

Probability Vs Likelihood

Some Example:

Example 1: Coin Toss



→ fair

$$\hookrightarrow P(H) = 0.5$$

→ mathematically

Bernoulli → 1 parameter
Distribution $\hookrightarrow P$

$$\text{Tail} \Rightarrow P(T) = 0.5$$

$$P_2 \rightarrow P(H) = 0.5$$

$$q = 1 - p = 0.5$$

$$P(x) = p^k + (1-p)(1-k)$$

$$k = 0, 1$$

Tail

* for Head $k=1$

$P(x) = p \rightarrow$ Probability of Head

$$P(x) = (1-p)$$

* If we know distribution then we also know parameter
and also know event and chance then probability

Likelihood

5 times coin flip

HHHHHH → what is possibility that $P(H) = 0.5$ and
coin is fair.

↑ we get 5 Head

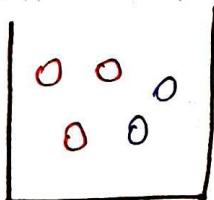
So, what is the chances that this is fair coin

$$\begin{matrix} H & H & H & H & H \\ \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ \text{or } & \text{or } & \text{or } & \text{or } & \text{or } 0.5 \end{matrix} \Rightarrow L = (0.5)^5$$

↳ very small

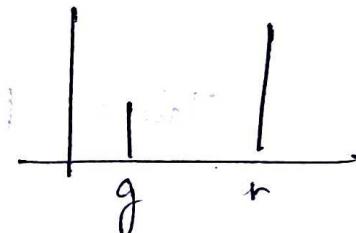
* The likelihood that the coin is fair given 5 times Head is ~~less~~ ~~very~~ ~~fair~~ ~~fair~~.

Example 2



$$P(A) = \frac{2}{5}$$

1 trial \rightarrow red
 \rightarrow green \rightarrow Bernoulli



ball is red $\rightarrow \frac{3}{5}$

$$pk + (1-p)(1-k)$$

$$k=0$$

$$(1-p) = 1 - \frac{2}{5} = \frac{3}{5}$$

likelihood

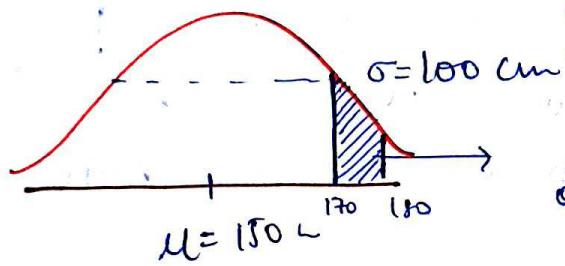
5 balls \rightarrow 4 4 4 4 4 $\rightarrow \left(\frac{2}{5}\right)^5 \rightarrow$ small

Example 3

continuous

↳ normal distribution

dist $\rightarrow (M, \sigma)$



Find area to get probability of getting height between 170-180 using μ ; σ

Likelihood

$$\boxed{100 \text{ cm}} \rightarrow \boxed{N(\mu, \sigma^2)}$$

\uparrow
event

$$L(\mu, \sigma^2 | X = 100) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\mu = 150 \\ \sigma = 10$$

$$= \frac{1.4 \times 10^{-7}}{\text{Area}}$$

Very small,

parameter \rightarrow event (prob)
 event \rightarrow parameter (likelihood)

Definition:

Probability: This is a measure of the chance that certain event will occur out of all the possible events. It's usually presented as a ratio or fraction, and it ranges from 0 (meaning the event will not happen) to 1 (meaning the event is certain to happen).

Likelihood: In statistical context, likelihood a function that measures the plausibility of a particular parameter value given some observed data. It quantifies how well a specific outcome support specific parameter values.

Maximum Likelihood Estimation (MLE)

coin toss

$$P = 0.5$$



$$\boxed{H H H H H} \rightarrow (0.5)^5 = L(P=0.5 \mid H H H H H)$$

$$\text{if } P = 0.6 \rightarrow (0.6)^5 = L(P=0.6 \mid H H H H H)$$

which one better $(0.5)^5$ or $(0.6)^5$

$$(0.5)^5 < (0.6)^5$$

| | |
|---|---|
| ○ | ○ |
| ○ | ○ |

$P(B) = \frac{2}{5}$

L: $\boxed{B B | B B B}$

$$L = (P=2/5 \mid H H H H H) = \left(\frac{2}{5}\right)^5$$

if:

| | |
|---|---|
| ○ | ○ |
| ○ | ○ |

$P(B) = \frac{3}{5}$

L: $(P=2/5 \mid B B B B B) = \left(\frac{3}{5}\right)^5$

which one better
better

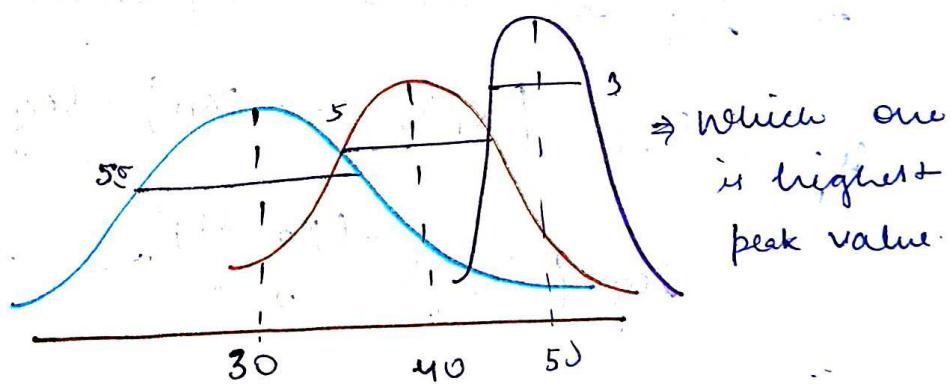
MLE for Normal Distribution

dataset = 100 or n

$\{$
 61
 32
 45
 63
 ...
 89

→ observable

* we know n data points follow normal distribution but which μ, σ value follow normal distribution



⇒ which one is highest peak value

$$L(\mu, \sigma | x_1, x_2, x_3, \dots, x_n) \quad \mu=110, \sigma=20$$

$$L(\mu=100, \sigma=10 | x_1=61, x_2=32, \dots, x_n=89) = 0.89$$

$$L(\mu=110, \sigma=20 | x_1=61, x_2=32, \dots, x_n=89) = 0.79$$

* Maximum likelihood is good 0.89

Let assume

we have only x_1 data

$$L(\mu, \sigma | x_1) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1 - \mu)^2}{2\sigma^2}}$$

If we have x_n , independent (x_i not depend on x_j) dataset and

$$L(\mu, \sigma | x_1, x_2, x_3, \dots, x_n) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1-\mu)^2}{2\sigma^2}} \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_2-\mu)^2}{2\sigma^2}} \times \dots \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}}$$

* log both side (log likelihood)

$$\log(L(\mu, \sigma | x_1, x_2, \dots, x_n)) = \boxed{\log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1-\mu)^2}{2\sigma^2}}\right) + \dots + \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}}\right)}$$

* focus on ① then apply all of them

$$\log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1-\mu)^2}{2\sigma^2}}\right) = \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \log e^{-\frac{(x_1-\mu)^2}{2\sigma^2}}$$

$$= \log(2\pi\sigma^2)^{-\frac{1}{2}} - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \log 2\pi - \frac{2}{n} \log \sigma - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$\log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1-\mu)^2}{2\sigma^2}}\right) = -\frac{1}{2} \log 2\pi - \log\left(\frac{(x_1-\mu)^2}{2\sigma^2}\right)$$

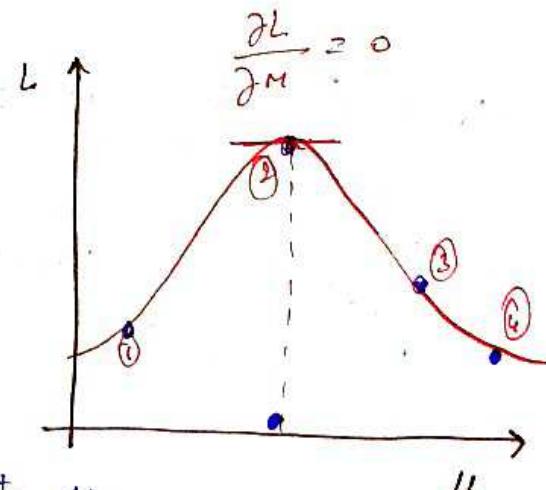
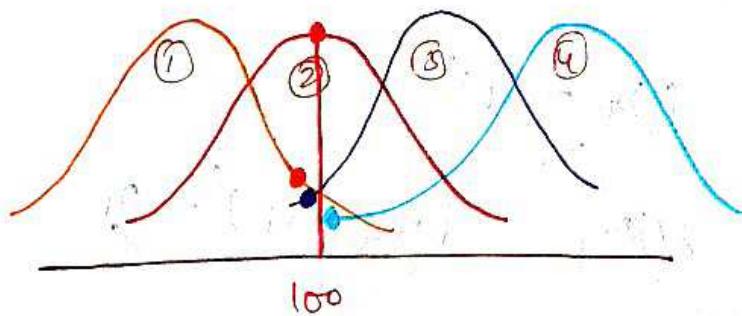
* Now apply for all eqn

$$= -\frac{1}{2} \log 2\pi - \log \sigma - \frac{(x_1 - \mu)^2}{2\sigma^2} - \frac{1}{2} \log 2\pi - \log \sigma -$$

$$\frac{(x_2 - \mu)^2}{2\sigma^2} - \dots - \frac{1}{2} \log 2\pi - \log \sigma - \frac{(x_n - \mu)^2}{2\sigma^2}$$

$$L(\mu) = \frac{n}{2} \log 2\pi - n \log \sigma - \frac{(x_1 - \mu)^2}{2\sigma^2} - \frac{(x_2 - \mu)^2}{2\sigma^2} - \dots$$

$$\frac{(x_n - \mu)^2}{2\sigma^2} - ①$$



eqn ① differentiate wrt μ

$$= 0 - 0 + \frac{\partial}{\partial \mu} \frac{(x_1 - \mu)}{2\sigma^2} + \frac{(x_2 - \mu)}{\sigma^2} + \dots + \frac{(x_n - \mu)}{\sigma^2} = 0$$

$$(x_1 - \mu) + (x_2 - \mu) + \dots + (x_n - \mu) = 0$$

$$x_1 + x_2 + x_3 + \dots + x_n - n\mu = 0$$

$$n\mu = x_1 + x_2 + \dots + x_n$$

$$\boxed{\mu = \frac{x_1 + x_2 + \dots + x_n}{n}}$$

eqn ① is differentiate by σ

$$\frac{\partial \log(L)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{2(x_1 - \mu)^2}{\sigma^3} + \frac{(x_2 - \mu)^2}{\sigma^3} + \dots + \frac{(x_n - \mu)^2}{\sigma^3}$$

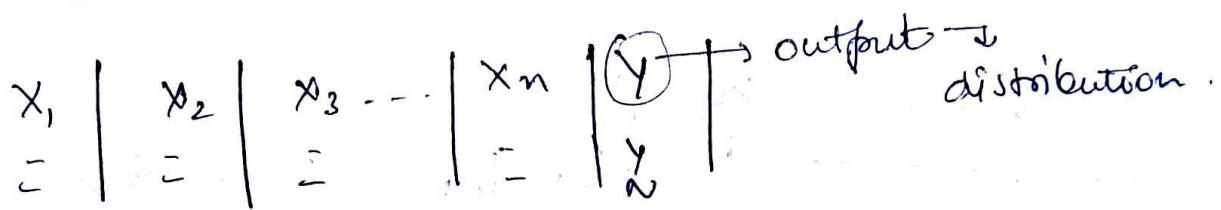
$$\frac{(x_1 - \mu)^2}{\sigma^3} + \frac{(x_2 - \mu)^2}{\sigma^3} + \dots + \frac{(x_n - \mu)^2}{\sigma^3} = \frac{n}{\sigma}$$

$$\frac{(x_1 - \mu)^2}{\sigma^2} + \frac{(x_2 - \mu)^2}{\sigma^2} + \dots + \frac{(x_n - \mu)^2}{\sigma^2} = n$$

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

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

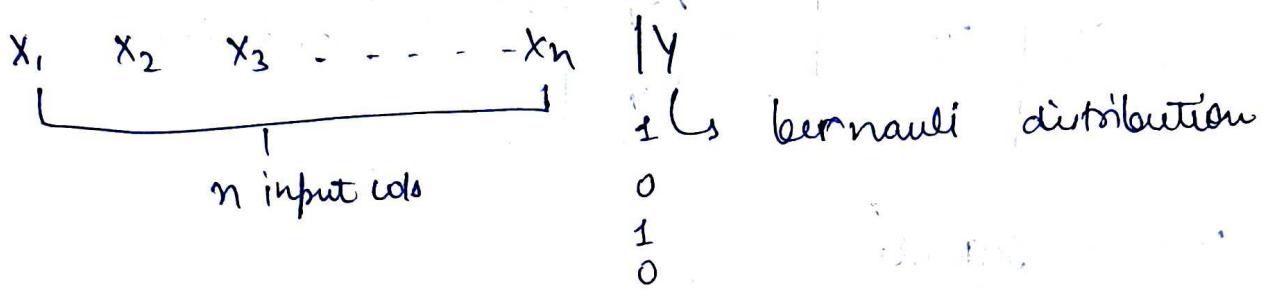
MLE in Machine Learning



$\gamma, N \rightarrow$ Bernauli

- 1.) Find out the distribution of Y/x
 - 2.) Decide to apply a ML model parametric in nature
 - ↳ logistic Regression
 - 3.) You randomly some value $\rightarrow \beta_0 \beta_1 \beta_2 \dots \beta_n$
 - 4.) Select a likelihood function.
 ex:- pmf = $p^k + (1-p)^{n-k}$
 - 5.) Value of $\beta_0 \beta_1 \beta_2 \beta_3 \dots \beta_n$
 - ↳ maximum value.

MLE in Logistic Regression



[1 query \rightarrow 1 or 0] \rightarrow bernoulli distribution ✓

[100 queries 1 or 0] \rightarrow Binomial distribution

1) Assume a distribution for y (Bernoulli distribution)

2) Logistic Regression \rightarrow parameter

$$\beta_0, \beta_1, \beta_2, \dots, \beta_n$$

3) Random Values $\rightarrow \beta_0, \beta_1, \dots, \beta_n$

* assuming all y 's are independent $y/x \rightarrow$ parameters β

$$L(Y/X; \beta) = P_k + (1-P)(1-k) \times P$$

$$L(Y/X; \beta) = P_1 Y_1 + (1-P_1)(1-Y_1)] \times P$$

for all rows of Y

$$\therefore k = Y$$

$$P = P(1)$$

$$1-P = P(0)$$

$$\begin{aligned}
 L(Y/X; \beta) &= P_1 Y_1 + (1-P_1)(1-Y_1) \times P_2 Y_2 + (1-P_2)(1-Y_2) \times \\
 &\quad \dots \dots P_n Y_n + (1-P_n)(1-Y_n)
 \end{aligned}$$

for easy mathematics

$$\frac{pk + (1-p)(1-k)}{L} \xrightarrow{\text{change}} p^y (1-p)^{(1-y)}$$

↳ output is same

$$L(Y/x; \beta) = p_1^{y_1} (1-p_1)^{1-y_1} \times p_2^{y_2} (1-p_2)^{1-y_2} \times \cdots \times p_n^{y_n} (1-p_n)^{1-y_n}$$

~~for~~ compact formula

$$L(Y/x; \beta) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

log both side

$$\log L(Y) = \sum_{i=1}^n \log \left(p_i^{y_i} (1-p_i)^{1-y_i} \right) = \sum_{i=1}^n \log p_i^{y_i} + \log (1-p_i)^{1-y_i}$$

$$\boxed{\log(L) = \sum_{i=1}^n y_i \log p_i + (1-y_i) \log (1-p_i)}$$

→ maximum → β
→ Gradient Descent

Some Important Question

1. Is MLE a general concept applicable to all machine learning algorithms.

Maximum Likelihood Estimation (MLE) is a general statistical concept that can be applied to many machine learning algorithms, particularly those that are parametric, but it's not applicable to all machine learning algorithms.

MLE is commonly used in algo. such as linear regression, logistic regression and neural network among others. These algorithms use MLE to find the optimal value of parameter that best fit the training data.

However, there are some machine learning algo. that don't rely on MLE. For example:

- ① Non-parametric methods
- ② unsupervised learning algo
- ③ Reinforcement learning

2. How is MLE related to the concept of loss function?

In machine learning, a loss function measures how well a model's predictions align with the actual values. The goal of training a machine learning model is often to find the model parameters that minimize the loss function.

Maximum Likelihood Estimation (MLE) is a method of estimating the parameters of a statistical model to maximize the likelihood function, which is conceptually similar to minimize a loss function, which is conceptually similar to minimize a loss function. In fact, for many common models, minimizing the loss function is equivalent to maximizing the likelihood function.

MLE and the concept of loss functions in machine learning are closely related. Many common loss functions can be derived from the principle of maximum likelihood estimation under certain assumptions about the data or the model. By minimizing these loss functions, we're effectively performing maximum likelihood estimation.

3. Then why does loss function exist, why don't we maximize likelihood.

The confusion arises from the fact that we're using two different perspective to look at the same problem.

In many machine learning algorithms, the aim is to minimize the difference between the predicted and actual values, and this is typically represented by a loss function. When we talk about minimizing the loss function, it's essentially the same as saying we're trying to find the best model parameters that gives us the closest prediction to the actual values.

On the other hand, when we looked at the problem from statistical perspective, we talk in terms the same as maximizing the likelihood function. In fact, many common loss function can be derived from the principle of MLE under certain assumption about data. Of maximizing the likelihood of seeing the observed data given the model parameters. This is represented by a likelihood function.

for many models, these two perspectives are equivalent - minimizing the loss function is the same as maximizing the likelihood function.

In fact, many common loss functions can be derived from the principle of MLE under certain assumptions about the data.

So why do we often talk about minimizing the loss function instead of maximizing the likelihood? There are a few reasons:

1. Computational Reasons! It's often easier and more computationally efficient to minimize a loss function than to maximize a likelihood function. This is particularly true when working with complex models like neural networks.

2. Generalization: The concept of a loss function is more general and can be applied to a wider range of problems. Not all machine learning problems can be framed in terms of maximizing a likelihood. For example, many non-parametric methods and unsupervised learning algorithms don't involve likelihoods.

3. Flexibility: Loss functions can be easily customized to the specific needs of a problem. For instance, we might want to give more weight to certain types of errors, or we might want to use a loss function that is robust to outliers.

4. Then why study about maximum likelihood at all?

The study of Maximum Likelihood Estimation (MLE) is essential for several reasons, despite the prevalence of loss functions in machine learning.

1. Statistical foundation: MLE provides a strong statistical foundation for understanding machine learning models. It gives a principled way of deriving the loss function used in many common machine learning algorithms, and it helps us understand why these loss function work and under what assumption.

2. Interpretability: The MLE framework gives us a way to interpret our model parameters. The MLEs are the parameters that make the observed data most likely under our model, which can be a powerful way of understanding what our model has learned.

3. Model Comparison: MLE gives us a way to compare different models on the same dataset. This can be done using tool like the Akaike information criterion (AIC) or the Bayesian information criterion (BIC), which are based on the likelihood function and can help us choose the best model for our data.

4. Generalization to Other Methods: MLE is a specific case of more general methods, like Expectation Maximization and Bayesian Inference which are used in more complex statistical modelling. Understanding MLE can provide a stepping stone to these more advanced topics.

5. Deeper Understanding: Lastly, understanding MLE can give us a deeper understanding of our models, leading to better institution, better model selection and ultimately, better performance on machine learning tasks.

Task → ① Likelihood function of softmax reg.

② MLE to linear Regression

$y \mid \hat{y}$ normal distribution
by hint

→ ③ → [mle]
OLS

Assumptions Logistic Regression

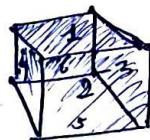
1. Binary Logistic Regression requires the dependent variable to be binary: That means the outcome variable must have two possible outcomes, such as "Yes" vs "no", "success" vs "failure", "spam" vs "not spam", etc.
2. Independence of observations: The observations should be independent of each other. In other words, the outcome of one instance should not affect the outcome of another.
3. Linearity of independent variable and log odds: Although logistic regression does not require the dependent and independent variables to be related linearly, it requires that the independent variables are linearly and related to the log odds.

4. Absence of multicollinearity

5. Large sample size.

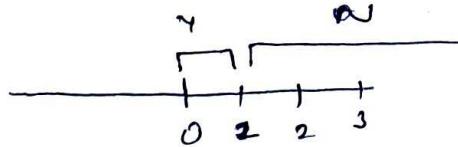
Odds and log(Odds)

odds: The odds of an event is the ratio of the probability of the event happening (P) to the probability of the event not happening ($1-P$). It's a way of expressing the likelihood of an event - If the odds are greater than 1, the event is more likely to happen than not and vice-versa



$$3 \text{ odds} \rightarrow P(3) = \frac{1}{6} \text{ and } P(13) = \left(\frac{5}{6}\right)$$

$$\text{Odds} = \frac{1/6}{5/6} = \frac{1}{5}$$



$$= 2yw$$

∞ not symmetric

That's why log. for symmetric

Another Interpretation of Logistic Regression

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

$$P_{\text{odd}} = 1 \quad \beta_1 = 1 \quad \beta_2 = 2$$

| Gpa | iq | place | \hat{y} |
|-----|----|-------|-----------|
| 8 | 80 | 1 | 0.73 |
| - | - | 5 | 0.37 |

$$P(u) = P$$

$$\hat{y} = P(u) = P = \frac{1}{1 + e^{-\beta x}}$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$P = \frac{1}{1 + e^{-\beta x}} \Rightarrow 1 + e^{-\beta x} = \frac{1}{P}$$

$$e^{-\beta x} = \frac{1}{P-1} \Rightarrow e^{-\beta x} = \frac{1-P}{P}$$

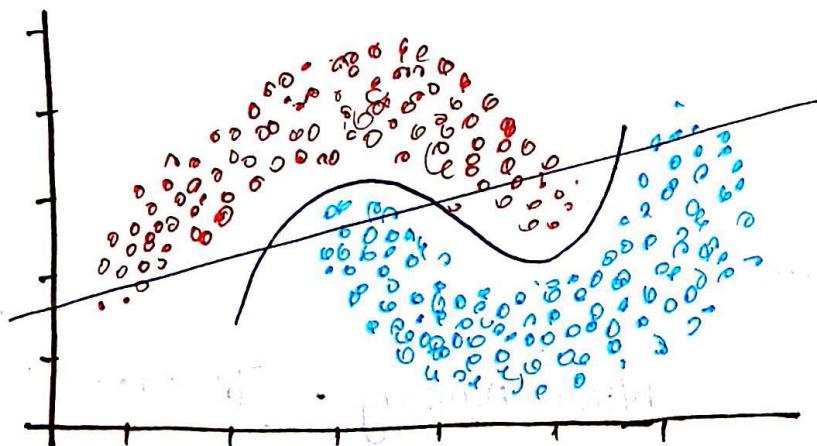
$$\frac{1}{e^{\beta x}} = \frac{1-P}{P} \Rightarrow \frac{P}{1-P} = e^{\beta x}$$

$$\log\left(\frac{P}{1-P}\right) = \beta x$$

$$\log\left(\frac{P}{1-P}\right) = e^{\beta x}$$

$$\log\left(\frac{P}{1-P}\right) = \beta x$$

Polynomial Features



Polynomial feature

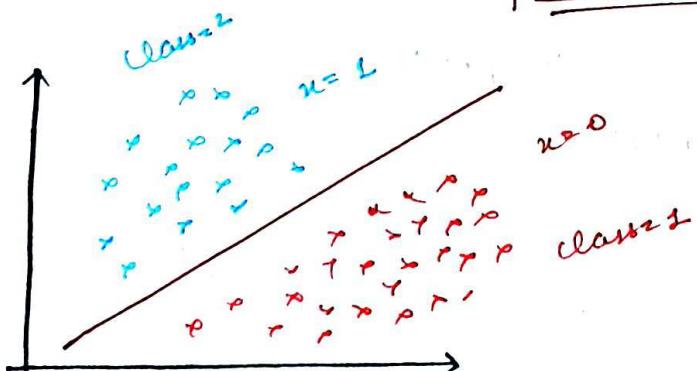
$x \rightarrow \text{degree}^2$

$$\begin{array}{c|c|c} x^0 & x^1 & x^2 \\ \uparrow & \uparrow & \uparrow \\ \beta_0 & \beta_1 & \beta_2 \\ & & \underline{\beta_3} \end{array}$$

degree $\uparrow^n \rightarrow$ Overfitting

degree $\ll \rightarrow$ underfitting

Performance Metrics



| Dataset | Actual | | Predicted |
|---------|--------|-------|-----------|
| | x_1 | x_2 | |
| - | - | 0 | 1 |
| - | - | 1 | 1 |
| - | - | 1 | 1 |
| - | - | 0 | 0 |
| - | - | 1 | 1 |
| - | - | 0 | 1 |
| - | - | 1 | 0 |