

<u>Supervised learning</u>	<u>Unsupervised learning</u>
<ul style="list-style-type: none"> (i) Algorithms are trained using labeled data. (ii) More accuracy. (iii) Number of classes is known. (iv) Uses offline analysis. 	<ul style="list-style-type: none"> (i) unlabeled data. (ii) Less accuracy (iii) not known. (iv) Uses real time analysis of data.
<u>Classification</u>	<u>Regression</u>
<ul style="list-style-type: none"> (i) Output variable must be a discrete value. (ii) Classification tries to find the decision boundary which can divide the dataset into different classes. (iii) Types: Binary classifier, Multi-class classifier 	<ul style="list-style-type: none"> (i) Output variable must be of continuous nature or real value. (ii) Regression tries to find the best fit line, which can predict the output. (iii) Types: Linear regression, Non-linear regression.
<u>Classification Algorithms :</u>	<u>Regression Algorithms :</u>
<ul style="list-style-type: none"> * Logistic regression (not a regression) * K-nearest neighbors (KNN) * Support Vector Machine (SVM) <ul style="list-style-type: none"> • Naive Bayes classification * Decision tree classification * Random forest classification 	<ul style="list-style-type: none"> • Simple linear regression • Multiple linear regression • Polynomial regression • Support vector regression • Decision tree regression • Random forest regression.

<u>Overfitting & Underfitting :</u> [Prerequisites: Variance, Bias]
<p>* <u>Bias</u>: refers to how correct the model is.</p> <ul style="list-style-type: none"> • Simple model with lot of mistakes \Rightarrow High bias. • Complicated model that does well on training data \Rightarrow low Bias.
<p>* <u>Variance</u>: describes how much a prediction could potentially vary if one of the predictors changes slightly.</p> <ul style="list-style-type: none"> • Simple model changes prediction slowly with predictor value \Rightarrow low variance • Complicated model that changes prediction widely, even predictor value changes slightly : High variance • Low Bias

too complex relative to amount of noise

- ① Oversfitting:
 - When a model is trained with so much data, it starts learning from noise and results in high variance.
 - Oversimplified model does not categorise the data correctly.
 - " " have high variance and low bias.
 - Model is too complex and unrealistic, low generalization capacity.

② Techniques to reduce overfitting:

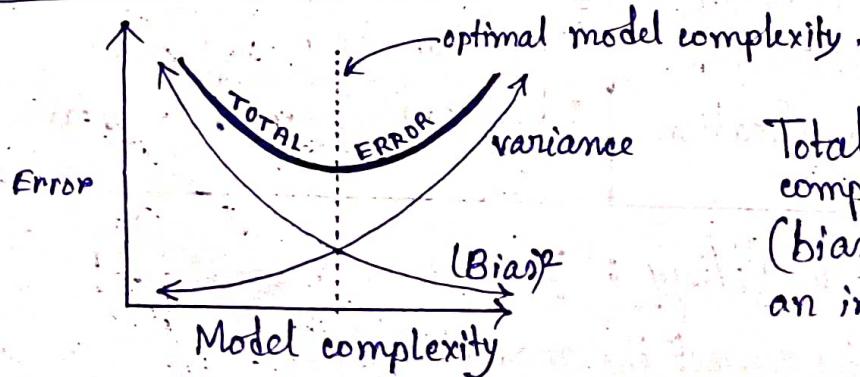
- 1) Increase training data
- 2) Reduce model complexity
- 3) Regularization : Ridge & Lasso
- 4) Use of dropout for neural networks.
- 5) Using cross validation.

- ③ Underfitting: If a statistical model can not capture underlying trend of the training data due to less data and more features to build the model, hence have low generalization capacity, high bias, low variance

④ Techniques to reduce underfitting:

- 1) Increase model complexity
- 2) Remove noise from data
- 3) Increase number of features performing feature engineering.
- 4) Increase number of epochs (number of passes of training dataset)
- 5) Increase duration of training.

* Bias-variance tradeoff:



Total error of a model is composed of three terms: $(\text{bias})^2$, variance, an irreducible error term.

- We modify the cost function to restrict the weights, which allows to trade excessive variance, potentially reducing overall error.

This trade comes in the form of regularization.

Linear Regression

A linear model that makes prediction by simply computing a weighted sum of the input features + a constant bias term.

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

\hat{y} = the predicted value, x_i = i^{th} feature value

w_0 = bias term

w_1, w_2, \dots, w_n = feature weights.

- ① Simple linear regression : $\hat{y} = w_0 + w_1 x$.
prediction is made by only one predictor variable 'x'.

- ② Multiple linear regression, more than one predictor variable is present. (x_1, x_2, \dots, x_n)

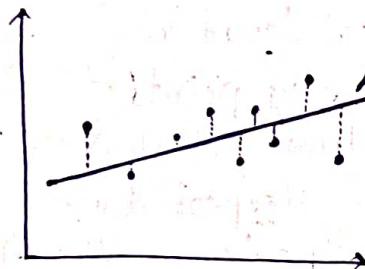
Vectorized form : $\hat{y} = h_w(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$

\mathbf{w} = model's parameter vector

\mathbf{x} = instance's feature vector.

* \mathbf{w}, \mathbf{x} are column vectors or column matrices.

$$\mathbf{w}^\top \mathbf{x} = [w_0 x_0 + w_1 x_1 + \dots + w_n x_n], x_0 = 1, w_0 = \text{bias term.}$$



(The unique line such that the sum of the squared vertical distances between the data points and the line is the smallest possible.)

→ Our goal is to find the line, i.e., find the weights w_1, w_2, \dots, w_n and bias term w_0 .

Average Loss function:

$$\text{MSE}(\mathbf{X}, h_w) = \frac{1}{m} \sum_{i=1}^m (w_1^\top \mathbf{x}^{(i)} - y^{(i)})^2$$

MSE = mean square error.

The squared vertical distance

→ Target is to minimize the loss function.

- ① for simple linear regression, $\hat{y} = w_0 + w_1 x$,

we have the closed form solution : (taking partial derivatives of loss function wrt w_0 & w_1)

$$w_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

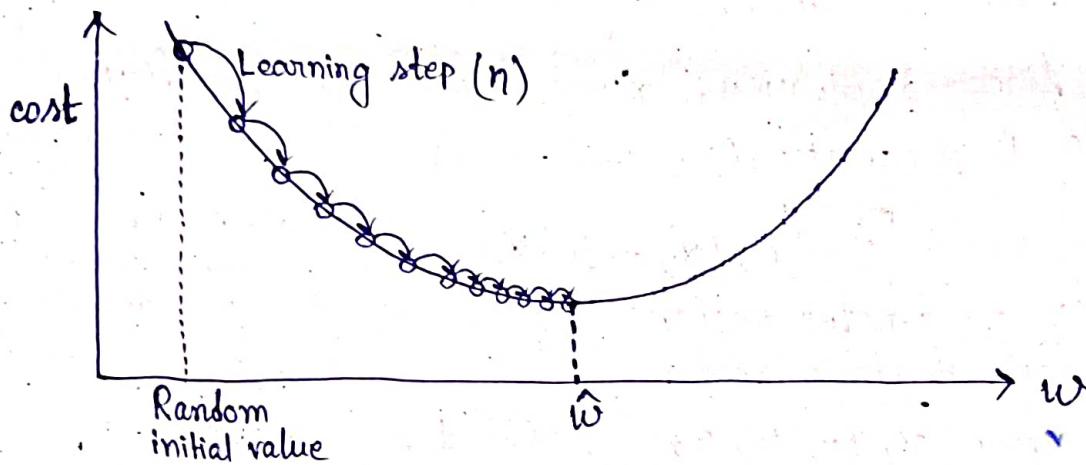
$$w_0 = \frac{\sum y - w_1 \sum x}{n}$$

① LMS (least minimum slope) Algorithm (delta rule)

LMS algo tunes w_0 & w_1 to produce the model that minimizes the avg loss function.

- We start with a initial guess of w_0, w_1 .
- then update the values of w_0 & w_1 .
- We use gradient descent method to update the values.

Gradient descent



The loss function is a convex function (quadratic) that has a single minima. We start at any point of this function space and update the values of w_0 & w_1 , as : we take a step in the direction of steepest decrease of Loss function, i.e, towards the -ve slope (or gradient).

The learning step size (η) is proportional to the slope of the cost function so that the steps gradually get smaller as the parameters approach the minimum.

Batch gradient descent: $MSE(x, h_w)$ is represented as $MSE(w)$

(i) Gradient vector of the cost function:

$$\nabla_w (MSE(w)) = \begin{bmatrix} \frac{\partial}{\partial w_0} MSE(w) \\ \frac{\partial}{\partial w_1} MSE(w) \\ \vdots \\ \frac{\partial}{\partial w_n} MSE(w) \end{bmatrix} = \frac{2}{m} X^T (XW - y)$$

* x, w are all matrices.

x = matrix of all features.

* η = learning step.

(ii) Gradient descent step: $w^{(next\ step)} = w - \eta \nabla_w MSE(w)$

** The calculation involves the full training dataset at each gradient descent step. So, it is terribly slow on very large training sets.

④ Calculation of partial derivative for a single feature

$$\begin{aligned}
 \frac{\partial \text{MSE}(w)}{\partial w_n} &= \frac{\partial}{\partial w_n} \left(\frac{1}{m} \sum_{i=1}^m (w^\top x^{(i)} - y^{(i)})^2 \right) \\
 &= \frac{\partial}{\partial w_n} \left(\sum_{i=1}^m (w^\top x^{(i)} - y^{(i)}) \right) \cdot \left\{ \frac{2}{m} \sum_{i=1}^m (w^\top x^{(i)} - y^{(i)}) \right\} \\
 &= x_n^{(i)} \cdot \left\{ \frac{2}{m} \sum_{i=1}^m (w^\top x^{(i)} - y^{(i)}) \right\} \\
 &= \boxed{\frac{2}{m} \sum_{i=1}^m (w^\top x^{(i)} - y^{(i)}) \cdot x_n^{(i)}}
 \end{aligned}$$

$$w_n^{(\text{next step})} = w_n - \eta \cdot \frac{\partial \text{MSE}(w)}{\partial w_n}$$

④ Notations

$x^{(i)}$ = vector of all feature values of i^{th} instance.

X = a matrix containing all the feature values of all instances.

$x_n^{(i)}$ = n^{th} feature value of i^{th} instance.

Example:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}, \quad X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_{p-1}^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_{p-1}^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_{p-1}^{(m)} \end{bmatrix}_{m \times p}, \quad \hat{W} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{p-1} \end{bmatrix}_{p \times 1}$$

$$\nabla_w (\text{MSE}(w)) = \frac{2}{m} \cdot X^\top (Xw - Y)$$

$$* X = \begin{bmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix}$$

* The Normal equation (Closed form soln.)

$$\hat{W} = (X^\top X)^{-1} X^\top Y$$

\hat{W} = the value of w that minimizes cost function.

y = vector of target values.

X = matrix containing all feature values of all instances.

⑩ Stochastic Gradient Descent :

Stochastic Gradient descent picks a random instance in the training set at every step and computes the gradients based on that single instance, that makes the algorithm makes work much faster and possible to train on huge training sets.

Instead of gently decreasing the cost function until it reaches the minima, the cost function will bounce up and down; decreasing only on average. Over time, it will end up very close to the minimum.

After reaching very close to minima, it still continues to bounce around and never settles down.

When the cost function is very irregular, this can help the algo to jump out of local minima and find the global minima.

⑪ Mini batch gradient descent :

At each step, instead of computing the gradients based on full training set, or based on just one instance, it computes the gradients on small random sets of instances called mini-batches.

* All the three gradient descents end up near the minimum, but batch gradient descent's path actually stops at the minimum, while both stochastic and mini-batch gradient descent continue to walk around the minima.

① Polynomial Regression:

A straight line will never fit this data properly.

Add higher degrees to the feature:

* If there are more than one feature (n features),

and degree ' n ', total features (including higher & order terms and combinations of them) = $(n+d)! / (n! \cdot d!)$

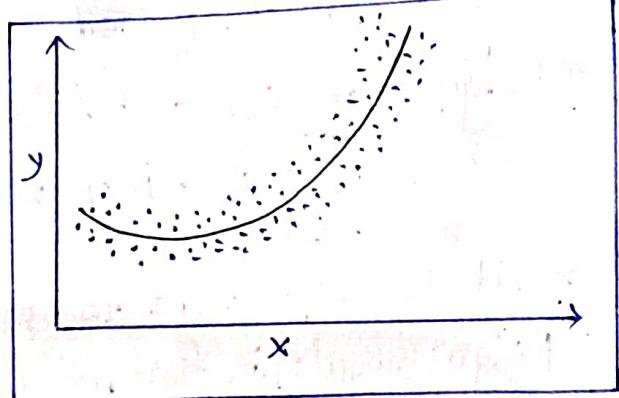
for one feature and degree = 2, $\hat{y} = w_0 + w_1 x + w_2 x^2$

for two features (x_1 & x_2) and degree = 2,

* Not only the higher order of the features are added, but the combinations are also added.

$$\hat{y} = w_0 + w_1 x_1^2 + w_2 x_1 + w_3 x_2^2 + w_4 x_2 + w_5 x_1 x_2$$

$$(2+2)! / (2! \cdot 2!) = 6 \text{ (including power)}$$



② Regularization:

(1) L1 regularization or Lasso regularization:

→ adds a penalty to the loss function.

Lasso regularization cost function:

$$J(w) = \text{MSE}(w) + \alpha \sum_{i=1}^n |w_i|$$

$\sum_{i=1}^n |w_i|$ = sum of absolute values of weights.

α = regularization parameter.

* Lasso regularization tends to eliminate the weights of the least important features. So, it automatically performs the feature selection and outputs a sparse model.

* To avoid gradient descent from bouncing around the optimum at the end while using Lasso, the learning step is needed to be reduced gradually so that the steps get smaller and smaller near the optimum and so it will converge.

(2) L2 regularization or Ridge regularization:

→ Adds sum of squared values of weights as penalty.

① Ridge regression cost function:

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

| it tries to keep
the model weights
as small as
possible |

* It is important to scale the data before performing Ridge regularization, as it is sensitive to the scale of the input features.

Ridge regression closed form solution :

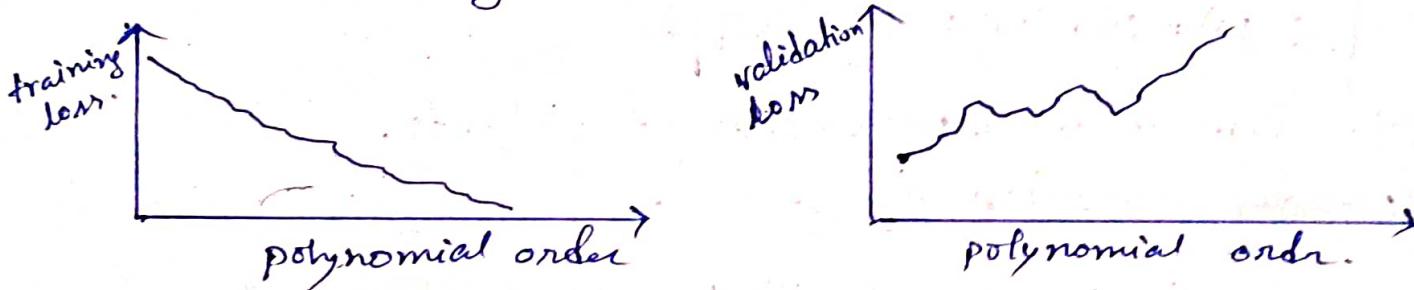
$$\hat{w} = (x^T x + \alpha I)^{-1} x^T y$$

I = identity matrix except with a '0' at top left, corresponding to bias-term.

* Elastic Net : it is a middle ground between ridge and Lasso regression.

$$J(w) = \text{MSE}(w) + \alpha \sum_{i=1}^n |w_i| + \frac{1-\rho}{2} \alpha \sum_{i=1}^n w_i^2$$

② Validation : Use of a second dataset to compute loss. Validation loss increases monotonically with the order of polynomial, i.e., model complexity.

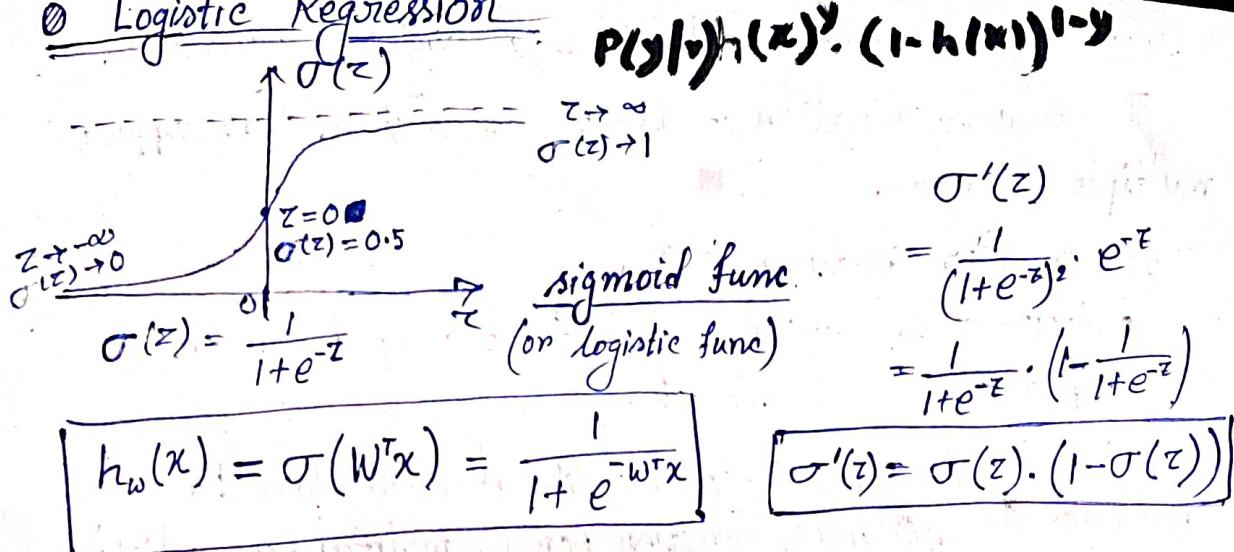


③ Cross validation : Every data pt gets to be in validation set exactly once and gets to be in training set ($K-1$) times.

→ This significantly reduces bias and variance as we are using most of the data for fitting and in validation set.

→ $K = 5$ to 10 is generally preferred.

① Logistic Regression



* Probability refers to the chance that a particular outcome occurs based on the values of parameters in a model.

* Likelihood refers to how well a sample provides support for particular values of a parameter in a model.

① Logistic regression model estimates the probability that an instance x belongs to the +ve class

① Cost function for a single training instance:

$$C(w) = \begin{cases} -\log(h_w(x)) & \text{if } y=1 \\ -\log(1-h_w(x)) & \text{if } y=0 \end{cases} \quad [\text{There are two classes}]$$

① Cost function over the whole training set is the avg cost over all training instances.

① Logistic regression cost function (log loss)

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

There is no known closed form solution to compute 'w'.

The cost function is convex. so, gradient descent method finds global minimum.

$$\frac{\partial}{\partial w_n} J(w) = \frac{1}{m} \sum_{i=1}^m (\sigma(w^T x^{(i)}) - y^{(i)}) x_n^{(i)}$$

- ② Softmax regression: (Multinomial logistic regression)
 This is generalization of logistic regression, that support multiple classes.

→ for a given instance (x):

- Softmax regression model first computes a score $s_k(x)$ for each class k .
- Then estimates the probability of each class by applying the softmax function (or normalized exponential) to the scores.
- Softmax score for class k :

$$s_k(x) = x^T w^{(k)}$$

- Each class has its own dedicated parameter vector $w^{(k)}$.
 → All these vectors are stored as rows in a parameter matrix 'W'.

$$W = \begin{bmatrix} (w^{(1)})^T \\ (w^{(2)})^T \\ \vdots \\ (w^{(n)})^T \end{bmatrix}$$

- ② The probability that the instance x belongs to class k : (Softmax function)

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

- k = number of classes
- $s(x)$ = vector containing the scores of each class for instance x
- $\sigma(s(x))_k$ = estimated probability that x belongs to class k , given score of each class for x .

→ Softmax regression predicts the class with the highest estimated probability.

- ③ Softmax regression classification prediction:

$$\hat{y} = \operatorname{argmax}_k \sigma(s(x))_k = \operatorname{argmax}_k s_k(x) = \operatorname{argmax}_k ((w^{(k)})^T x)$$

- 'argmax' operator returns the value of a variable that maximizes a func.

① Cross entropy cost function:

It penalizes the model when it estimates a low probability for a target class. It is used to measure how well a set of estimated class probabilities matches the target classes.

$$J(W) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(p_k^{(i)})$$

- $y_k^{(i)}$ is the target probability that the i^{th} instance belongs to class k . In general, it is either 0 or 1.
- for $K=2$, the cost function is equivalent to logistic regression's cost func.

② Cross entropy gradient vector for class k :

$$\nabla_{w^{(k)}} J(W) = \frac{1}{m} \sum_{i=1}^m (\hat{p}_k^{(i)} - y_k^{(i)}) x^{(i)}$$

Now, the gradient vector for every class can be computed & then we gradient descent to find the parameter matrix 'W', that minimizes the cost function.

Decision Tree

Gini impurity:

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$

$p_{i,k}$ = the ratio of class k instances among the training instances in the i^{th} node.

① CART cost function for classification:

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

CART = classification and regression tree.

Training set is split into two subsets using a single feature 'K' and threshold ' t_k ', m = number of instances.

CART algo searches for pair (k, t_k) that produces purest subsets.

Once the CART algo successfully searches for the sets, it splits the subsets using same logic recursively until it reaches the maximum depth, or it cannot find a split that will reduce impurity.

② Entropy: A sets entropy is 0 when it contains instances of only one class.

Entropy of i^{th} node: $H_i = - \sum_{k=1}^n p_{i,k} \cdot \log_2(p_{i,k})$, $p_{i,k} \neq 0$

③ Information Gain(feature) = Entropy(dataset) - Entropy(feature)

most of the times, they lead to similar trees.

Gini impurity tends to isolate the most frequent class in its own branch of the tree.

Entropy tends to produce slightly balanced tree.

④ CART cost function for regression:

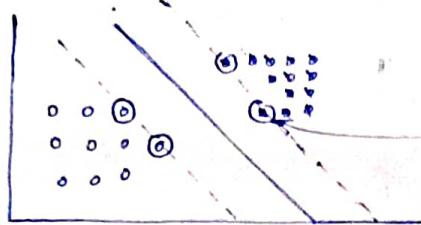
$$J(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}$$

$$\text{MSE}_{\text{node}} = \sum_{i \in \text{node}} (\hat{y}_{\text{node}} - y^{(i)})^2$$

$$\hat{y}_{\text{node}} = \frac{1}{m_{\text{node}}} \cdot \sum_{i \in \text{node}} y^{(i)}$$

Support Vector Machine

→ capable of linear, non-linear classification, regression, and outlier detection.



support vectors: instances located on the edge of the street.

① Hard and soft margin classification:

In hard margin classification, all instances must be off the ~~the~~ street to one side for a class. It is sensitive to outliers.

Soft margin classification finds a good balance between keeping the street as large as possible and limiting the margin violations.

② Kernel trick:

The kernel trick makes possible to get same result even if many polynomial features are added with very high degree polynomials. So, there is no combinatorial explosion of the number of features.

③ Similarity features:

Similarity function measures how much each instance resembles a particular landmark.

Gaussian Radial Basis function:

$$\phi_y(x, l) = \exp(-\gamma \|x - l\|^2)$$

l = landmark

x = value of instance.

decision func: $w^T x + b = w_1 x_1 + \dots + w_n x_n + b$

prediction $\hat{y} = 0$ if $w^T x + b < 0$

1 if $w^T x + b \geq 0$

① Hard margin SVM classifier objective:

$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} w^T w$$

$$\text{subject to } t^{(i)}(w^T x^{(i)} + b) \geq 1 \quad \text{for } 1, 2, \dots, m.$$

② Soft margin SVM classifier objective:

$$\underset{w,b,\zeta}{\text{minimize}} \quad \frac{1}{2} w^T w + C \sum_{i=1}^m \zeta^{(i)}$$

$$\text{subject to } t^{(i)}(w^T x^{(i)} + b) \geq 1 - \zeta^{(i)}$$

ζ = slack variable ≥ 0

$\zeta^{(i)}$ measures how much the i^{th} instance is allowed to violate the margin.

③ Kernelized SVM

→ Second degree polynomial mapping:

$$\phi(x) = \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

for two 2D vectors a, b

④ Kernel trick for 2nd degree polynomial mapping:

$$\phi(a)^T \phi(b) = (a^T b)^2$$

$$K(a, b) = (a^T b)^2 \quad \boxed{\text{is 2nd degree polynomial kernel}}$$

* Kernel is a function capable of computing the dot product $\phi(a)^T \phi(b)$ based only on the original vectors a, b without computing the transformation ϕ .

Common kernels:

Linear: $K(a, b) = a^T b$

Polynomial: $K(a, b) = (ra^T b + r)^d$

Gaussian RBF: $K(a, b) = \exp(-r||a-b||^2)$

Sigmoid: $K(a, b) = \tanh(ra^T b + r)$

① Making prediction with a kernelized SVM :

$$\begin{aligned}
 h_{\hat{\omega}, \hat{b}}(\phi(x^{(n)})) &= \hat{\omega}^\top \phi(x^{(n)}) + \hat{b} \\
 &= \left(\sum_{i=1}^m \hat{\alpha}^{(i)} t^{(i)} \phi(x^{(i)}) \right)^\top \phi(x^{(n)}) + \hat{b} \\
 \star \rightarrow &= \sum_{i=1}^m \hat{\alpha}^{(i)} t^{(i)} (\phi(x^{(i)})^\top \phi(x^{(n)})) + \hat{b} \\
 &= \sum_{i=1}^m \hat{\alpha}^{(i)} t^{(i)} k(x^{(i)}, x^{(n)}) + \hat{b} \\
 &\quad \hat{\alpha}^{(i)} > 0
 \end{aligned}$$

Using kernel trick to compute bias term \hat{b} :

$$\begin{aligned}
 \hat{b} &= \frac{1}{n_s} \sum_{i=1}^m \left(t^{(i)} - \hat{\omega}^\top \phi(x^{(i)}) \right) \\
 &\quad \alpha^{(i)} > 0 \\
 &= \frac{1}{n_s} \sum_{i=1}^m \left(t^{(i)} - \left(\sum_{j=1}^m \hat{\alpha}^{(j)} t^{(j)} \phi(x^{(j)}) \right)^\top \phi(x^{(i)}) \right) \\
 &= \frac{1}{n_s} \sum_{i=1}^m \left(t^{(i)} - \sum_{j=1}^m \hat{\alpha}^{(j)} t^{(j)} k(x^{(i)}, x^{(j)}) \right)
 \end{aligned}$$

* If you are starting to get a headache, it's perfectly normal; it's an unfortunate side effect of the kernel trick.

* Linear classifier SVM cost function :

$$J(\omega, b) = \frac{1}{2} \omega^\top \omega + C \sum_{i=1}^m \max(0, 1 - t^{(i)}(\omega^\top x^{(i)} + b))$$

① Dual form of the linear SVM objective:

$$\underset{\alpha}{\text{minimize}} \quad \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^{(i)} \alpha^{(j)} t^{(i)} t^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \alpha^{(i)}$$

subject to $\alpha^{(i)} > 0$

after we find α ,

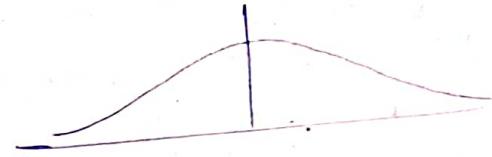
② the dual solution to the primal solution:

$$\hat{w} = \sum_{i=1}^m \hat{\alpha}^{(i)} t^{(i)} x^{(i)}$$

$$\hat{b} = \frac{1}{n_s} \sum_{i=1}^m (t^{(i)} - (\hat{w}^T x)^{(i)})$$

③ Normal distribution / Gaussian distribution : (standard)

$$f_z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \quad z \in \mathbb{R}$$



④ Binomial distribution:

$$f(k, n, p) = {}^n C_k p^k (1-p)^{n-k} = \Pr(X=k)$$

⑤ Maximum likelihood estimation:

generalised gaussian distribution with mean = μ & $SD = \sigma$,

$$N(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

the parameters μ and σ are represented together as

$$\theta = \{\mu, \sigma\}$$

for a dataset of given size 'n',

$$L = P(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n P(x_i | \theta)$$

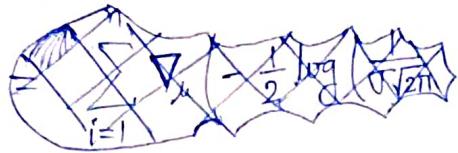
$$\hat{\theta}_{MLE} = \arg \max_{\theta} \prod_{i=1}^n P(x_i | \theta) \Rightarrow \frac{\partial}{\partial \theta} \prod_{i=1}^n P(x_i | \theta) = 0$$

maximizing $\prod_{i=1}^n P(x_i | \theta)$ means maximizing its logarithm.

$$\Leftrightarrow \frac{\partial}{\partial \theta} \log \left(\prod_{i=1}^n P(x_i | \theta) \right) = 0 \Rightarrow \frac{\partial}{\partial \theta} \left(\sum_{i=1}^n \log P(x_i | \theta) \right) = 0$$

$$\text{Now, } \sum_{i=1}^n \frac{\partial}{\partial \theta} \log(P(x_i|\theta)) = 0$$

$$\therefore \sum_{i=1}^n \nabla_\mu \left(\log \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \right) \right)$$



$$= \sum_{i=1}^n \nabla_\mu \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_i-\mu)^2}{2\sigma^2} \right)$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n \nabla_\mu (x_i - \mu)^2$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

$$\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\Rightarrow \boxed{\hat{\mu}_{MLE} = \frac{1}{n} \sum_i x_i}$$

Again,

$$\sum_{i=1}^n \nabla_\sigma \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_i-\mu)^2}{2\sigma^2} \right) = 0$$

$$\Rightarrow -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_i (x_i - \mu)^2 = 0$$

$$\Rightarrow \hat{\sigma}_{MLE} = \left(\frac{1}{n} \sum_i (x_i - \mu)^2 \right)^{1/2}$$

Example:

Suppose the weights of randomly selected college students are normally distributed with unknown μ and σ^2 .

Random sample: 115 122 130 127 149 160, 152 138
149 180

- Q. Identify the likelihood function
find maximum likelihood estimator $\hat{\theta}$.

$$L = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

Joint density
func

- o log likelihood func. $\log L = \sum_{i=1}^n \log f(x_i | \theta)$

- o Probability distribution functions

- ① Binomial distribution $\Rightarrow B(n, p) = \frac{n!}{k!(n-k)!} p^k q^{n-k}$,
 $[q = 1-p]$

mean = np

variance = npq

- ② Bernoulli distribution $\Rightarrow \text{Bern}(p) = [p^k (1-p)^{1-k}]$
 mean = p , variance = pq

- ③ Poisson Distribution $\Rightarrow P_0(\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$

mean = variance = λ

- ④ Normal / Gaussian distribution $\Rightarrow N(\mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 mean = μ , variance = σ^2

- ⑤ Standard normal: $\mu = 0, \sigma = 1$

scaling: $z = (x - \mu)/\sigma$, $f_z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$