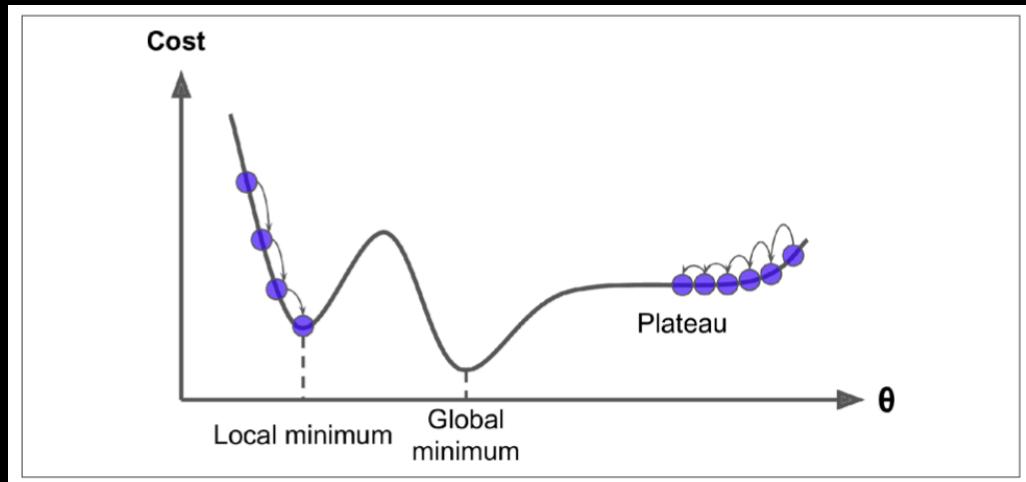
↓
Take one step at a time each attempted to
decrease the cost function

↓
Until the algorithm converges to
minimum

Size of step ↘ determined by
learning rate hyperparameter
not too small and not too large

Problem with GD:

- Local minima
- Plateau terrain
- Cut off



But fortunately the MSE cost function is a convex func:

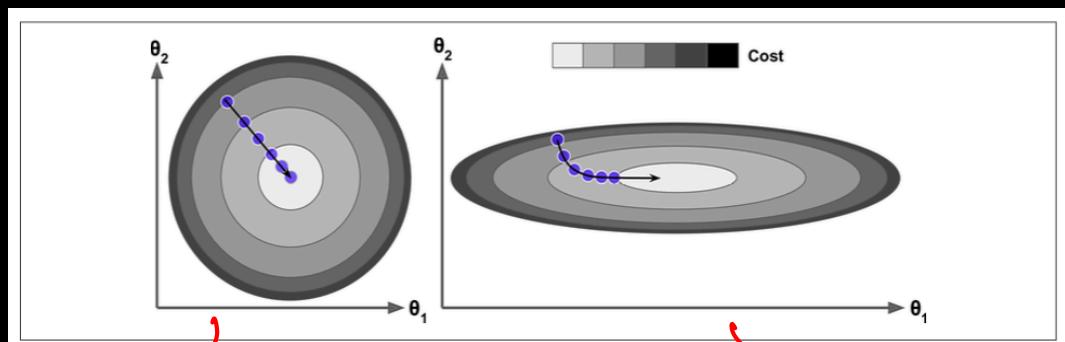
- It will only have one minima and it will be global minima
- It will be smooth

(Convexity + Smoothness)

The cost function has a shape of bowl

can be elongated bowl if the features have different scale

* Before using GD we should ensure all features have same scale



will take smaller time as it is scaled

will take longer time as features are not scaled

More parameters the mode have

more dimensions this space has

more hard to search

BATCH GD

$$MSE(\theta) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} MSE(\theta) = \frac{2}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

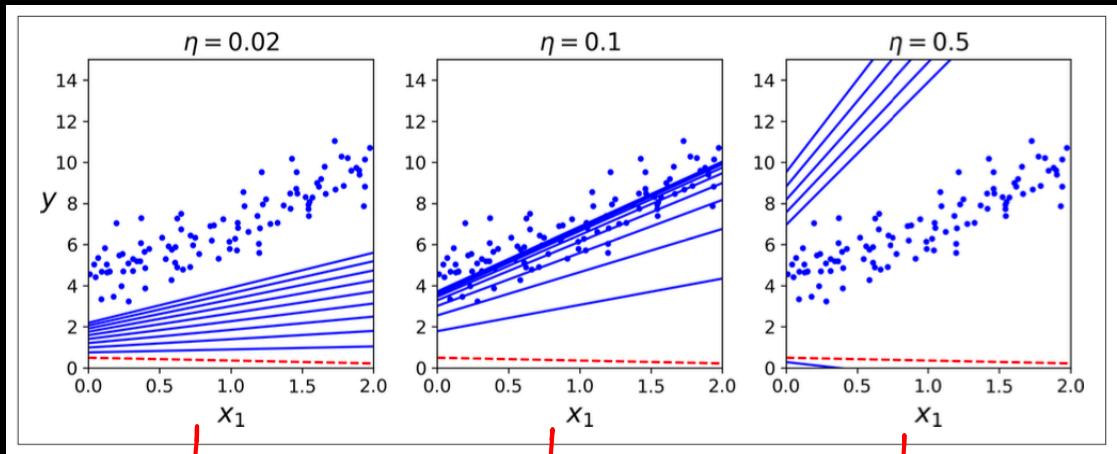
gradient vector of the cost fun^c:

$$\nabla_{\theta} MSE(\theta) = \left[\begin{array}{c} \frac{\partial}{\partial \theta_0} MSE(\theta) \\ \frac{\partial}{\partial \theta_1} MSE(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} MSE(\theta) \end{array} \right] = \frac{2}{m} X^T (X\theta - y)$$

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} \text{MSE}(\theta)$$

— gradient decent Step

η = learning rate



too slow

fit

too fast

How to set number of iteration:

A simple solution is set the no. of iteration to large number.

How to know when to stop:

when norm of gradient vector

becomes smaller than a tiny number ϵ (the tolerance)

$$\|\nabla_{\theta} \text{MSE}(\theta)\| < \epsilon (1 \times 10^{-6})$$

Convergence Rate:

For a batch GD the convergence rate is

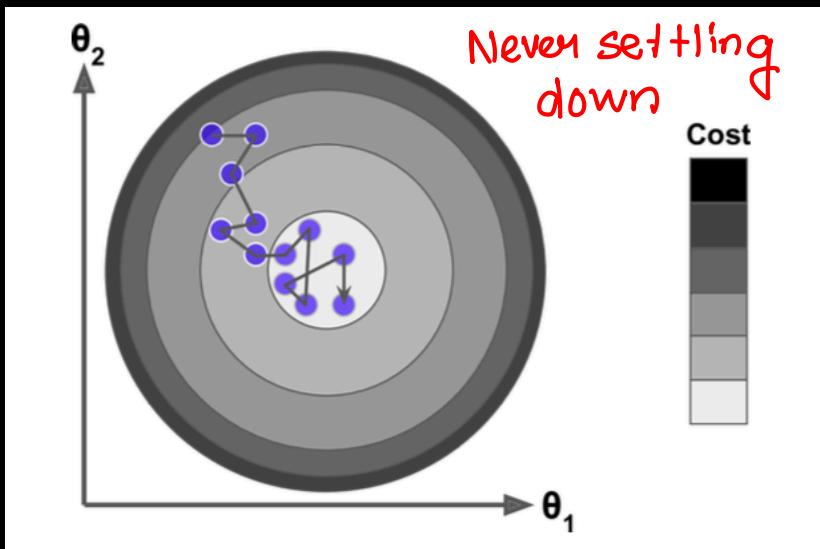
$$O\left(\frac{1}{\epsilon}\right)$$

To get within optimal solution, it may

take $\frac{1}{\epsilon}$ iteration

STOCHASTIC GD

- The batch gd uses whole training set and computes gradient at each step
- But stochastic gd uses random instances
- Hence it is faster to compute
- The algorithm is less regular than batch gradient
 - ↓
 - Until gently reaching minima it will bounce up and down
- Once the algorithm stops, the final parameter values are good but not optimal.



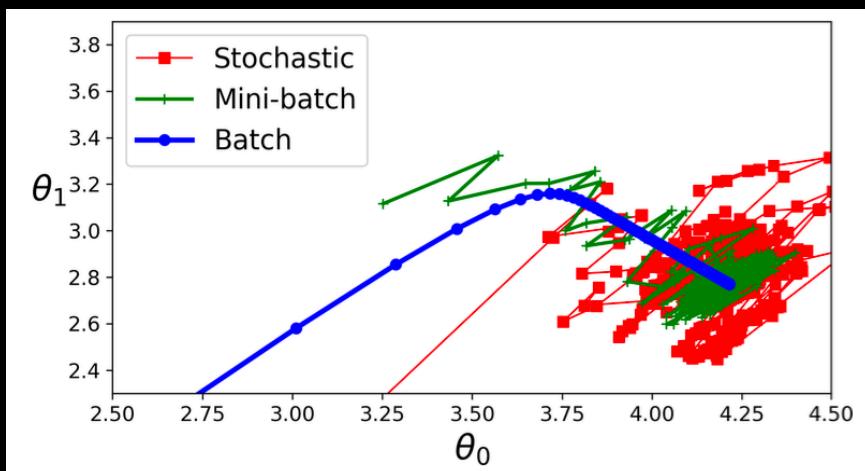
- This randomness less to jump out of local minimas but it means the algorithm can never settle down.
- One solution is reducing the learning rate.
 - The steps start out large but then gets smaller and smaller
 - This process is known as simulated annealing

- The function that determines the learning rate at each iteration is called the learning schedule.
- While using SGD the training instances should be distributed identically and independently (IID).
- If the instances are sorted by label, the SGD will optimize for one label and then for another

- SGD depends heavily on data shuffling to make unbiased parameter updates.
- Without shuffling, batch gradients become highly biased, especially early in training.
- This is especially dangerous for **imbalanced datasets** or **online learning**.

MINI BATCH GD

- Mini Batch GD computes gradients on small random sets of instances called mini-batches.
- One advantage of MBGD is that you get a performance boost over SGD.
- It may be harder for it to escape local minima



⇒ Batch one take more time but stops once reaching the global minima

⇒ Mini Batch and stochastic continue to walk around

Table 4-1. Comparison of algorithms for Linear Regression

Algorithm	Large m	Out-of-core support	Large n	Hyperparams	Scaling required	Scikit-Learn
Normal Equation	Fast	No	Slow	0	No	N/A
SVD	Fast	No	Slow	0	No	LinearRegression
Batch GD	Slow	No	Fast	2	Yes	SGDRegressor
Stochastic GD	Fast	Yes	Fast	≥ 2	Yes	SGDRegressor
Mini-batch GD	Fast	Yes	Fast	≥ 2	Yes	SGDRegressor

POLYNOMIAL REGRESSION

Polynomial Features all number of features upto a given degree.

Ex:- If there are two features:

a and b

for Polynomial degree = 3

a^3, b^3, b^2, a^2

and also, ab, a^2b, ab^2

For Polynomial Feature (degree = d)
and n feature:

$\frac{(n+d)!}{d! n!}$ feature array.

It basically creates $\frac{(n+d)!}{d! n!}$ features

and do linear regression between them.

LEARNING CURVES

How to know which curve to apply?

•) If training set error low } underfitting
 validation set error high }

If both high } overfitting

•) Another way is to look at learning curves:

Train the model several times on
different sized subsets of the training set

* If model underfitting: No use of adding
more training examples

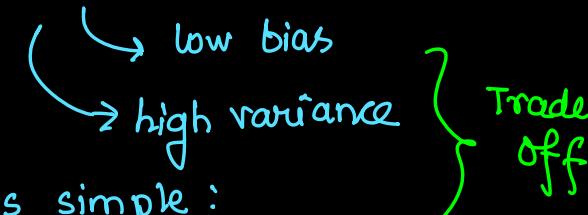
If model overfitting: Feed it more training
data until Val error reaches train error

BIAS-VARIANCE

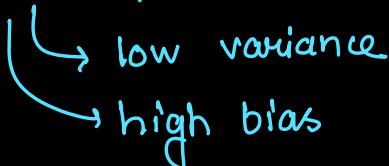
Model's generalization Error:

- $$\text{Bias} + \text{Variance} + \text{Irreducible Error}$$
- due to wrong assumptions
 - (Ex: Assume data is linear but actually quadratic)
 - A high bias model is most likely to overfit
 - due to excessive sensitivity to small variations in training data
 - A model with high degree of freedom is likely to have high variance leading to overfitting
- due to noisiness of data
→ only way to reduce this is to cleanup the data

* If model is complex :



If model is simple :



RIDGE REGRESSION

- A good way to reduce overfitting is to regularize the model.

 ↳ constrain it
 lower degree of freedom

- Tikhonov Regularization is the regularized version of Linear Reg

- A regularized term: $\alpha \sum_{i=1}^n \theta_i^2$
is added to CF

 ↳ only added during training

 ↳ once model is trained

 ↳ use unregularized parameter

- * CF used during training will not always be used for testing:

Reason: good training cost function should have optimization friendly derivative but the testing performance measure should be as close as possible to final objective

Ridge Regression Cost Function

$$J(\theta) = \text{MSE}(\theta) + \alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2$$



if $\alpha=0$,
Ridge Reg \approx Lin Reg

$\alpha = \text{very large}$

Model goes through avg. (simpler model)

this α hyperparameter
controls how we want to regularize

If we define w as a vector
of feature weights (θ_1 to θ_n):

then the regularized term is to

$$\frac{1}{2}(\|w\|_2)^2$$

for gradient decent add αw

* Imp to Scale before apply Ridge Reg

α Value	Effect on Model	Bias	Variance
Very Low	Overfit	Low	High
Medium	Balance fit	-	-
Very High	Underfit	High	Low

Ridge Regression Closed Form Solⁿ:

$$\hat{\theta} = (X^T X + \alpha I)^{-1} X^T y$$

With Ridge: the model is forced to keep coefficients small, leading to smaller curve

without it high degree polynomial
wiggles a lot (overfit)

- * Instead of shrinking coefficients the Ridge Reg can also constrain weights
And solve Constrained Optimization Problem

Minimize $\sum_{i=1}^m (y^i - \hat{y}^i)^2$ subject to $\sum_{j=1}^n \theta_j^2 \leq t$

constrain

* Bayesian Interpretation of Ridge:

$$P(\theta) \propto \exp\left(-\frac{\lambda}{2} \|\theta\|_2^2\right)$$

LASSO REGRESSION

↳ Least Absolute Shrinkage and Selection Operator Regression

→ uses ℓ_1 norm just like Ridge uses ℓ_2

Lasso Reg. Cost Function:

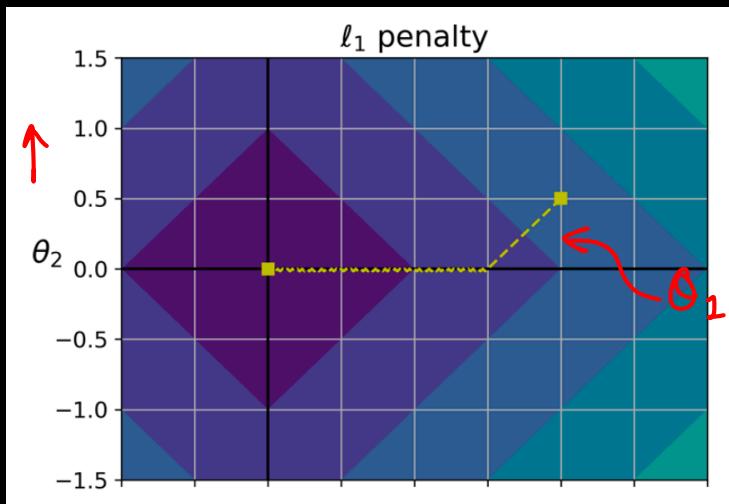
$$J(\theta) = \text{MSE} \theta + \alpha \sum_{j=1}^n |\theta_j|$$

→ Lasso Reg tends to eliminate the weights of least important features (set \rightarrow zero)

Or one can say Lasso Reg automatically performs feature selection and outputs a sparse Model

Lasso Vs Ridge Regularization:

Ex: for $\theta_1 = 2$ and $\theta_2 = 0.5$

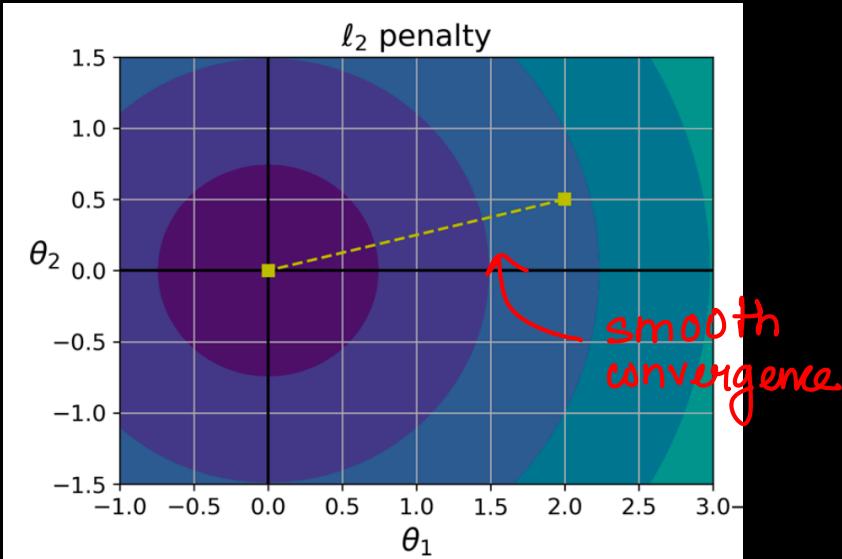


Loss func: $\ell_1 = |\theta_1| + |\theta_2|$

Diamond shaped counter
(Manhattan norm)

$$\text{gradient} = \frac{d}{d\theta_j} |\theta_j| = \begin{cases} 1 & \theta_j > 0 \\ -1 & \theta_j < 0 \\ \text{undef} & \theta_j = 0 \end{cases}$$

\therefore gradient decent moves in straight path based on gradient direction



Loss fun² :

$$\ell_2 = \theta_1^2 + \theta_2^2$$

circular contours

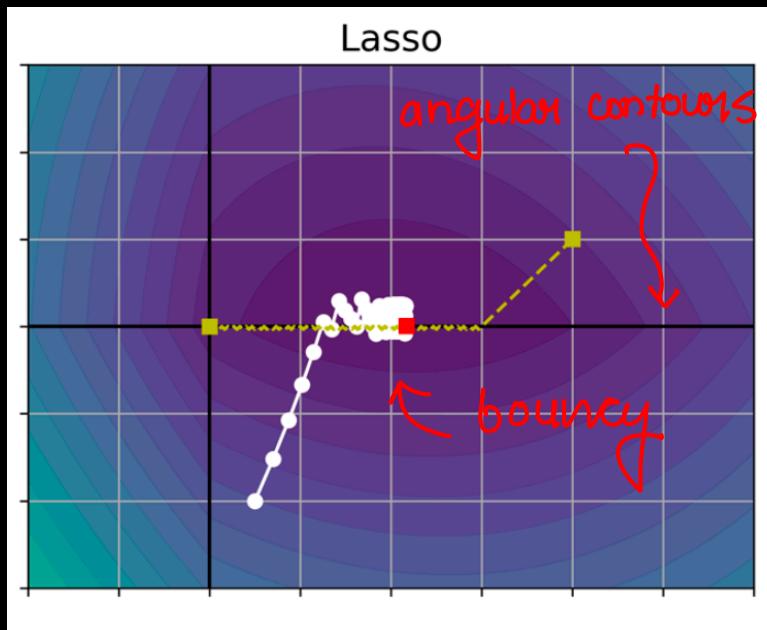
Euclidean Norm

$$\text{gradient} = \frac{d}{d\theta_j} \ell_2 = 2\theta_j$$

Smoothly decays to zero as
 $\theta_j \rightarrow 0$

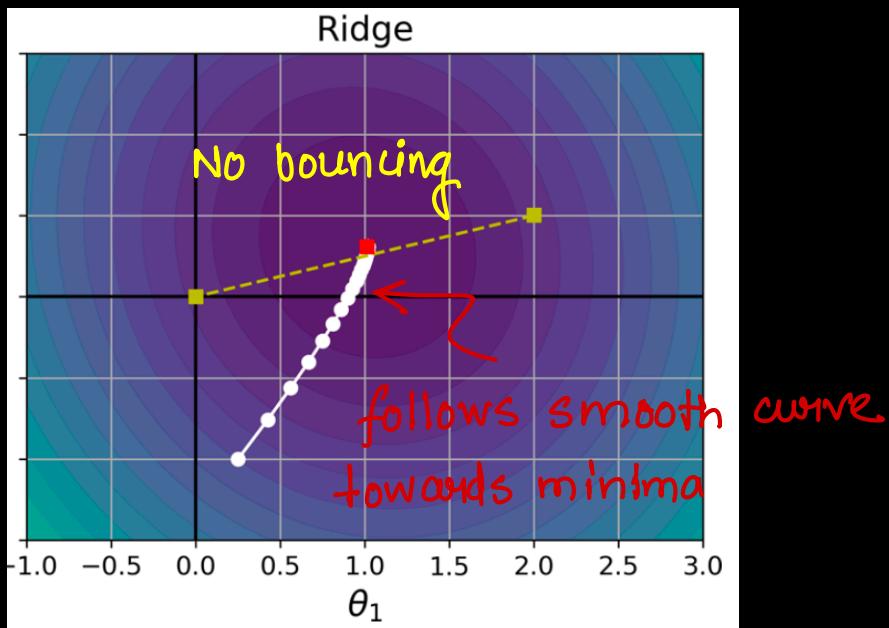
Feature	L1 (Lasso)	L2 (Ridge)
Derivative	Constant ± 1 or subgradient at 0	Smooth: $2\theta_j$
Behavior near zero	Sharp, aggressive pull toward 0	Gentle, smooth shrinkage
Effect	Sets some weights exactly to zero	Keeps all weights, shrinks all

•) Cost fun^c with Lasso Reg fun^c
 ↳ (actually, Smooth bowl)
 (MSE + ℓ_1)



•) Analogy ball rolls into a corner
 then along edges
 Constant bounces at corners

•) Cost func with Ridge Reg func:
 $(MSE + l_2)$



•) Smooth convergence

Property	Lasso (L1)	Ridge (L2)
Penalty Shape	Diamond (Manhattan norm)	Circle (Euclidean norm)
Encourages Sparsity?	<input checked="" type="checkbox"/> Yes, zeroes out weights	<input type="checkbox"/> No, shrinks but keeps all
Feature Selection?	<input checked="" type="checkbox"/> Yes	<input type="checkbox"/> No
Gradient Behavior	Constant, bounces at corners	Smooth, decays near optimum
Path to Optimum	Axes-aligned, piecewise	Smooth, radial inward
Good for High-Dim?	<input checked="" type="checkbox"/> When only few features matter	<input checked="" type="checkbox"/> When all features contribute

Lasso-Regression Subgradient Vector

$$g(\theta) = \nabla_{\theta} \text{MSE}(\theta) + \alpha \begin{pmatrix} \text{sgn}(\theta_1) \\ \text{sgn}(\theta_2) \\ \vdots \\ \text{sgn}(\theta_n) \end{pmatrix}$$

$$\text{sgn}(\theta_i) = \begin{cases} -1 & \theta_i < 0 \\ 0 & \theta_i = 0 \\ +1 & \theta_i > 0 \end{cases}$$

ELASTIC NET

•) Middle ground btw Ridge and Lasso:

•) We can control the mix ratio:

when $\gamma=0 \Rightarrow$ Ridge Regression

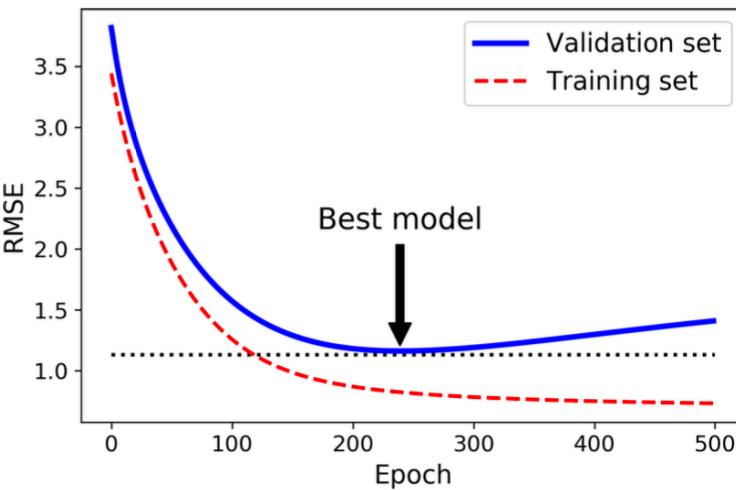
$\gamma=1 \Rightarrow$ Lasso Regression

Elastic Net Cost Function:

$$J(\theta) = \text{MSE}(\theta) + \gamma \alpha \sum_{i=1}^n |\theta_i| + \frac{1-\gamma}{2} \alpha \sum_{i=1}^n \theta_i^2$$

Lasso gives sparsity
Ridge gives stability } Net give both

* Early Stopping



•) RMSE decreases for both training and validation



But after sometime the model overfits



Hence it is imp to stop when the model hit minimum RMSE

•) `warm_start = True`

allows incremental training across epochs

`Clone()`

Best model stored using `Clone`

LOGISTIC REGRESSION

- It is used to estimate the probability that an instance belongs to a particular class.
- If prob (>50%) it belongs to that class (+ve class) or it does not (-ve class)

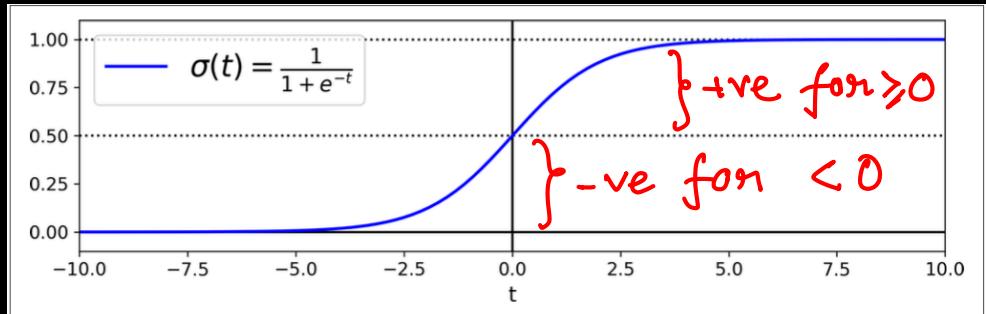
Logistic Reg Model (Vectorized Form)
Estimated Prob:

$$\hat{p} = h_{\theta}(x) = \sigma(x^T \theta)$$

$\sigma(\cdot)$ = sigmoid func used to get value btw 0 to 1

Logistic Func'

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$



• Once the Log Reg estimated the prob

$$\hat{p} = h_{\theta}(x)$$

$x \in$ positive class

It predicts \hat{y} easily

Log Reg Model Prediction:

$$\hat{y} = \begin{cases} 0 & , \hat{p} < 0.5 \\ 1 & , \hat{p} \geq 0.5 \end{cases}$$

So,

$$\tau(\underline{x^T \theta})$$

value < 0

$\tau(<0)$

= < 0.50

= -ve class

value ≥ 0

$\tau(\geq 0)$

= ≤ 0.50

= +ve class

Score + is also known as logit

$$\text{logit}(P) = \log\left(\frac{P}{1-P}\right)$$

Also called log odds,

since log of ratio between
+ve class and -ve class

$$(t = X^T \theta)$$

$$p = \sigma(t) = \frac{1}{1+e^{-t}}$$

$$\Rightarrow \frac{p}{1-p} = \frac{1+e^{-t}}{1-e^{-t}}$$

$$\Rightarrow = \frac{1+e^{-t}}{e^{-t}}$$

$$\Rightarrow = e^t$$

$$\Rightarrow \log\left(\frac{p}{1-p}\right) = \log(e^t)$$

$$\Rightarrow \text{logit}(p) = t$$

Cost Func of single training instance:

$$C(\theta) = \begin{cases} -\log(\hat{p}) & \text{if } y=1 \\ -\log(1-\hat{p}) & \text{if } y=0 \end{cases}$$

Logistic Regression Cost Func:
(log loss)

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{p}^{(i)}) + (1-y^{(i)}) \log(1-\hat{p}^{(i)})]$$

Log Cost Func partial derivative:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (\sigma(\theta^\top x^{(i)}) - y^{(i)}) x_j^{(i)}$$

* Decision Boundary:

- A surface that separates different classes in a classification problem
- It represents a threshold
- A line in 2D
A plane in 3D
and hyperplane in higher dimension

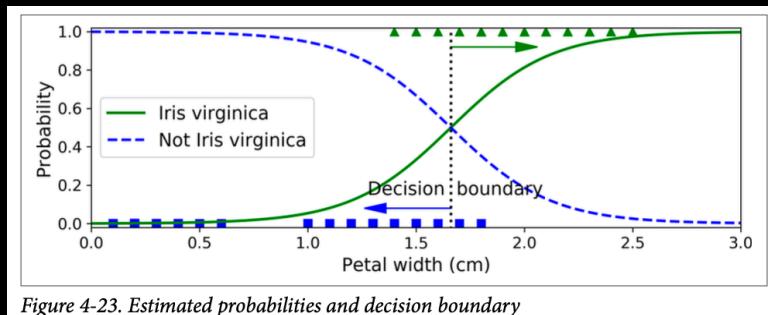
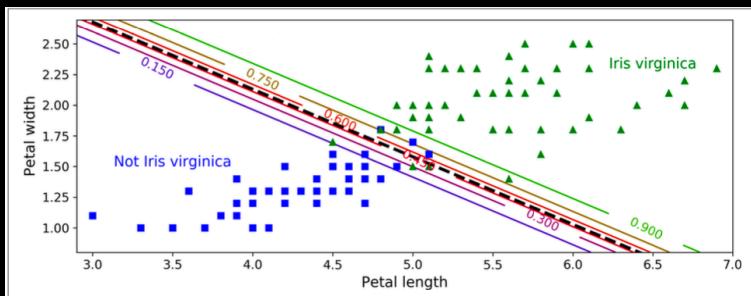


Figure 4-23. Estimated probabilities and decision boundary



[Linear Decision]
Boundary]

SOFTMAX REGRESSION

- also known as multiNominal Logistic Regression

SoftMax Score for class k:

$$S_k(x) = x^T \theta^k$$

- Each class has its own dedicated parameter

vector $\theta^{(k)}$

All these vectors are stored as rows
in a parameter matrix

• When given instance x

↓

SoftMax Regression first calculates score $s_k(x)$

for each class k

↓

then estimates the prob of each class by apply the softmax fun^c
(also called Normalized Exp)

↓

Once computed probability classifier predicts the class with highest estimated probability

SoftMax function:

$$\hat{p}_x = \sigma(s(x))_k = \frac{\exp(s_k(x))}{\sum_{j=1}^k \exp(s_j(x))}$$

k is number of class

$s(x)$ vector containing the scores

$\sigma(s(x))$ estimated prob that the instance x belongs to class k given the score of each instance

SoftMax Reg classifier Prediction:

$$\hat{y} = \operatorname{argmax}_k \sigma(s(x))_k$$

$$= \operatorname{argmax}_k s_k(x)$$

$$= \operatorname{argmax} ((\theta^{(k)})^T x)$$

-) argmax operator returns the value of variable for which the function maximizes
 -) In this operator the value of k for which the estimated prob $\sigma(s(x))_k$ maximize
- * SoftMax should only be used with mutually exclusive classes
It is not multilabel or multiorput

ENTROPY

- Measures the uncertainty in a prob dist.

Entropy for discrete dis P:

$$H(P) = - \sum_x P(x) \log P(x)$$

High Entropy = more uncertainty
(uniform dist)

Low Entropy = less uncertainty
(one class has prob 1)

Cross-Entropy:

$$H(P, Q) = - \sum P(x) \log Q(x)$$

P ≈ True distribution

Q ≈ Predicted Distribution

- Used as loss function in Classification

KL Divergence:

$$D_{KL}(P || Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

= Cross Entropy - Entropy

- Not symmetric

Cost Entropy Cost Function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log (\hat{p}_k^{(i)})$$

For ($K=2$),

The equation becomes log loss

.) goal is to minimize $J(\theta)$,
cost function



typically done by gradient descent

.) Used in  SoftMax Classifier
Binary Classifier