

Week 6: 4.1 - 4.5

Objectives:-

PAGE No.	
DATE	/ /

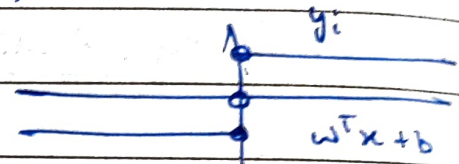
- Write linear classⁿ eqⁿ
- List reasons for knowing linear class.
- Write Bayesian decision fn
- Derive perceptron algorithm
- Derive grad desc for logistic reg.
- Derive loss fn for primal SVM

• What is linear classification?

Binary classⁿ

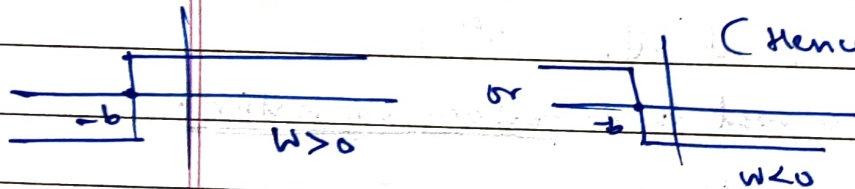
$x_i \in \mathbb{R}^D$ for each x_i given $t_i \in \{-1, +1\}$

Find $y_i = \text{Sign}(W^T x_i + b)$ $y_i \neq t_i$



Example 1D case:-

$$y = (Wx + b) \text{ Sign} \quad \text{WLOG assume } W > 0 \Rightarrow \text{Sign depends on}$$
$$= \text{Sign}(x + b/W) \quad [\text{if constant is } > 0] \quad W, b.$$



(Hence $\|W\| = 1$)

Conceptually a simple threshold.

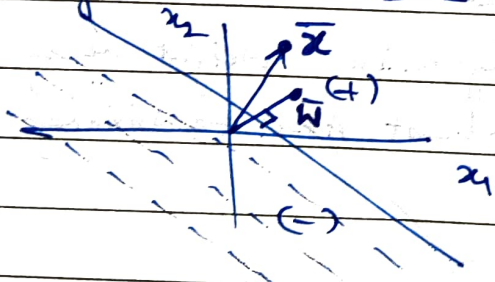
Example 2-D case:-

dimension

$$y = \text{Sign}(W_1 x_1 + W_2 x_2 + b)$$

$$\|W\| = 1$$

a line



$$\text{Consider } W_1 x_1 + W_2 x_2 + b = 0$$

$$W^T x + b = 0$$

One side +ve, one side -ve

Taking dot product with \bar{W} and seeing if product > 0 or not
Contains if we scale \bar{W} and $b \iff$ [Again $\|W\| = 1$]
Orienting hyperplane by $\|W\|$ around unit sphere/circle.

(4.2) Why study linear classifiers?

• Simplest classifiers
(one of two)

• Many non-linear problems
can be linearized

• Natural outcome of
for familiar class of
class conditional densities.

Linearizing non linear problems:-

- Add derived features:

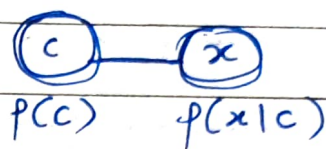
- Powers
- Interaction terms
- kernels

or use pretrained Neural Networks.

PAGE No.	
DATE	/ /

Bayesian decision rule for classification:

C_1 cats | C_2 dogs $\cdot P(C|x)$

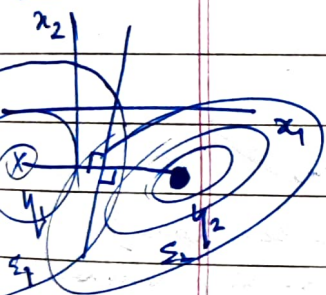


$$P(C|x) = \frac{P(x|C) P(C)}{P(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

We may find $\frac{P(C_1|x)}{P(C_2|x)}$ and if $> 1 \Rightarrow C_1$
 $< 1 \Rightarrow C_2$ [for example]

(Assuming equal risk) \leftarrow

Gaussian class conditional in D-Dimensions:



$P(x|C_1)$ $P(x|C_2)$ both modelled as gaussian

$$P(x|C_j) = \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu_j)^T \Sigma^{-1} (x - \mu_j)\right)$$

D dimensions

[General setting]

Boundary may or may not be linear.

We wish to look at:

$$\log\left(\frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)}\right) \text{ compared with } 0.$$

$$\ln[P(x|C_1)] + \ln[P(C_1)] - \ln[P(x|C_2)] - \ln[P(C_2)]$$

(If σ_i & σ_j are equal, we can cancel out the constant)

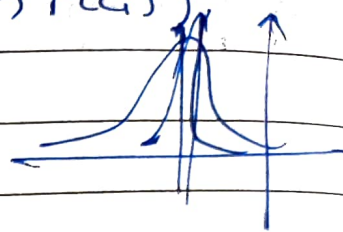
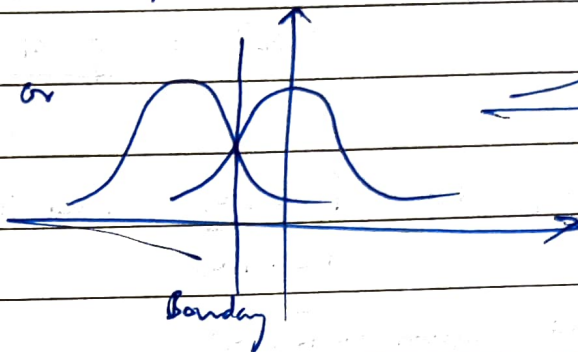
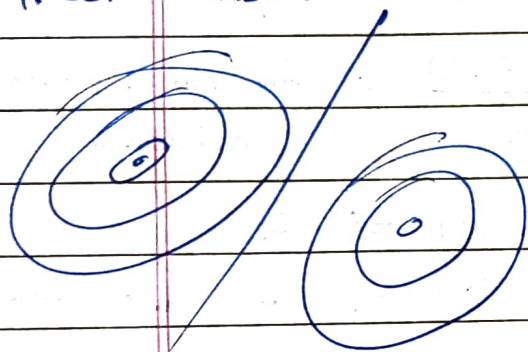
$$\begin{aligned} \Rightarrow \ln(P(x|C_1)) + \ln(P(C_1)) - \ln(P(x|C_2)) - \ln(P(C_2)) \\ = -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + K_1 \\ - \frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) - K_2 \end{aligned}$$

We get a straight line if $\Sigma_1 = \Sigma_2 = \Sigma$.

In general we will not get straight line

4.3 Algorithm to build a Bayesian Classifier:

1. Assume parametric forms of $p(x|c_i)$
2. Set $p(c_i) = N_i/N$
3. Find the parameters of $p(x|c_i)$ using methods such as ML.
4. Set decision criteria $y = \arg \max (P(x|c_i) P(c_i))$



4.4 Perceptron learning Algorithm & loss:-

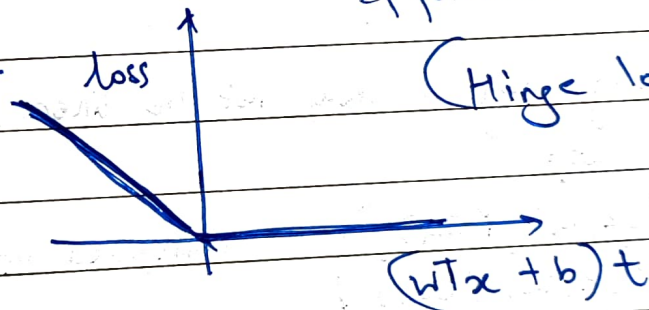
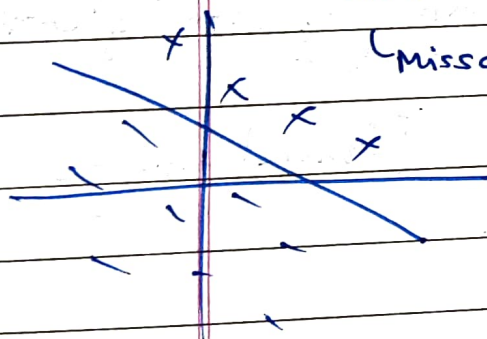
$$L(W) = - \sum_{i \in M} (W^T x_i + b) t_i$$

$\rightarrow > 0$ for $+ve$ < 0 for $-ve$

(Missclassified)

(We wish to achieve zero loss)

If points are not lin. sep. we have zero loss.



(Hinge loss)

Hinge loss:-

$$\max (0, -(W^T x_i + b) t_i)$$

Convex fns
Can do gradient descent

$$W \leftarrow W - \eta \nabla_W L(W)$$

$$W \leftarrow W + \eta x_i t_i \quad i \in M \text{ (Misclassified)}$$

Infinite solutions if we have a lin. sep problem

4.5 Linear classification - Logistic Regression:

Circling back to Bayesian Decisions:

$$a = \log \left(\frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)} \right)$$

PAGE No.	
DATE	/ /

We want to estimate $P(C_1|x) = \sigma(a) = \frac{1}{1 + \exp(-a)}$

[Now 'a' is a linear fn of x only if both are Gaussian and have the same Σ .]

$$a = \frac{P(x|C_1) P(C_1)}{\sum_i P(x|C_i) P(C_i)}$$

Gradient Descent for logistic regression:

$p(t|w)$ [t ∈ {0,1}] maximize

$$= \prod_{i=1}^N y_i^{t_i} (1-y_i)^{1-t_i} \quad (\text{iid hence product})$$

When $t=0$ \swarrow $t=1$

$y_i = P(t_i=1 | x_i, w)$
(monotonic)

We take log transformation as log is ~~convex~~ & then we may deal with \sum terms instead of \prod terms

We get:

$$\underset{w}{\text{max}} \left[- \sum_{i=1}^N t_i \log(y_i) + (1-t_i) \log(1-y_i) \right]$$

$y_i = w^T x_i + b$ if assume Gaussian & same Σ

Cross Entropy loss between t_i & y_i (Minimize cross entropy)



$$- \int \log(q) p \, dx$$

$$\nabla_w L = \sum_i (y_i - t_i) x_i \quad \text{This + gradient descent = solution.}$$

$$y_i = \sigma(w^T x_i + b)$$

$$+ \lambda \|w\|_2^2 \quad (\text{Hence regularization})$$

[Generalized Linear Models]

can be done.