

right you 23

(*) Advance ML (*)

classmate

Date _____
Page _____

Def-A :- For $x_1, \dots, x_n \in X$

The dichotomies generated by \mathcal{H} on these points is defined by

$$\mathcal{H}(x_1, \dots, x_n) = \{h(x_1), \dots, h(x_n) \mid h \in \mathcal{H}\}$$

NOTE:- Higher the VC dimension, richer the hypothesis space...

what is dichotomy?

max. dichotomies possible = 2^n

Theorem:-

\Rightarrow If $m_{\mathcal{H}}(K) \leq 2^K$ for some value K then

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{K-1} N \binom{N}{i}$$

samples learning

what is growth function of hypothesis set \mathcal{H} ?

what is VC-dimension

Higher VC dim is good for test error!
 \rightarrow NO, it may overfit.

\Rightarrow By definition of VC-dimension...

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{\text{dvc}} N \binom{N}{i}$$

Higher the VC dim.
Higher the generalization error

$\Rightarrow \text{dvc} \geq N \Leftrightarrow$ there exists data

of size N such that

\mathcal{H} shatters the data.

\Rightarrow The VC-Generalization bound for any tolerance $\epsilon > 0$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8 \cdot \ln 4m_{\mathcal{H}}(2N)}{N}}$$

20th Jan

CLASSMATE

Date _____
Page _____

TASKS:- (1) Replicate the results. (dataset public?) .

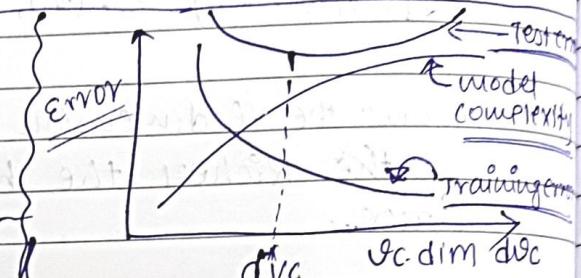
(2) DO something more → Tweak the algo...
→ Hyperparameter Tweaking...

(3) if code already available then

~~then~~ replicate it to another framework ...

$$m_H(N) \leq \sum_{i=1}^{dvc} \binom{N}{i}$$

$$m_H(N) \leq N^{dvc} + 1$$



• The VC bound is loose! -

- Hoeffding's inequality has some slack.
- $m_H(N)$ gives worst case estimate.
- Bounding $m_H(N)$ by polynomial of order.
- VC dimension is useful for comparing generalization performance of models.

(*) Bias and variance :-

$$E_{\text{out}}(g^{(D)}) = E_x[(g^{(D)}(x) - f(x))^2]$$

Expected value w.r.t (x)

$g^{(D)}$ is independent on (D) .

$$E_D[E_{\text{out}}(g^{(D)})]$$

$$= E_D[E_x[(g^{(D)}(x) - f(x))^2]]$$

$$= E_x[E_D[(g^{(D)}(x) - f(x))^2]]$$

$$= E_x[ED[g^{(D)}(x)^2] - 2ED[g^{(D)}(x)]f(x) + f(x)^2]$$

$$\text{let } ED[g^{(D)}(x)] = \bar{g}(x)$$

$$ED[E_{\text{out}}(g^{(D)})]$$

$$= E_{\text{out}}[ED(g^{(D)})]$$

$$= E_x[ED[g^{(D)}(x)^2] - 2\bar{g}(x)f(x) + f(x)^2]$$

$$= E_x[ED[g^{(D)}(x)^2] - \bar{g}(x)^2 + \bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2]$$

$$= E_x[ED[g^{(D)}(x)^2] - \bar{g}(x)^2 + \bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2]$$

- Rule of thumb:-

$$N \approx \log 1/dc$$

25th January 23



= Ex

$$= \text{Ex} [g^*(x)^2] - g(x)^2$$

$$= \text{Ex} [(g^{(t)}(x) - \bar{g}(x))^2]$$

→ Variance

$$\begin{aligned} g(x)^2 - 2\bar{g}(x)f(x) + f(x)^2 \\ = (\bar{g}(x) - f(x))^2 \end{aligned}$$

→ Bias

$$\text{Ex} [\text{Eout}(g^{(t)})]$$

$$= \text{Ex} [\text{bias}(x) + \text{var}(x)]$$

⇒ Here, along with the algorithm also matters...

(*) convex optimization in ML :-

outline:-

• convex optimization,

• Primal-Dual formulation

↳ KKT conditions

(*) The PLA Algorithm :-

$$w(0) = (0, 0, 0, 0, 0)$$

while there is a misclassified point $x(t), y(t)$

$$w(t+1) = w(t) + y(t) * (t).$$

⇒ if function is convex function then Gradient Descent will give global minimum.

optimization of ML

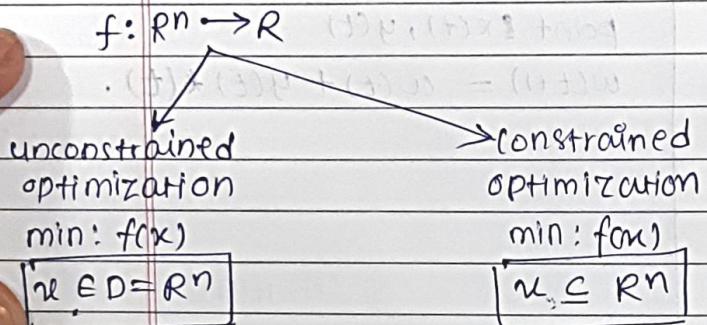
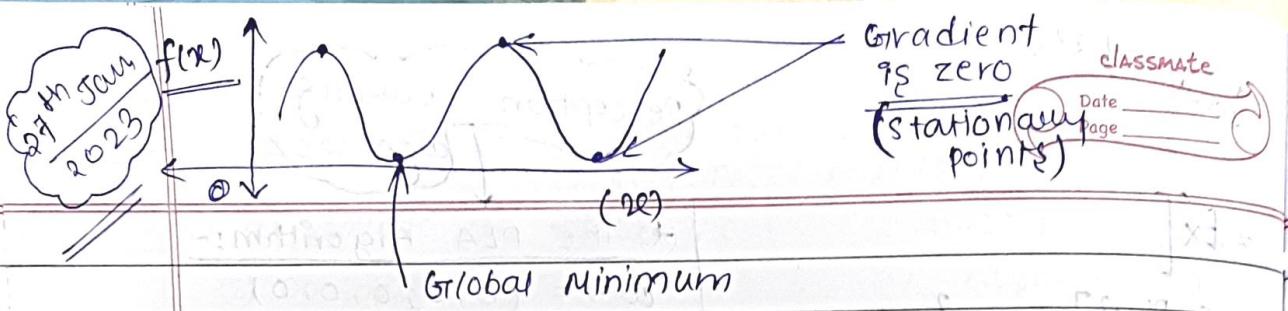
continuous

convex

Non-convex

discrete
(combinatorial)

↳ submodularity,



Example :-

Linear Regression

$A \in \mathbb{R}^{n \times d}$

$B \in \mathbb{R}^n$

$$\min_{x \in D} \|Ax - b\|^2$$

unconstrained.

Ridge Regression.

$$\min_{x \in D} \|Ax - b\|^2 + \lambda \|x\|^2$$

{ unconstrained.

$$\min_{x} \|Ax - b\|^2 \text{ such that}$$

$$\|x\|_2^2 \leq R$$

{ constrained.

(i) Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ continuous and differentiable fn.

For local optimum :

zero

vector:

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}_{n \times 1} = 0$$

$$\nabla^2 f(x) = H(x)$$

matrix ($n \times n$)

{ if Hessian is PSD:
local minimum

{ else if Hessian is NSD:
local maximum

Hessian Matrix

$$\frac{\partial^2 f(x)}{\partial x_1^2}, \dots, \frac{\partial^2 f(x)}{\partial x_n^2}$$

(*) Positive / Negative Semi-Definite:-

$\{\text{PSD}\} \Rightarrow$ all eigen values

$\{\text{NSD}\}$



$$\boxed{x^T A x} \geq 0 \quad \text{all } x \in \mathbb{R}^n$$

$$\boxed{x^T A x} \leq 0 \quad \forall x$$

$1 \times n \quad n \times n \quad n \times 1$

$\boxed{1x1}$

Quadratic form.

$\boxed{Ax \neq 0}$

if $x^T A x > 0$ then

it is $\textcircled{P.D.}$ matrix.

example: $A = I$ (Identity matrix) \Rightarrow all eigen values are greater than zero.

$$\Rightarrow x^T A x > 0$$

$$\Rightarrow x^T I x > 0$$

$$\Rightarrow x^T x > 0$$

$$\Rightarrow \|x\|_2^2 > 0$$

$$\boxed{A > B} \quad \uparrow \quad A \text{ is PD w.r.t } B.$$

$$\text{means } \boxed{A - B > 0}$$

(*) constrained optimization problem :-

$$f_0: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\min f_0(x)$$

such that

$$f_i(x) \leq 0 \quad i \in [1, m]$$

$$f_j(x) = 0 \quad j \in [1, l]$$

classmate

Date _____

Page _____

8th. Feb 23

Lagrangian form of C.O.P

$$f_0(x) + \sum_{i=1}^m \lambda_i^0 f_i(x) + \sum_{j=1}^l \mu_j^0 f_j(x)$$

$$\|Ax - b\|_2^2 \rightarrow \|Ax - b\|_2^2 + \lambda \|x\|_2^2$$

$$\text{s.t. } \|x\|_2^2 \leq R$$

λ^0, μ^0 is Lagrangian multipliers.

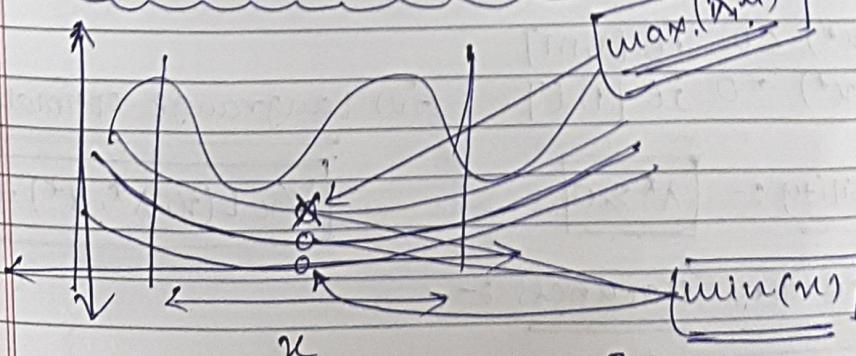
$$L(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i^0 f_i(x) + \sum_{j=1}^l \mu_j^0 f_j(x)$$

what if x is feasible point?

For feasible x , $L(x, \lambda, \mu) \leq f_0(x)$

$$\min_x L(x, \lambda, \mu) \leq f_0(x^*)$$

$$\text{s.t. } \lambda \geq 0$$



$$\max_{\lambda, \mu} \left[\min_x L(x, \lambda, \mu) \right] \leq f_0(x^*)$$

8th
Feb
2023

$$\min_u L(x, \lambda, u) = \underline{g(\lambda, u)}$$

classmate

Date _____

Page _____

$$\max_{(x,u)} g(\lambda, u) \leq f_0(x^*)$$

dual function.

$$g(\lambda, u) = \min_u L(x, \lambda, u)$$

$\Rightarrow (\lambda^*, u^*)$ are point of maximum. maximum

$$\Rightarrow g(\lambda, u) \leq g(\lambda^*, u^*)$$

$$\Rightarrow g(\lambda, u) \leq g(\lambda^*, u^*) \leq f_0(x^*)$$

maximum value of the dual function is equal to less than or equal to primal function

$$f_0(x^*) - g(\lambda^*, u^*) \geq 0$$

weak duality

theorem

duality gap.

\Rightarrow Primal Problem is convex problem.

(usually)
(not always)

strong duality hold

$$f_0(x^*) - g(\lambda^*, u^*) = 0$$

constraint should satisfy some condition

(*) KKT conditions :-

x^* \leftarrow point of primal optimal
 λ^*, u^* \leftarrow point of dual optimal
with zero duality gap.

(i) Primal Feasibility :-

$$f_i^*(u^*) \leq 0 \quad i \in [1, m]$$

$$f_j^*(x^*) = 0 \quad j \in [1, l]$$

iv) Lagrange optimality :-

(ii) Dual feasibility :-

$$\lambda^* \geq 0$$

$$\nabla_x L(x, \lambda^*, u^*) = 0$$

(iii) Complementary slackness :-

$$\lambda_i^* f_i^*(x^*) = 0 \quad \forall i$$

8th Feb
2023

Any tuple (x, λ, u) is called KKT point

if it satisfies the KKT conditions

when you have $(x^*, \lambda^*, u^*) \Rightarrow$ KKT point.

classmate

Date _____

Page _____

if duality gap = 0

x^* } primal optimal

λ^* } dual optimal

u^*

(x^*, λ^*, u^*) is KKT

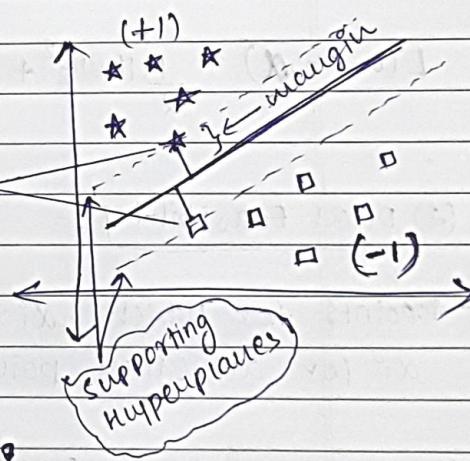
point

If Problem is convex

9th Feb. 2023

(*) Hard Margin SVM :-

support
vectors.



How New data point will be classified?

$\Rightarrow x_{\text{new}} \leftarrow$ new data point.

\rightarrow if $(w^T x_{\text{new}} + w_0) \geq 1$ then positive class.
else negative class ...

How to detect misclassification?

if $y_i(w^T x_i + w_0) < 1$ then point is misclassified...

$$\hookrightarrow 0 < 1 - y_i(w^T x_i + w_0).$$

Hinge

$$h[(x_i, y_i), (w, w_0)] = \max(0, 1 - y_i(w^T x_i + w_0))$$

9th Feb.
2023



classmate

Date _____
Page _____

$$\text{Min}_{(w, w_0)} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } y_i(w^\top x_i + w_0) \geq 1 \quad \forall i \in [1, n]$$

primal form of

Hard Margin.

= KKT conditions :-

(1) primal Feasibility :-

$$1 - y_i(w^\top x_i + w_0) \leq 0$$

w^*, w_0^*

Primal variables.

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^\top x_i + w_0))$$

(2) Dual Feasibility :- $\alpha_i \geq 0 \quad \forall i \in [1, n]$

→ points for which $\alpha_i > 0$ are the support vectors.

$\alpha_i = 0$ for all other points is zero.

(3) complementary slackness :-

$$\alpha_i (1 - y_i(w^\top x_i + w_0)) = 0 \quad \forall i \in [1, n]$$

(4) Lagrange optimality

$$\begin{aligned} L(w, w_0, \alpha) &= \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^\top x_i + w_0)) \\ &= \frac{1}{2} w^\top w + \sum_{i=1}^n \alpha_i (1 - y_i(w^\top x_i + w_0)) \end{aligned}$$

$$\nabla_w L(w, w_0, \alpha) = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i$$

9th
Feb
2023

$$L(w, w_0, \alpha) = -\sum_{i=1}^n \alpha_i y_i = 0$$

classmate

Date _____
Page _____

$$\boxed{\sum_{i=1}^n \alpha_i y_i = 0}$$

$$\begin{aligned} \alpha_i & \text{ for } i \in C^+ \\ \alpha_j & \text{ for } j \in C^- \\ \sum \alpha_i &= \sum \alpha_j \end{aligned}$$

Support Vectors :-

$$\{x_i | \alpha_i^* > 0\}$$

$$\min_{(w, w_0)} \frac{1}{2} \|w\|_2^2$$

such that $1 - y_i(w^T x_i + w_0) \leq 0 \quad \forall i \in [1, n]$

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + w_0))$$

$$\boxed{g(\alpha) = \min_{(w, w_0)} L(w, w_0, \alpha)}$$

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i \quad w_0 \text{ satisfies } \sum_{i=1}^n \alpha_i y_i = 0$$

$$\begin{aligned} g(\alpha) &= \frac{1}{2} w^{*T} w_0^* + \sum_{i=1}^n \alpha_i (1 - y_i(w^{*T} x_i + w_0^*)) \\ &= \frac{1}{2} w^{*T} w_0^* + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i w^{*T} x_i - \sum_{i=1}^n \alpha_i y_i w_0^* \end{aligned}$$

$$= \frac{1}{2} w^{*T} w_0^* + \sum_{i=1}^n \alpha_i - w^{*T} w^*$$

$$= \sum_{i=1}^n \alpha_i - \frac{1}{2} w^{*T} w^*$$

$$= \sum_{i=1}^n \alpha_i - \frac{1}{2} \left(\left(\sum_{i=1}^n \alpha_i y_i x_i \right)^T \left(\sum_{j=1}^n \alpha_j y_j x_j \right) \right)$$

$$\text{Max } g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

s.t. $x_i > 0$ \uparrow dual form of hard margin

$$x_i \sum_{i=1}^n \alpha_i y_i = 0$$

10th Feb 23

(*) Kernel SVM :-

$$x \in \mathbb{R}^d$$

$$x \rightarrow z \in \mathbb{R}^h$$

$$\text{s.t. } [n > d]$$

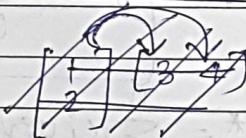
$$x_i \Rightarrow z_i \Rightarrow \phi(x_i)$$

$$\text{label}(x_{\text{test}}) = \text{sign} \left[w^* T \phi(x_{\text{test}}) + w_0 \right]$$

$$w^* = \sum_{i=1}^m \alpha_i y_i \phi(x_i)$$

$$\text{label}(x_{\text{test}}) = \text{sign} \left[\sum_{i=1}^m \alpha_i y_i \phi(x_i)^T \phi(x_{\text{test}}) + w_0 \right]$$

- How to implement SVM in higher dimension without actually getting $\phi(x)$:-



Mercer's Theorem :-

Kernel fun :-

$$K: X \times X \rightarrow \mathbb{R}$$

A symmetric function $K(x_i^o, x_j^o) = K(x_j^o, x_i^o)$

can be represented as $K(x_i^o, x_j^o) = \phi(x_i^o)^T \phi(x_j^o)$
for some (ϕ) .

IFF,

$$K = \begin{pmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ ; & ; & ; & \\ K(x_m, x_1) & K(x_m, x_2) & \dots & K(x_m, x_n) \end{pmatrix}_{(m \times n)}$$

Gram matrix.

• Example of kernels :-

- Linear Kernel :-

$$K(x_i^o, x_j^o) = x_i^o \cdot x_j^o$$

- Polynomial Kernel :-

$$K(x_i^o, x_j^o) = (1 + x_i^o \cdot x_j^o)^t ; t > 0$$

- Gaussian Kernel (RBF kernel) :-

$$K(x_i^o, x_j^o) = \exp\left(-\frac{\|x_i^o - x_j^o\|_2^2}{2\sigma^2}\right)$$

NOTE! A necessary and sufficient condition to make our custom kernel is that the Gram Matrix should be positive definite.

(*) Soft SVM :-

Regularization term

$$\text{L1-SVM} : - \min_{(w, w_0, \varepsilon)} \frac{1}{2} \|w\|_2^2 + \frac{1}{2} \sum_{i=1}^n \varepsilon_i$$

$$\text{L2-SVM} : - \min_{(w, w_0, \varepsilon)} \frac{\lambda}{2} \|w\|_2^2 + \frac{1}{2} \sum_{i=1}^n \varepsilon_i^2$$

$$\text{L1C SVM} : - \min_{(w, w_0, \varepsilon)} \frac{C}{n} \sum_{i=1}^n \varepsilon_i + \frac{\lambda}{2} \|w\|_2^2$$

$$\text{L2C SVM} : - \min_{(w, w_0, \varepsilon)} \frac{C}{n} \sum_{i=1}^n \varepsilon_i^2 + \frac{\lambda}{2} \|w\|_2^2$$

Example :- let $x \in \mathbb{R}$.

$$\begin{aligned}
 K(x_i^o, x_j^o) &= \exp\left(-\frac{(x_i^o - x_j^o)^2}{2\sigma^2}\right) \\
 &= \exp\left(-\frac{x_i^{o2}}{2\sigma^2} - \frac{x_j^{o2}}{2\sigma^2} + \frac{2x_i^o x_j^o}{2\sigma^2}\right) \\
 &= \exp\left(\frac{-x_i^{o2}}{2\sigma^2}\right) \cdot \exp\left(\frac{-x_j^{o2}}{2\sigma^2}\right) \cdot \exp\left(\frac{2x_i^o x_j^o}{2\sigma^2}\right) \\
 &= \exp\left(\frac{-x_i^{o2}}{2\sigma^2}\right) \exp\left(\frac{-x_j^{o2}}{2\sigma^2}\right) \cdot \cancel{\exp\left(\frac{1 + x_i^o x_j^o + x_i^{o2} x_j^{o2} + \dots}{2\sigma^2}\right)}
 \end{aligned}$$

10th Feb.
2023

$$= \exp\left(-\frac{x_i^2}{\sigma^2}\right) \cdot \exp\left(-\frac{x_j^2}{\sigma^2}\right) \cdot \left[1 \frac{x_i}{\sigma} \frac{x_i^2}{\sqrt{2}\sigma^2} \dots\right] \left[1 \frac{x_j}{\sigma} \frac{x_j^2}{\sqrt{2}\sigma^2} \dots\right]$$

classmate
Date _____
Page _____

$$= \exp\left(-\frac{x_i^2}{\sigma^2}\right) \left[1 \frac{x_i}{\sigma} \frac{x_i^2}{\sqrt{2}\sigma^2} \dots\right] \cdot \exp\left(-\frac{x_j^2}{\sigma^2}\right) \left[1 \frac{x_j}{\sigma} \frac{x_j^2}{\sqrt{2}\sigma^2} \dots\right]$$

$\hat{\Phi}(x_i)$ $\hat{\Phi}(x_j)$

(*) Low Rank Approximation :-

$$\min_{A \in \mathbb{R}^{n \times d}} \|A - A'\|_F^2$$

subject to $\text{rank}(A) = k$

$$\|Ux\|_2 = \|x\|_2$$

$$\|Vx\|_2 = \|x\|_2$$

↑ Orthogonal.

$$A_k = A = \sum_{i=1}^k \sigma_i u_i v_i^T$$

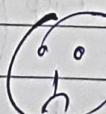
solution

(*) SMOTE :- Synthetic Minority Over-Sampling Technique.

- Avoid overfitting due to exact replicas of minority class samples.
- subset of minority class is taken.
- New synthetic data samples are created from subsets and added.



No overfitting.
Reduced overfitting



Mainly introduce noise.

27th
Feb
2023

(3. x | θ) || prior = (L(θ), $P(\theta)$)

classmate

Date _____
Page _____

$$P(\theta | D, \alpha) \propto P(D | \theta, \alpha) \cdot P(\theta | \alpha)$$

Posterior

Likelihood prior

Max

Probability
[0, 1]

Likelihood
 ≥ 1 .
(can be)

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$P(\theta | \{x_1, y_1\}, \{x_2, y_2\}, \dots, \{x_n, y_n\}, \alpha)$$

$$\propto P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, \theta, \alpha)$$

$$\cdot P(\theta | \alpha)$$

$$P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, \theta) = L(\theta)$$

$$= \prod_{i=1}^n P(y_i | x_i, \theta)$$

\Rightarrow Form of prior

similar to posterior

then it is called as conjugate prior.

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i, \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P(y_i | x_i, \theta)$$

\Rightarrow we use log likelihood
since, it is monotonically increasing.

$$= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n -\log P(y_i | x_i, \theta)$$

\hookrightarrow negative log likelihood.

classmate
Date _____
Page _____

~~loss at single point~~

$$l(\theta; (x_i, y)) = -\log P(y|x_i, \theta)$$

~~loss at single point~~

$$\text{Lreg}(\theta; (x_i, y)) = (\theta^T x_i - y)^2$$

$$P(y_i|x_i, (w, b)) = N(y_i|w^T x_i, \frac{1}{\beta})$$

$$\epsilon \sim N(0, \frac{1}{\beta})$$

$$\text{Likelihood} = \prod_{i=1}^n P(y_i|x_i, w, b)$$

$$= \prod_{i=1}^n N(y_i|w^T x_i, \frac{1}{\beta})$$

$$= \prod_{i=1}^n \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\beta(y_i - w^T x_i)^2\right)$$

for multivariate,

$$f_x(u) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (u - \mu)^T \Sigma^{-1} (u - \mu)\right)$$