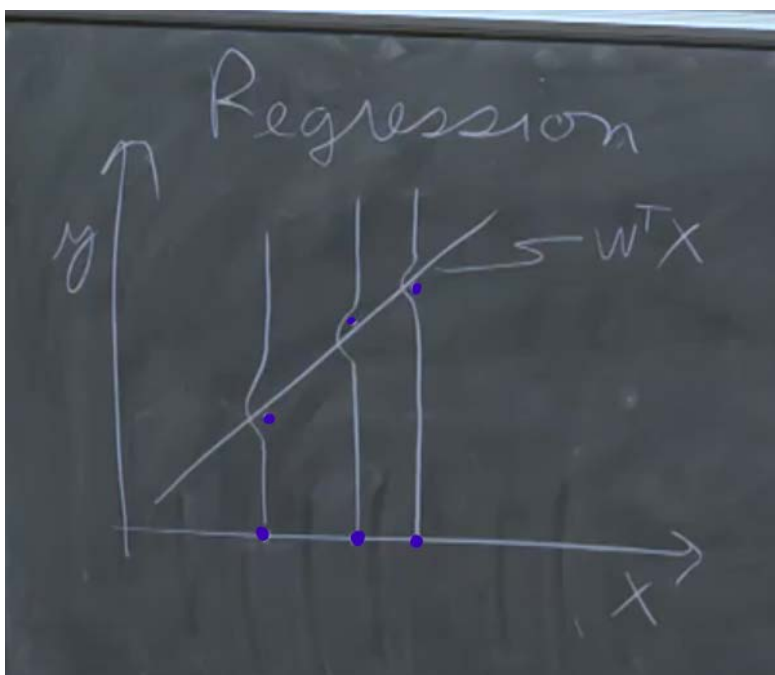
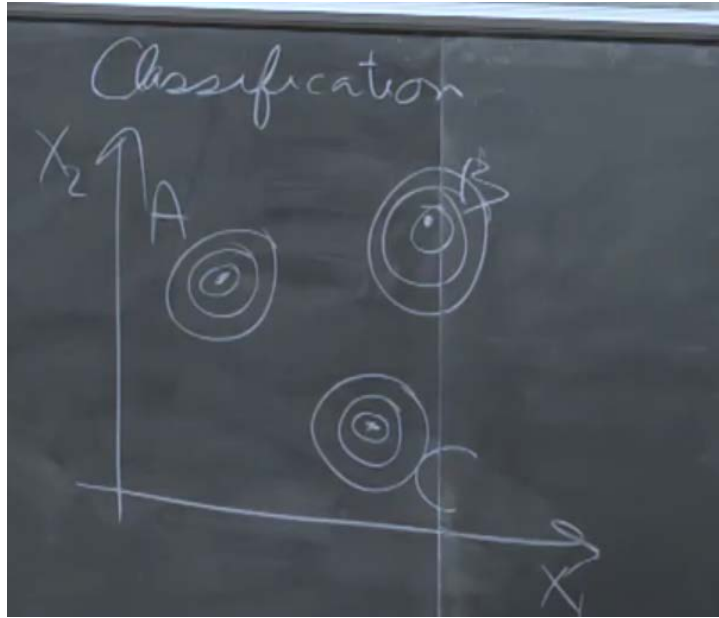


بسم الله الرحمن الرحيم

محاضرة ال Mixture of Gaussians .. ده برضو كلاسيكيشن يعني ... بس باستخدام الجاوسيان مكشيز ... وده model في family اسمها Generative models .. احنا اشفنا بالنسبة لل regression شفنا statistical linear regression model .. دلوقت هنعمل حاجه زي كذا لل classification .. تعال الاول نشوف يعني ايه Generative models .. في ال regression احنا كان هدفنا نقلل الدستانس بالنسبة للخط .. عندك شوية داتا وبتعمل فيتنج كيرف .. ومن ناحية الأوبتيمايزيشن .. هدفك تقلل المسافات .. في الريجرشن .. القصه كلها ان ممكن ن suspect that the underlying function is linear و هي ديه الترو و ديه اسمها $w^T x$.. و الإنبوتس هي الإكس أكسيس و الأوتبوتس هي الواي أكسيس .. لو الفانكشن كانت صحيحه .. يبقى الخط هيعدي علي النقط كلها . بس الحقيقه ان الداتا اللي احنا لاحظناها اتقاست من شوية سنسورز و السنسورز مش جامده جداً يعني ويتضيف للنويز .. فالنويز ديه ممكن يبقى ليها جاوسين ديستريوشن .. فلو عندك شوية داتا بوينت إكس .. حتي لو انا عارف ال true underlying function .. اللي انا would likely observe هو ان الاوتبوت sampled from a Gaussian distribution .. فاللي اترسم هو Gaussian علي ال vertical line و فعلياً انت تقدر تجيب measurement at the line ... فانا لو عاوز أجيرويت شوية داتا .. sample from the Gaussian .. وهجيب نقطه .. مش هتبقا علي نفس الخط .. بس هتبقا قريبه من الخط ... فهنا نقدر نقول ان ده generative model عشان انا اقدر استخدم الموديل عشان ا simulate ال creation of data set .. طب ايه الهدف من حاجه زي كذا .. في أليكشنز كثيره هدفها حاجه زي كذا انها ت generate data .. وهنا بتحتاج generative model .. فهم بقدر يجينيريتوا داتا زي اللي عندك في التريننج سيت .. بس ففي الصورة عندنا هي بس حاجه بسيطه بتوضح الفكره يعني



في حالة الكلاسيكيشن بقا .. عندك مثلاً 3 كلاسيز .. A, B and C ... في الحاله ديه .. استخدم أول أكسيس عشان ت represent x_1 اللي هو أول كومبوننت .. و الإكس 2 علي الفيرتيكال لاين .. خد بالك ان الأكسيس مختلفه .. هنا الأكسيس انبوت انبوت .. الافتراضيه اللي عندي ان لو عندي كلاسيكيشن .. لما بقيس الإنبوت هو مش هيبقا accurate .. بس بعد كذا انا هتوقع ان كل داتا بوينت في كلاس سي مثلاً .. centered عند النقطه بتاعت C .. بس ممكن يكون في Gaussian distribution اللي ممكن ي correspond to noise in the measurements of the input .. وبعدين كنتيجته لكدا عندك داتا بت follow Gaussian distribution centered بالطريقه اللي الدكتور راسمها كذا ..



بناءً على فكره ان الداتا بتاعتنا بتتبع في مناطق مختلفه .. ولكل كلاس في ديستريبيوشن .. انت تقدر تبني a probabilistic model .. اللي هو تقدر تقول فيه .. لما مكنتش ببص على أي داتا بوينت انا مجرد بخمن ايه هو الكلاس بتاع الداتا بوينت اللي جايلي ديه فساعتها ببيقا عندي prior distribution ونسميه probability of c و ال c هنا هي الكلاس يعني بس ديه prior يعني قبل ما نبص على أي داتا بوينت .. أول ما هيكون عندك داتا بوينت .. ايه احتمالية انك تختار نقطه من الكلاس C فبيقا عندك class conditional distribution اقدر استخدمو عشان أ generate points from that class ..

انت كذا تقدر تحسب ال Posterior باستخدام Bayes theorem .. وده ال prediction for the class of this data point .. في السلايد ديه لازم تاخذ بالك .. ان ده مش Bayesian learning احنا بنستخدم bayes theorem عشان ن infer ون estimate ايه هي البروبابلي بتاعت الكلاس ل given data point .. انما لما كنا بنتكلم على Bayesian learning .. ده كان في ال context بتاع اننا عندنا شوية بارمترز بتاعت الموديل وبعدين بحاول اني أ estimate البارمترز بتاعت الموديل باستخدام bayes theorem ... فهنا احنا متكلمناش خالص على البارمترز ... كل اللي عملناه اننا اتكلمنا على ال inference عشان تحسب ال posterior probability of a class ..

Probabilistic Generative Model

- $\Pr(C)$: prior probability of class C
- $\Pr(x|C)$: class conditional distribution of x
- Classification: compute posterior $\Pr(C|x)$ according to Bayes' theorem

$$\Pr(C|x) = \frac{\Pr(x|C) \Pr(C)}{\sum_C \Pr(x|C) \Pr(C)}$$

$$= k \Pr(x|C) \Pr(C)$$

طيب في الموديل ده هنحط شوية assumptions .. عندنا finite number of classes وهي حاجه categorical .. فيالتالي ال natural type of prior نقدر نستخدمو هو ال multinomial distribution .. وده distribution over finite number of possible outcomes that are categorical .. زي مثلا .. عندك coin اقلب الكوين .. تلاقي head or tail .. ده هيديك bernolli distribution .. انما لو عملت generalization هتلاقي عندك dice .. فتقدر تستخدم ال multinomial distribution ... هنا هنستخدم الرمز π_k ... عشان k ال class ال probability of class k .. وده مجرد رقم ما بين الصفر والواحد ..

هنفترض كمان ان عندنا داتا .. والداتا ديه ليها d dimensions .. فدلوقت انت عاوز بيقا عندك a class conditional distribution : بقلنا :

For a given class, where the data is likely to be?

وهنا هنحط الافتراض ان ال class conditional distribution is a Gaussian distribution .. تمام .. فانا عندي الكلاس بتاعي .. $P(x|Class\ c)$

وده اللي احنا رسمناه علي البورد .. احنا كمان هنفترض ان كل Gaussian distribution هبيقا ليه نفس ال covariance matrix اللي هي capital sigma .. تحت الفرضيه ديه .. هنقول العلاقه اللي مكتوبه ديه .. والافتراض ده مهم جداً .. خليك فاكراه ..

Assumptions

- In classification, the number of classes is finite, so a natural prior $\Pr(C)$ is the multinomial

$$\Pr(C = c_k) = \pi_k$$

- When $x \in \mathbb{R}^d$, then it is often OK to assume that $\Pr(x|C)$ is Gaussian.
- Furthermore, assume that the same covariance matrix Σ is used for each class.

$$\Pr(x|c_k) \propto e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}$$

دلوقت بعد ما عرفنا ال prior و ال likelihood .. فال likelihood is a Gaussian distribution .. نقدر نحسب ال inference to estimate the property of each class given a data ... Bayesian learning ... وده مش posterior to do inference .. point

Posterior Distribution

$$\begin{aligned} \Pr(c_k|x) &= \frac{\pi_k e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}}{\sum_k \pi_k e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}} \\ &= \frac{\pi_k e^{-\frac{1}{2}(x^T \Sigma^{-1} x - 2\mu_k^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k)}}{\sum_k \pi_k e^{-\frac{1}{2}(x^T \Sigma^{-1} x - 2\mu_k^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k)}} \end{aligned}$$

Consider two classes c_k and c_j

$$= \frac{1}{\pi_j e^{\mu_j^T \Sigma^{-1} x - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j} + \pi_k e^{\mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k}}$$

Posterior Distribution

$$= \frac{1}{1+e^{-\left(\mu_k^T - \mu_j^T\right) \Sigma^{-1} x + \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j - \ln \frac{\pi_k}{\pi_j}}}$$

$$= \frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x} + w_0)}}$$

where $\mathbf{w} = \Sigma^{-1}(\mu_k - \mu_j)$

and $w_0 = -\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j + \ln \frac{\pi_k}{\pi_j}$

كما بقي عندنا ال posterior في فورم كويسه .. هل حد عارف ايه الفانكشن ديه .. الفانكشن ديه مستخدمه كتير في النويرال نتورك .. هوب ديه السجمويد فانكشن ..

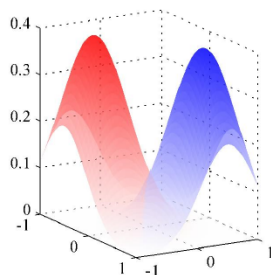
Logistic Sigmoid

- Let $\sigma(a) = \frac{1}{1+e^{-a}}$
└───────────> Logistic sigmoid
- Then $\Pr(c_k|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$
- Picture:

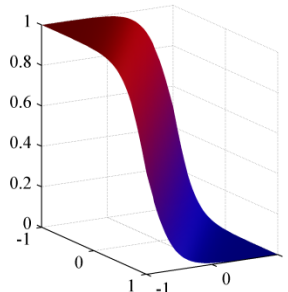
الفانكشن ديه مهمه ... هنرسمها ..

Logistic Sigmoid

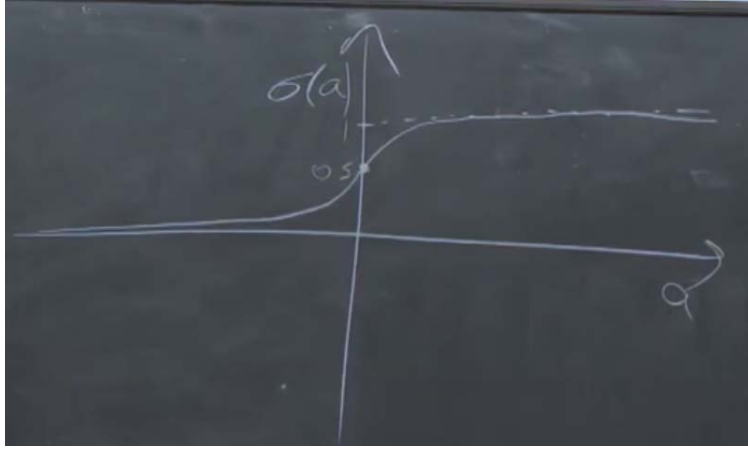
class conditionals



posterior



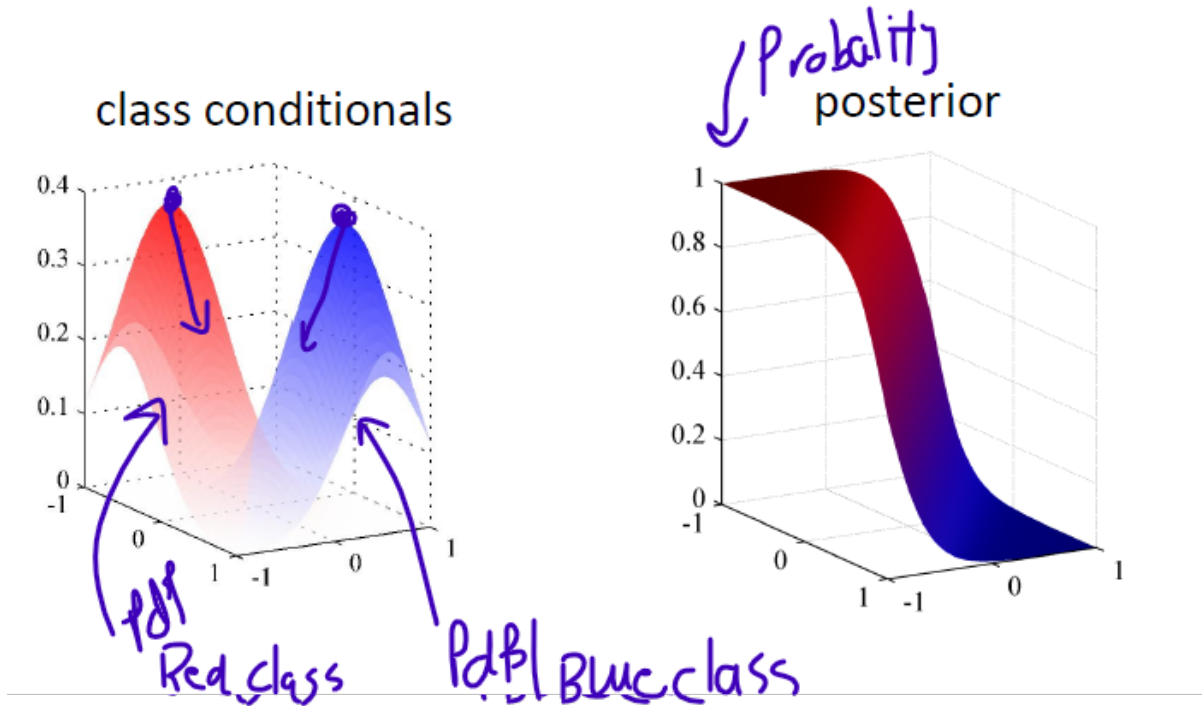
الإنيتوت هو رييل نامير .. والأوتيتوت ما بين الصفر والواحد



وعشان كذا هي مشهوره في ال Neural networks .. هنا بتحشر الأوتيتوت ما بين الصفر والواحد عشان تقدر تتخيلو علي انو probability ..

طبيب تعال نشوف الصورة اللي في لسللايز الي فيها الرسومات الملونه .. انت عندك اول كيرف في كيرف أحمر .. ده أول كلاس .. و ده ال Gaussian distribution للداتا اللي جايه منو .. وفي كلاس ثاني هو الازرق ده ... لو هتيجي تحسب ال posterior .. هتلاقي ان ليها احتمالية اعلي للون الأحمر مثلاً يعني ..

Logistic Sigmoid



دلوقت احنا نقدر نعمل prediction .. لما بنحسب ال posterior بيقا عندنا distribution over the classes و بالتالي اقدر اعمل بريدكشن .. في شوية ألكيشنز احنا مش بنعوز نرجع ايه هي ال probability for every class بس اللي احنا عاوزينو ... اهو هو ده الكلاس اللي انا قدرت اعمله prediction حتي لو انا مش واثق ف نفسي .. وحالياً لو انت عاوز تختار كلاس ... الطبيعي انك تقول اختار الكلاس اللي ليها بروبابلتي اعلي من ال 0.5 ..

طبيب دلوقت احنا لما هنعمل predictions زي كذا .. احنا بنحتاج نعرف ال boundary ما بين ال 2 كلاسيز عند ال $K = 1$.. عندك الدايجرام اللي اسمو voronoi diagram اللي بيديك الباونديرز لما بتستخدم إقليدين ديستنس .. لما بتستخدم ال mixture of Gaussian .. الكلاس باونديري .. هو ان الباونديري محدد لما بتكون احتمالية كل كلاس قد بعض .. فهتبقا عند ال 0.5 .. بعد شوية simplifications .. بنوصل ان عندنا linear separator

Prediction

$$\begin{aligned} \text{best class} &= \operatorname{argmax}_k \Pr(c_k | \mathbf{x}) \\ &= \begin{cases} c_1 & \sigma(\mathbf{w}^T \mathbf{x} + w_0) \geq 0.5 \\ c_2 & \text{otherwise} \end{cases} \end{aligned}$$

$$\text{Class boundary: } \sigma(\mathbf{w}_k^T \bar{\mathbf{x}}) = 0.5$$

$$\Rightarrow \frac{1}{1 + e^{-(\mathbf{w}_k^T \bar{\mathbf{x}})}} = 0.5$$

$$\Rightarrow \mathbf{w}_k^T \bar{\mathbf{x}} = 0$$

\therefore linear separator

لما بتكون السجمويد ب 0.5 .. لو بصيت علي الإكس أكسيس .. هنا عندها الإنبوت بتاع السجمويد ببيقا بصفر .. وبكدا انت فهمت انت ازاي وصلت للإكسبرشن اللي هو $\mathbf{w}^T \mathbf{x} = 0$.. عشان ال activation = 0 وده linear separator ... عشان هو linear expression ... طيب لو عندنا أكثر من 2 كلاسز .. هنعمل نفس الحسابات .. وهنفترض ان كل كلاس لييه class conditional distribution اللي هو Gaussian و ليهم نفس ال covariance matrix اللي هي Capital sigma ... فيدل ما بنقسم البسط علي المقام هنسيب البسط والمقام زي ما هم .. ونحاول ن simplify طريقه الكتابه نفسها .. فديه هي ال softmax ..

Multi-class Problems

- Consider Gaussian conditional distributions with identical Σ

$$\begin{aligned} \Pr(c_k | \mathbf{x}) &= \frac{\Pr(c_k) \Pr(\mathbf{x} | c_k)}{\sum_j \Pr(c_j) \Pr(\mathbf{x} | c_j)} \\ &= \frac{\pi_k e^{-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k)}}{\sum_j \pi_j e^{-\frac{1}{2}(\mathbf{x} - \mu_j)^T \Sigma^{-1} (\mathbf{x} - \mu_j)}} \\ &= \frac{\pi_k e^{-\frac{1}{2}(-2\mu_k^T \Sigma^{-1} \mathbf{x} + \mu_k^T \Sigma^{-1} \mu_k)}}{\sum_j \pi_j e^{-\frac{1}{2}(-2\mu_j^T \Sigma^{-1} \mathbf{x} + \mu_j^T \Sigma^{-1} \mu_j)}} \\ &= \frac{e^{\mu_k^T \Sigma^{-1} \mathbf{x} - \frac{1}{2}\mu_k^T \Sigma^{-1} \mu_k + \ln \pi_k}}{\sum_j e^{\mu_j^T \Sigma^{-1} \mathbf{x} - \frac{1}{2}\mu_j^T \Sigma^{-1} \mu_j + \ln \pi_j}} = \frac{e^{\mathbf{w}_k^T \bar{\mathbf{x}}}}{\sum_j e^{\mathbf{w}_j^T \bar{\mathbf{x}}}} \Rightarrow \text{softmax} \\ \text{where } \mathbf{w}_k^T &= \left(-\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \ln \pi_k, \mu_k^T \Sigma^{-1} \right) \end{aligned}$$

السوفت ماكس هو الجينيريللايزيشن بتاع السجمويد ... طب هو ليه اسمو softmax ...

Softmax

- When there are several classes, the posterior is a **softmax** (generalization of the sigmoid)

- Softmax distribution: $\Pr(c_k|\mathbf{x}) = \frac{e^{f_k(\mathbf{x})}}{\sum_j e^{f_j(\mathbf{x})}}$

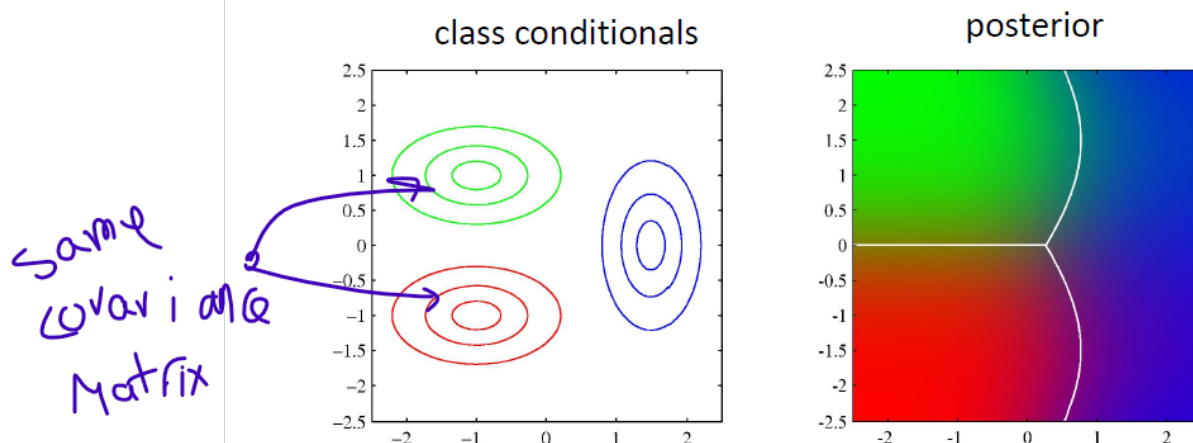
- Argmax distribution:

$$\Pr(c_k|\mathbf{x}) = \begin{cases} 1 & \text{if } k = \operatorname{argmax}_j f_j(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases}$$
$$= \lim_{\text{base} \rightarrow \infty} \frac{\text{base}^{f_k(\mathbf{x})}}{\sum_j \text{base}^{f_j(\mathbf{x})}}$$
$$\approx \frac{e^{f_k(\mathbf{x})}}{\sum_j e^{f_j(\mathbf{x})}} \quad (\text{softmax approximation})$$

تعال نقول انا بحسب posterior و عاوز ا associate probability 1 for the class that has the heighest function and zero otherwise

فهو هيدي واحد للي ليه اعلي قيمه .. وصفر للباقي .. ومن هنا جيه الماكسمم .. انو بياخد ال argmax .. لو عاوز ا soften الكلام ده .. من هنا هو ال softmax هو عبارته عن approximation .. لانو ممكن تكتب ال argmax بطريقه مختلفه يعني ... روح للمالانهايه .. و خذ البيز و حط الفانكشن لكل كلاس ك exponent و اقسم علي ال base مرفوع للفانكشن بتاعت كل كلاس و اجمعهم كلهم فلما البيز تروح للإنتي .. هتقسم علي 1 .. وبتاخذ البسط بس ... فهنا انت عاوز بببيز فاينايته .. فعشان كذا البيز هنا هو ال exponent

Softmax



المهم الدكتور بيقولك ان الشكل بتاع الجاوسين بيتحدد عن طريق ال covariance matrix ... الباوندري هنا شكلو زي الخوخه .. الباوندري linear separator ما بين الاثنين اللي ليهم نفس الجاوسين covariance matrix ... انما لما مكنتش في نفس الكوفاريانس ماتركس .. انت بقا عندك non-linear ... boundary

الخطوه الجايه ... احنا شغنا ازاى نعمل ال inference .. هنحسب ال posterior .. بس البارمترز بتاعت الموديل بتيجي منين .. فهدفنا هنا هو البارمترز استميشن

Parameter Estimation

- Where do $\Pr(c_k)$ and $\Pr(x|c_k)$ come from?
- Parameters: $\pi, \mu_1, \mu_2, \Sigma$

$$\begin{aligned} \Pr(c_1) &= \pi, & \Pr(x|c_1) &= k_{\Sigma} e^{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)} \\ \Pr(c_2) &= 1 - \pi, & \Pr(x|c_2) &= k_{\Sigma} e^{-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2)} \end{aligned}$$

where k_{Σ} is the normalization constant that depends on Σ

- Estimate parameters by
 - Maximum likelihood
 - Maximum a posteriori
 - Bayesian learning

بالنسبه للرننج بقا نقدر نستخدم ال Max Likelihood ... هنا هدفنا نلاقي البارمترز بتاعت الموديل اللي تماكسيمائز ال likelihood بتاع الداتا .. فاحنا هنبدا بال likelihood فانكشن بتاعت الداتا سبت ...

Maximum Likelihood Solution

- Likelihood:

$$L(\mathbf{X}, \mathbf{y}) = \Pr(\mathbf{X}, \mathbf{y} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_n [\pi N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{y_n} [(1 - \pi) N(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-y_n}$$

class #1 ← class #2
عنا البرهان نبتأ
أشهر

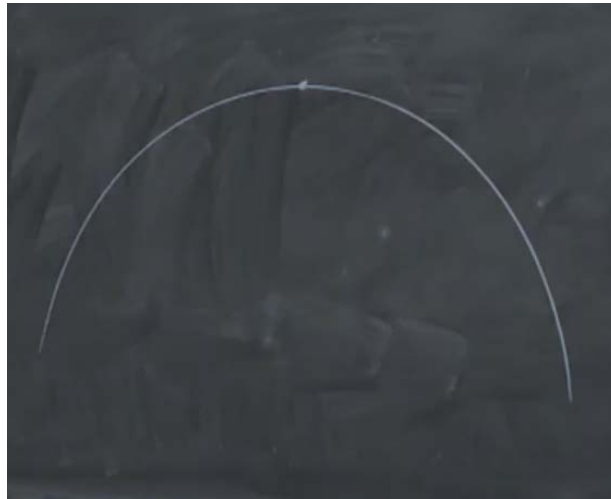
- ML hypothesis:

$$\langle \pi^*, \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*, \boldsymbol{\Sigma}^* \rangle =$$

$$\underset{\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}}{\operatorname{argmax}} \sum_n y_n \left[\ln \pi + \ln k_{\boldsymbol{\Sigma}} - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \right] + (1 - y_n) \left[\ln(1 - \pi) + \ln k_{\boldsymbol{\Sigma}} - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \right]$$

optimization problem

الأوبجكتيف هنا بتاع ال maximizing log-likelihood بي correspond to a curve ليه الشكل اللي الدكتور رسمه ده .. حاجة concave function التوب بتاعها هو الماكسيمم .. فاشتق وساوي بالصفر وهيجبك النقطة اللي انت عاوزها للماكسيم



Maximum Likelihood Solution

- Set derivative to 0

$$0 = \frac{\partial \ln L(X, y)}{\partial \pi}$$

$$\Rightarrow 0 = \sum_n y_n \left[\frac{1}{\pi} \right] + (1 - y_n) \left[-\frac{1}{1 - \pi} \right]$$

$$\Rightarrow 0 = \sum_n y_n (1 - \pi) + (1 - y_n)(-\pi)$$

$$\Rightarrow \sum_n y_n = \pi (\sum_n y_n + \sum_n (1 - y_n))$$

$$\Rightarrow \sum_n y_n = \pi N \quad (\text{where } N \text{ is the \# of training points})$$

$\therefore \frac{\sum_n y_n}{N} = \pi$

sum data points belong to class \rightarrow π \leftarrow *This is sample mean (Empirical mean)*

الإكسبرشن الأخير ده .. هو انك تجمع الداتا بوينتس اللي نفس الكلاس بتاعك و تقسم علي عددهم ...

Maximum Likelihood Solution

$\mu_1 N_1 \rightarrow 0 = \frac{\partial \ln L(X, y)}{\partial \mu_1}$

$$\Rightarrow 0 = \sum_n y_n [-\Sigma^{-1} (x_n - \mu_1)]$$

$$\Rightarrow \sum_n y_n x_n = \sum_n y_n \mu_1$$

$$\Rightarrow \sum_n y_n x_n = N_1 \mu_1$$

$\therefore \frac{\sum_n y_n x_n}{N_1} = \mu_1$

Empirical Mean

Similarly: $\frac{\sum_n (1 - y_n) x_n}{N_2} = \mu_2$

where N_1 is the # of data points in class 1

N_2 is the # of data points in class 2