

Today

- Classification
- Binary least squares classifier (separating plane)
- Extensions to multiple classes
- One vs rest / One vs one / K classes discriminant
- Probabilistic classifiers (in particular logistic regression)

General idea: Assign \mathbf{x} feature vectors one of K distinct classes C_1, \dots, C_k

→ the result of this is a division of the input space into decision regions where boundaries are called decision boundaries.

We will start with simple models for which the decision boundaries are linear functions of the input vector

→ Datasets where classes can be separated exactly by linear decision surfaces are called **linearly separable**.

→ One of the most compact representation for class labels
is to take $t^{(i)} = \{0, 1\}$

For multi-class problems, we will use a one-hot
encoding scheme. E.g if we have K classes, then a
feature vector from class C_2 will be represented as

$$\vec{t} = (0, 1, 0, \dots, 0) \in \mathbb{R}^K$$

Discriminant function takes a feature vector as input
and assign it to one of the K classes

We will start by considering linear discriminant of the

form

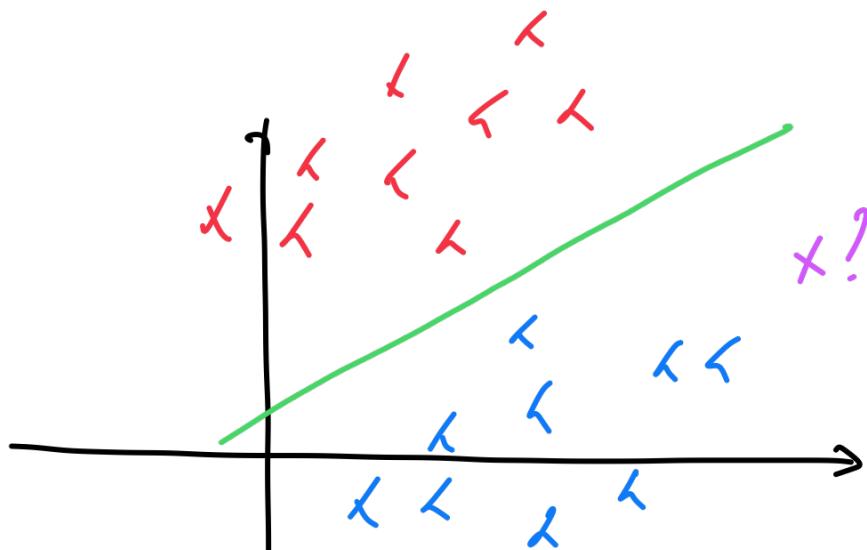
$$y(x) = \beta^T x + \beta_0$$

Decision boundary given by $y(x) = 0$

Assuming $t^{(i)} = +1$ for $x^{(i)} \in C_1$ & $t^{(i)} = -1$ for $x^{(i)} \in C_2$

a new point will be assigned to C_1 if $y(x) > 0$

C_2 $y(x) < 0$

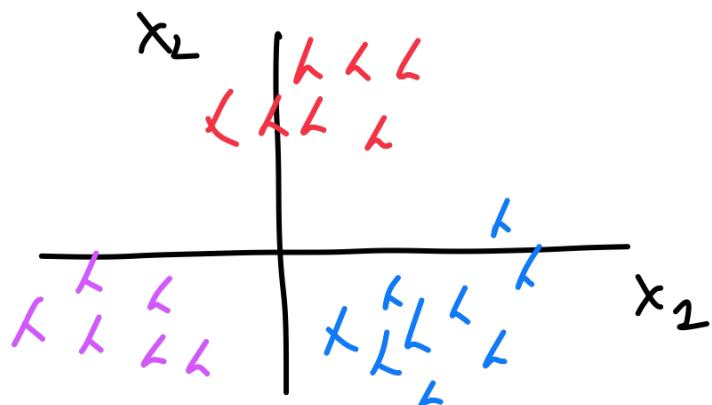


General objective:
learn the discriminant
function $y(x)$ that returns
a positive value for the red points
and a negative value for the blue prototypes

→ One of the simplest approach at learning our classifier is through least squares

$$\beta^* = \underset{\beta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_d x_d^{(i)}))^2$$

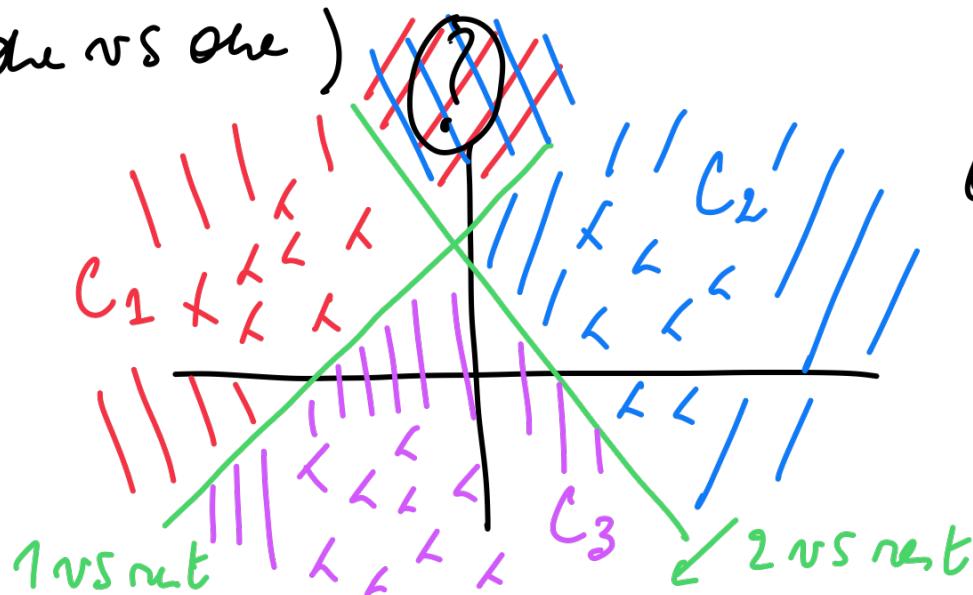
How about the Multiple class setting?



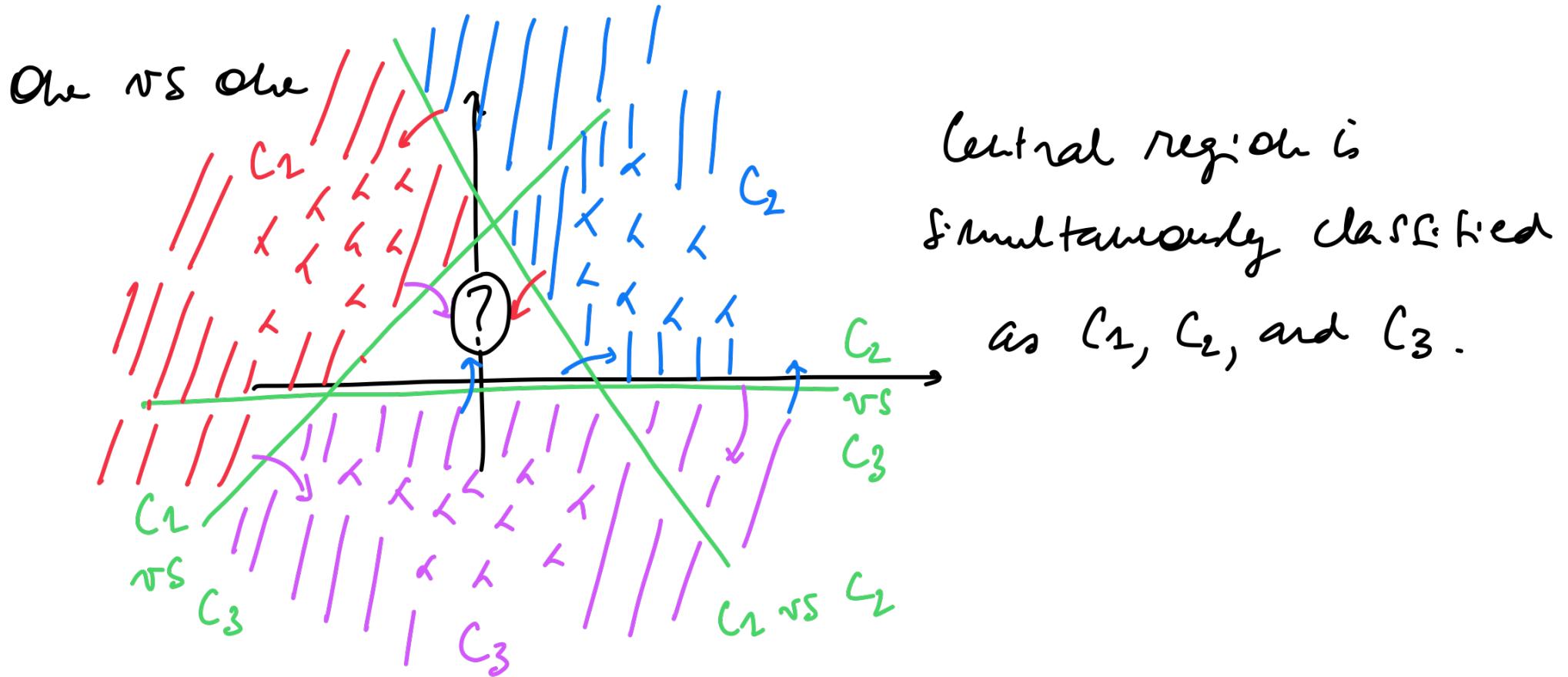
→ Solution 1 : Learn $k-1$ binary classifiers which separate each of the classes C_1, \dots, C_{k-1} from the rest of the data \mathcal{X} .

(one vs all / one vs the rest)

→ Solution 2 Learn a set of $\frac{k(k-1)}{2}$ classifiers separating class C_i from class $C_j \rightarrow$ class of a new point is determined based on the outcomes of the $\frac{k(k-1)}{2}$ classifiers following a



majority vote



→ In order to resolve those ambiguities, we can rely on a single K class discriminant function (learning K linear discriminant functions at once.)

→ encode the targets using a 1 of K coding scheme

$$t^{(i)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{from which we can define the target matrix}$$

$$T = \begin{bmatrix} | & | & | \\ t^{(1)} & t^{(2)} & \dots & t^{(N)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{K \times N}$$

We then consider a set of K separating planes β_k

which we show as the rows of B

$$B = \begin{bmatrix} \overrightarrow{\beta}_1^T \\ \vdots \\ \overrightarrow{\beta}_K^T \end{bmatrix}$$

where each $\overrightarrow{\beta}_k$ is a vector

$$\overrightarrow{\beta}_k = (\beta_{k0}, \beta_{k1}, \dots, \beta_{kD}) \in \mathbb{R}^{D+1}$$

We can then write the deviations between the targets and the predictions from the discriminant functions as

as

$$\underline{T} - \underline{B} \cdot \underline{\tilde{X}}^T = \begin{bmatrix} | & | & | \\ t^{(1)} & t^{(2)} & \dots & t^{(N)} \\ | & | & | \end{bmatrix} - \begin{bmatrix} -\vec{\beta}_1^T \\ \vdots \\ -\vec{\beta}_K^T \end{bmatrix} \begin{bmatrix} | & | \\ \tilde{x}^{(1)} & \tilde{x}^{(2)} \dots \tilde{x}^{(N)} \\ | & | \end{bmatrix} \underbrace{\underline{\tilde{X}}^T}_{\underline{\tilde{X}}^T}$$

and one can then minimize the residual sum of squares less

as before

$$\frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K (t_k^{(i)} - \vec{\beta}_k^T \underline{\tilde{X}}^{(i)})^2 \quad (*)$$

Minimizing $(*)$ is equivalent to looking for a set of planes where the plane associated to the k^{th} class outputs a value as close as possible to 1 for points that are in C_k and 0 otherwise.

$$\begin{aligned}
 (*) &= \text{sum of squared entries of the matrix } T - B \tilde{X} \\
 &= \text{Tr}((T - B \tilde{X}^T)^T (T - B \tilde{X}^T))
 \end{aligned}$$

→ Computing the derivatives of $(**)$ with respect to B and setting it to 0, we get a matrix equation of the normal equations which we can solve for B

$$\hat{B}^T = \left(\begin{matrix} \tilde{X}^T & \tilde{X} \\ \tilde{X} & \tilde{X} \end{matrix} \right)^{-1} \tilde{X}^T T T^T$$

So far: Supervised learning

- linear regression (w/ OLS loss)
- $\{x^{(i)}, t^{(i)}\}_{i=1}^N \rightarrow$ gradient descent (including batch, stochastic)
 - Normal equations
- regularization (Best subset selection,
 l_p regularization including
Ridge & LASSO)
- Statistical intuition (including MLE / MAP)
- Evolution of the MSE and Bias / Variance
Decomposition as a function of Model Complexity

Linear Classification ($t^{(i)} \in \{1, 2, \dots, K\}$)

→ Least Squares binary classifier

→ Extensions to multiple classes through one vs rest
or one vs one and limitations

→ Single K class discriminant classifier (least squares)
(Relying on the hot encoding of the $t^{(i)}$)

$$\vec{t}^{(i)} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ if } x^{(i)} \in C^1 \quad T = \begin{bmatrix} 1 & 1 & 1 \\ \vec{t}^{(1)} & \vec{t}^{(2)} & \dots & \vec{t}^{(N)} \\ 1 & 1 & 1 \end{bmatrix}$$

Classification Matrix

$$= \begin{bmatrix} -\vec{\beta}^{(2)}- \\ | \\ -\vec{\beta}^{(K)}- \end{bmatrix} \quad \tilde{x} = \begin{bmatrix} -\tilde{x}^{(1)}- \\ | \\ -\tilde{x}^{(N)}- \end{bmatrix}$$

→ K classes discriminant: solve for B in

$$\min_B \sum_{i,j} (T - B \tilde{X}^T)_{ij}^2$$

Today: — Probabilistic classifiers

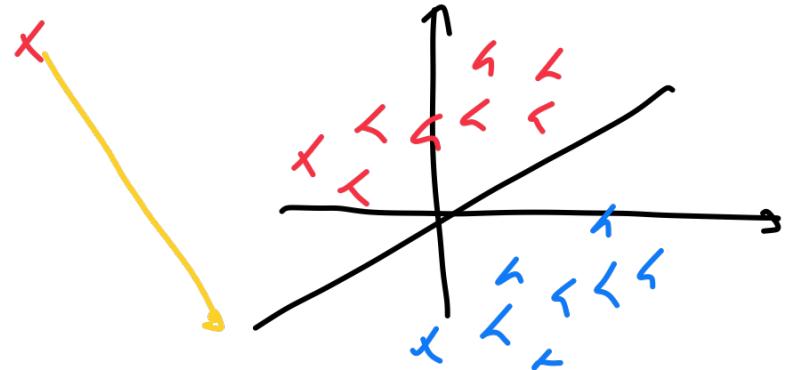
Discriminative → logistic regression
 → definition
 → learning through MLE and
 binary cross entropy

Generative → Gaussian Discriminant Analysis
 linear Discriminant Analysis

logistic Regression

Recall the least squares loss

$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \tilde{x}^{(i)})^2$$



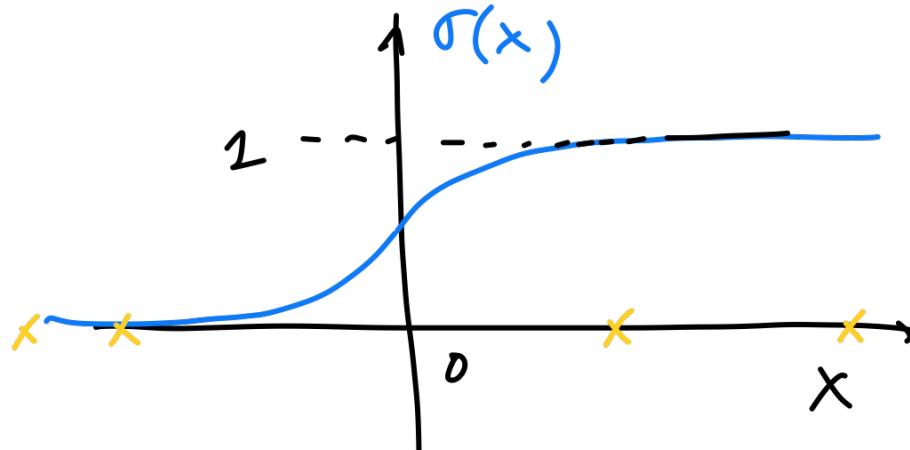
→ 2 issues: output is real (unbounded) while he would prefer a binary or confidence measure

→ risk of misclassification in the presence of outliers

- One fix is to introduce an activation function.
- One possibility to improve interpretability is to choose a function that outputs a value in $[0, 1]$ (can then be considered as measure of how "confident" we are of our classification).
- $y(x) = 1 \rightarrow$ means $x \in C^1$ with 100% certainty
- $y(x) = 0 \quad // \quad x \in C^0 \quad //$
- $y(x) \in (0, 1) \rightarrow$ non zero risk of misclassification

An example of such function is the sigmoid function $\sigma(x)$

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



→ Whenever we are with respect to the decision boundary, the decision remains bounded.

From this, we can define our logistic regression classifier as

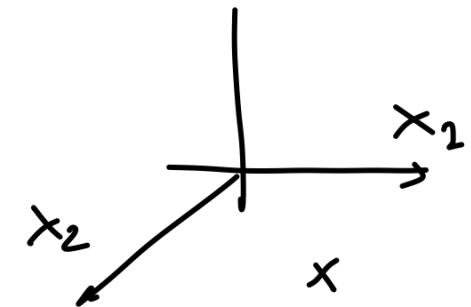
$$h_{\beta}(x) = \sigma(\beta^T \tilde{x}) = \frac{1}{1 + e^{-\beta^T \tilde{x}}}$$

→ Since our classifier outputs a value in $[0, 1]$ we can interpret this output as a probability.

letting

$$P(t^{(i)} \in C_1 | x^{(i)}) = h_{\beta}(x^{(i)})$$

$$P(t^{(i)} \in C_0 | x^{(i)}) = 1 - h_{\beta}(x^{(i)})$$



We can (1) consider the probability of observing our complete set of training pairs \$\{x^{(i)}, t^{(i)}\}_{i=1}^N\$ and then look for the weight vector \$\beta\$ which maximizes this probability

$$P(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\}_{i=1}^N; \beta) = \prod_{i=1}^N P(t(x^{(i)}) = t^{(i)} | x^{(i)}; \beta)$$

provided that \$t^{(i)}\$'s are independent

$$= \prod_{i=1}^N \left\{ \begin{array}{ll} P(t(x^{(i)}) = 1 | x^{(i)}; \beta) & \text{if } t^{(i)} = 1 \\ P(t(x^{(i)}) = 0 | x^{(i)}; \beta) & \text{if } t^{(i)} = 0 \end{array} \right. (*)$$

$$(*) = \prod_{i=1}^N p(t(x^{(i)}) = 1 | x^{(i)}; \beta)^{t^{(i)}} p(t(x^{(i)}) = 0 | x^{(i)}; \beta)^{(1-t^{(i)})}$$

→ taking the log and maximizing, we get

$$\hat{\beta}_{MLE} = \underset{\beta}{\operatorname{argmax}} \log \left(\prod_{i=1}^N h_{\beta}(x^{(i)})^{t^{(i)}} (1 - h_{\beta}(x^{(i)}))^{(1-t^{(i)})} \right)$$

$$= \underset{\beta}{\operatorname{argmax}} \sum_{i=1}^N t^{(i)} \log(h_{\beta}(x^{(i)})) + (1-t^{(i)}) \log(1-h_{\beta}(x^{(i)}))$$

$$= \underset{\beta}{\operatorname{argmin}} - \sum_{i=1}^N t^{(i)} \log(h_{\beta}(x^{(i)})) + (1-t^{(i)}) \log(1-h_{\beta}(x^{(i)}))$$

Log Loss = Binary Cross Entropy

→ the minimizer to $(*)$ can then be found through gradient descent. Taking the gradient of the binary cross entropy loss with respect to β ,

$$\begin{aligned} & \frac{\partial}{\partial \beta} \left\{ -t^{(i)} \log(h_{\beta}(x^{(i)})) - (1-t^{(i)}) \log(1-h_{\beta}(x^{(i)})) \right\} \\ &= -t^{(i)} \frac{\frac{\partial}{\partial \beta} h_{\beta}(x^{(i)})}{h_{\beta}(x^{(i)})} - (1-t^{(i)}) \cdot \underbrace{-\frac{\frac{\partial}{\partial \beta} h_{\beta}(x^{(i)})}{1-h_{\beta}(x^{(i)})}}_{(*)*} \end{aligned}$$

$$\frac{\partial}{\partial \beta} h_{\beta}(x^{(i)}) = \frac{\partial}{\partial \beta} \sigma(\beta^T \tilde{x}^{(i)}) = \frac{\partial}{\partial \beta} \left(\frac{1}{1+e^{-\beta^T \tilde{x}^{(i)}}} \right) (*)$$

$$(*) = \frac{1}{(1 + e^{-\beta^T \tilde{x}^{(i)}})^2} \cdot \left(+ \tilde{x}^{(i)} \cdot e^{-\beta^T \tilde{x}^{(i)}} \right)$$

$$= \frac{1}{(1 + e^{-\beta^T \tilde{x}^{(i)}})} \left(1 - \frac{1}{1 + e^{-\beta^T \tilde{x}^{(i)}}} \right) \cdot \tilde{x}^{(i)}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \Rightarrow \frac{\partial}{\partial \beta}(\sigma(\beta^T \tilde{x}^{(i)})) = \sigma(\beta^T \tilde{x}^{(i)}) (1 - \sigma(\beta^T \tilde{x}^{(i)})) \dots \tilde{x}^{(i)}$$

Substituting this in $(**)$ we get

$$\begin{aligned} \frac{\partial l}{\partial \beta} &= - \frac{t^{(i)}}{\sigma(\beta^T \tilde{x}^{(i)})} \cdot \cancel{\sigma(\beta^T \tilde{x}^{(i)})} (1 - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}^{(i)} \\ &\quad + \frac{(1 - t^{(i)})}{1 - \sigma(\beta^T \tilde{x}^{(i)})} \cancel{\sigma(\beta^T \tilde{x}^{(i)})} (1 - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}^{(i)} \end{aligned}$$

$$\begin{aligned}\frac{\partial l}{\partial \beta} &= -t^{(i)} (1 - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}^{(i)} + (1 - t^{(i)}) \underbrace{\sigma(\beta^T \tilde{x}^{(i)})}_{\text{in red}} \tilde{x}^{(i)} \\ &= (t^{(i)} - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}^{(i)}\end{aligned}$$

From which we can find β_{MLE} as

$$\beta^{(k+1)} = \beta^{(k)} - \eta (t^{(i)} - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}^{(i)}$$

gradient steps for binary cross entropy on the logistic regression classifier.

Probabilistic classifiers

→ Discriminative classifier

(given a feature vector $x^{(i)}$, a discriminative model outputs a decision)

→ learns a model for $p(t|x)$

→ Generative classifier

Idea: learn a representation for the objects from the two classes then classify a new example based on how it compares to our model of what typical instances should look like

→ learns a model for $p(x|t)$

Today

Probabilistic classifiers

- Discriminative classifiers $p(t|x)$
- Generative classifiers $p(x|t)$
- Gaussian Discriminant Analysis
(GDA) (QDA)
- Quadratic Discriminant Analysis
- Linear Discriminant Analysis (LDA)
- Comparison between logistic regression
& LDA

→ Perceptrons, perceptron learning rule, convergence

→ Neural Networks

Generative Classifier

1) learn a model for $p(x|t)$

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

$$\{x^{(i)}, t^{(i)}\}_{i=1}^N$$

2) to classify a new prototype x , use Bayes rule

$$p(t|x) = \frac{p(x|t)p(t)}{p(x)} = \frac{p(x|t)p(t)}{\sum_{\forall c} p(x|t=c)p(t=c)}$$

in order to determine the class of a new point $x^{(i)}$ we
maximize $p(t=t_c | x^{(i)})$ over t_c

$$t(x^{(i)}) = \underset{t_c}{\operatorname{argmax}} p(t=t_c | x^{(i)}) = \underset{t_c}{\operatorname{argmax}} \frac{p(x^{(i)} | t=t_c)p(t=t_c)}{\sum_{\forall c'} p(x^{(i)} | t=t_{c'})p(t=t_{c'})}$$

the denominator is constant across the t_c 's we can just focus on the numerator and define our class as

$$t(x^{(i)}) = \arg \max_{t_c} p(x^{(i)} | t = t_c) p(t = t_c)$$

→ To determine the class of a new prototype $x^{(i)}$ based on a generative classifier, we therefore need a model for $p(x|t)$ and for $p(t)$

→ One of the simplest model for $p(x|t)$ is the Multivariate Normal distribution

For $p(t)$ a reasonable approach is to use a Bernoulli

$$p(t=0) = \pi_0 \quad \text{and} \quad p(t=1) = \pi_1 = 1 - \pi_0$$

$$p(x|t=t_c) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\mu_c)^T \Sigma_c^{-1} (x-\mu_c))$$

Together these models describe the Gaussian Discriminant Analysis generative classifier

In logistic regression we trained our model following a Maximum likelihood approach on $p(t|x)$

When training generative classifiers we will again use a Max likelihood approach but this time on $\{p(\{x, t\})\}$

For the GDA model, recall that we have

$$p(x|t=t_c) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp(- (x-\mu_c)^\top \tilde{\Sigma}^{-1} (x-\mu_c))$$

$$p(t=t_c) = \Pi_c$$

$$p(\{x, t\}) = p(x|t) p(t)$$

As before, if we assume that our examples are independent
we can write $p(\{x^{(i)}, t^{(i)}\}_{i=1}^N)$ as

$$p(\{x^{(i)}, t^{(i)}\}_{i=1}^N) = \prod_{i=2}^N p(\{x^{(i)}, t(x^{(i)})=t^{(i)}\}) \quad (*)$$

$$P(\{x^{(i)}, t(x^{(i)}) = c\}) = \begin{cases} P(x^{(i)} | t(x^{(i)}) = c) p(t(x^{(i)}) = c) & \text{if } c = t^{(i)} \\ 0 & \text{otherwise} \end{cases}$$

$$P(\{x^{(i)}, t^{(i)}\}_{i=1}^N) = \prod_{i=1}^N \left(P(x^{(i)} | C_0) p(C_0) \right)^{1-t^{(i)}} \times \left(P(x^{(i)} | C_1) p(C_1) \right)^{t^{(i)}}$$

QUADRATIC DISCRIMINANT

ANALYSIS

$$P(\{x^{(i)}, t^{(i)}\}_{i=1}^N) = \prod_{i=1}^N \left(\frac{1}{(2\pi)^{D/2} |\Sigma_0|^{1/2}} \exp(-\frac{1}{2}(x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0)) \pi_0 \right)^{1-t^{(i)}} \times \left(\frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} \exp(-\frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)) (1 - \pi_0) \right)^{t^{(i)}}$$

When taking $\Sigma_0 = \Sigma_1$ we recover the simpler LINEAR DISCRIMINANT ANALYSIS MODEL.

$$p(\{x^{(i)}, t^{(i)}\}_{i=1}^N) = \prod_{i=1}^N \left(\frac{1}{(2\pi)^{D/2} |\Sigma|^{\frac{1}{2}}} \exp(-(\bar{x}^{(i)} - \mu_0)^T \Sigma^{-1} (\bar{x}^{(i)} - \mu_0)) \pi_0 \right)^{t^{(i)}} \\ \times \left(\frac{1}{(2\pi)^{D/2} |\Sigma|^{\frac{1}{2}}} \exp(-(\bar{x}^{(i)} - \mu_1)^T \Sigma^{-1} (\bar{x}^{(i)} - \mu_1)) (1 - \pi_0) \right)^{1-t^{(i)}}$$

Taking the log and maximizing we get

$$LL = \log(p(\{x^{(i)}, t^{(i)}\}_{i=1}^N))$$

$$\begin{aligned} \Sigma^*, \mu_0^*, \mu_1^*, \pi_0^* &= \operatorname{argmax}_{\Sigma, \mu_0, \mu_1, \pi_0} \sum_{i=1}^N -\log ((2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}) \\ (\star) &+ \sum_{i=2}^N -\frac{1}{2} (x^{(i)} - \mu_0)^T \overbrace{\Sigma^{-1}}^{\text{---}} (x^{(i)} - \mu_0) (1 - t^{(i)}) \\ (\star\star) &+ \sum_{i=2}^N -\frac{1}{2} (x^{(i)} - \mu_1)^T \overbrace{\Sigma^{-1}}^{\text{---}} (x^{(i)} - \mu_1) t^{(i)} \\ &+ \sum_{i=1}^N (1 - t^{(i)}) \log(\pi_0) + \sum_{i=1}^N t^{(i)} \log(1 - \pi_0) \end{aligned}$$

In the LDA setting, we get closed form expressions for each of the $\mu_0, \mu_1, \Sigma, \pi_0$ by computing the derivatives and setting them to 0.

Starting with π_0 , taking the derivative of the log likelihood function, we get

$$\frac{\partial \underline{LL}}{\partial \pi_0} = \underbrace{\sum_{i=1}^N (1 - t^{(i)})}_{\text{Number of examples in class } C_0} \cdot \frac{1}{\pi_0} + \underbrace{\sum_{i=1}^N t^{(i)}}_{\text{Number of examples in class } C_1} \cdot \frac{-1}{1 - \pi_0} = 0$$

$$N_0 \cdot \frac{1}{\pi_0} - N_1 \cdot \frac{1}{1 - \pi_0} = 0$$

Number of
examples in class
 C_0

Number of
examples in class
 C_1

$$\Leftrightarrow N_0(1 - \pi_0) = N_1 \pi_0 \quad \pi_0 = \frac{N_0}{N_0 + N_1}$$

For μ_0, μ_1 (the two means) we follow a similar reasoning on (*) and (**)

$$\frac{\partial}{\partial \mu_0} - \sum_{i=1}^N (x_i - \mu_0)^T \Sigma^{-1} (x_i - \mu_0) (1 - t^{(i)})$$

$$= \sum_{i=2}^N \cancel{1} (1 - t^{(i)}) \cancel{\sum_{i=1}^N (x_i - \mu_0)} = 0$$

$$\frac{\partial}{\partial \mu_1} - \sum_{i=1}^N (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) t^{(i)}$$

$$= \sum_{i=1}^N \cancel{2} t^{(i)} \cancel{\sum_{i=1}^N (x_i - \mu_1)} = 0$$

$$\begin{aligned} & \frac{\partial}{\partial x} x^T A x \\ &= 2Ax \end{aligned}$$

} Setting the derivatives to zero

$$\sum_{i=1}^N (1-t^{(i)}) x^{(i)} = \sum_{i=1}^N (1-t^{(i)}) \mu_0$$

$$\sum_{i=1}^N t^{(i)} x^{(i)} = \sum_{i=2}^N t^{(i)} \mu_1$$

$\sum_{i \in C_0} x^{(i)}$ N_0
 $\sum_{i \in C_1} x^{(i)}$ N_1

$$N_0 \mu_0 = \sum_{i \in C_0} x^{(i)} \quad \Rightarrow \quad \mu_0 = \frac{1}{N_0} \sum_{i \in C_0} x^{(i)}$$

$$N_1 \mu_1 = \sum_{i \in C_1} x^{(i)} \quad \mu_1 = \frac{1}{N_1} \sum_{i \in C_1} x^{(i)}$$

Finally For the Covariance Σ , using the fact that

$$\frac{\partial}{\partial A} \log(|A|) = (A^{-1})^T \quad (*)$$

Computing the derivative of the First 3 terms in the log likelihood w.r.t Σ^{-1} , we get

$$\begin{aligned} \frac{\partial}{\partial \Sigma^{-1}} \left\{ -\frac{N}{2} \underbrace{\log(|\Sigma|)}_{\text{using}} + \sum_{i=1}^N (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) (1 - t^{(i)}) \right. \\ \left. + \sum_{i=1}^N (x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1) t^{(i)} \right\} \end{aligned}$$

using
 $\log(|\Sigma|) = \log(|\Sigma|^T) = -\log(|\Sigma^T|)$ (*)

$$(*) = \frac{\partial}{\partial \Sigma^{-1}} \left(N \log(|\Sigma|) - \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu}_0)^T \Sigma^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_0) (1 - t^{(i)}) - \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1) t^{(i)} \right)$$

using (*)

as well as $\frac{\partial}{\partial A} (X^T A X) = X X^T$ we get

$$N \sum - \left\{ \sum_{i=1}^N (1 - t^{(i)}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_0) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_0)^T \right\}_{i \in C_0} + \left\{ \sum_{i=1}^N t^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T \right\}_{i \in C_1} = 0$$

Final solution for the covariance

$$\Sigma = \frac{1}{N} \left\{ \sum_{i \in C_0} (x^{(i)} - \mu_0) (x^{(i)} - \mu_0)^T + \sum_{i \in C_1} (x^{(i)} - \mu_1) (x^{(i)} - \mu_1)^T \right\}$$

Today

Generative classifiers

- linear Discriminant Analysis
- Connection between LDA & logistic regression
- Naive Bayes Classifier
- Perceptron algorithm (including training, convergence)
- Neural Networks (including training through back propagation)

I) Connection between LDA & logistic regression

in LDA

$$p(t=c|x, \theta) \propto \frac{\exp\left(-\frac{1}{2}(x-\mu_c)^T \Sigma^{-1} (x-\mu_c)\right) \pi_c}{\sum_{c'} \exp\left(-\frac{1}{2}(x-\mu_{c'})^T \Sigma^{-1} (x-\mu_{c'})\right) \pi_{c'}}$$

$\theta = (\pi_0, \mu_0, \mu_1, \Sigma)$
(vector of parameters)

c = class, in binary setting $c \in \{0, 1\}$

$$p(t=c|x, \theta) \propto \frac{\exp(\mu_c^T \Sigma^{-1} x - \frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu_c^T \Sigma \mu_c) \pi_c}{\sum_{c'} \exp(\mu_{c'}^T \Sigma^{-1} x - \frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu_{c'}^T \Sigma \mu_{c'}) \pi_{c'}}$$

$$p(t=c | x, \theta) \propto \frac{\exp(\mu_c^T \Sigma^{-1} x - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \pi_c) \exp(-\frac{1}{2} x^T \Sigma^{-1} x)}{\sum_{c'} \exp(\mu_{c'}^T \Sigma^{-1} x - \frac{1}{2} \mu_{c'}^T \Sigma^{-1} \mu_{c'} + \log \pi_{c'}) \exp(-\frac{1}{2} x^T \Sigma^{-1} x)}$$

(*)

Assuming 2 classes $c \in \{0, 1\} \rightarrow \theta = (\pi_0, \pi_1, \mu_0, \mu_1, \Sigma)$

and letting

$$\beta_0 = \Sigma^{-1} \mu_0 \quad \gamma_0 = -\frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \log \pi_0$$

$$\beta_1 = \Sigma^{-1} \mu_1 \quad \gamma_1 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \log \pi_1$$

Substituting $\beta_0, \beta_1, \gamma_0, \gamma_1$ in the class posterior (*) we get

$$p(t=c | x, \theta) = \frac{\exp(\beta_0^T x + \gamma_0)}{\exp(\beta_0^T x + \gamma_0) + \exp(\beta_1^T x + \gamma_1)}$$

$$p(t=0 \mid X, \theta) = \frac{1}{1 + \exp((\beta_1 - \beta_0)^T X + \gamma_1 - \gamma_0)}$$
$$= \sigma((\beta_0 - \beta_1)^T X + (\gamma_0 - \gamma_1))$$

→ LDA can read as a particular instance of a logistic classifier

II Naive Bayes Classifier

→ Generative classifier

$$p(x|y=c, \theta) = \prod_{j=1}^D p(x_j|y=c, \theta_{jc})$$

→ Naive Bayes assumes that the features are conditionally independent given the class label

Depending on the nature of the features we should use

Gaussian Naive Bayes (Real valued features)

$$p(x|y=c, \theta) = \prod_{j=1}^D N(x_j|\mu_{jc}, \sigma_{jc}^2) \quad (\text{LDA with diagonal covariance})$$

Bernoulli Naïve Bayes (binary features)

$$p(x|y=c, \theta) = \prod_{j=1}^D \text{Ber}(x_j | \mu_{jc})$$

Categorical Naïve Bayes (Categorical features $x_j \in \{1, \dots, k\}$)

$$p(x|y=c, \theta) = \prod_{j=1}^D \text{Cat}(x_j | p_{jc}) \quad p_{jc} \text{ probabilities}$$

$$\text{with } \sum_k p_{jc} = 1$$

$$p_{jc} \geq 0$$

III Perceptron

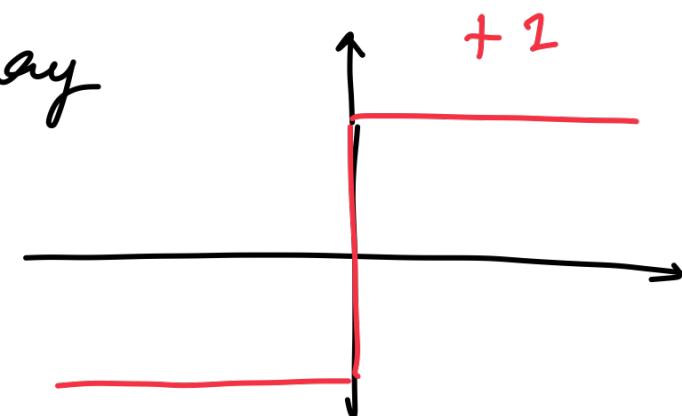
→ combination of linear model with Heaviside activation

$$y(x) = \sigma(\beta^T \tilde{x}) \quad \text{with} \quad \sigma(x) = \begin{cases} 1 & x > 0 \\ -1 & x \leq 0 \end{cases}$$

→ cannot be trained in the usual way

as the gradient vanishes almost

everywhere



$$\ell(\beta) = - \sum_{i \in \text{Misclassified}} t^{(i)} (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})$$

(*)

$\{x^{(i)}, t^{(i)}\}_{i=1}^N$

$t^{(i)} \in \{\pm 1\}$

We then minimize $\ell(\beta)$ with respect to β and then recompute the new set of misclassified examples.

For a point that is strongly misclassified, we have

$$\text{Sign}(f^{(i)}) \neq \text{Sign}(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_d x_d^{(i)})$$

hence a large positive contribution to $(*)$

Minimizing $(*)$ will thus "correct" the plane β based on the misclassified points

→ the loss $(*)$ can be minimized using an SGD approach over the misclassified points.

→ In this setting the SGD iterations are given by

$$\beta \leftarrow \beta - \gamma \underset{\beta}{\operatorname{grad}} \ell(\beta)$$

$$\beta \leftarrow \beta + \gamma \tilde{x}^{(i)} t^{(i)} \quad (\text{perceptron learning rule})$$

→ Refining Perceptron algorithm applies the rule iteratively on the misclassified points

for every misclassified $\tilde{x}^{(i)}$ do

$$| \qquad \beta \leftarrow \beta + \gamma \tilde{x}^{(i)} t^{(i)}$$

