

لما بنيجي نتكلم علي ال linear regression .. in terms of statistical learning .. القصه كلها ان لما بيكون عندنا شوية داتا .. هنقول ان الداتا جايه من شوية قياسات ممكن يبقا فيها نويز .. فانت عندك إنبوت  $x$  .. و أوتبوت  $y$  .. فانت مش هتلاحظ الأوتبوت بتاع ال function  $f$  .. وديه ال function اللي انت عاوز تتعلمها . انت هت observe ... a corrupted version of  $y$  .. اللي هي  $f(x) + \text{noise}$  .. زي مثلاً ل weather prediction .. وانت عندك historical data زي مثلاً wind pressure temperature .. و كل يوم عندك ال measured temperature .. بس ال thermometer مش accurate .. بيحطلك noise ... فيشكل عام لأي نوع من انواع الداتا لازم تفترض ان هيبقا في corruption نتيجة وجود ال noise .. فانت محتاج تاخدها في الاعتبار .. لما بتيجي ت model الكلام ده .. واحنا بنتكلم علي ال linear regression .. هنفترض ان ال function  $f$  هي linear function .. و النويز هنفترض انها Gaussian noise . انت ممكن تحط noise تحت أي distribution .. بس Gaussian هو اللي منتشر يعني .. ده كذا ال setting اللي هنتكلم عليه ...

فكدا انت تقدر تعبر عن likelihood distribution بقولنا .. لكل إنبوت في الداتا سبت بتاعتنا .. ايه هو ال likelihood اللي هيقس أوتبوت معين اللي هو  $y$  .. فالدكتور بيقول ان ال likelihood distribution . نقدر ن inferه من الموديل اللي عندنا

## Noisy Linear Regression

- Assume  $y$  is obtained from  $x$  by a deterministic function  $f$  that has been perturbed (i.e., noisy measurement)

Model  $\rightarrow y = f(\bar{x}) + \epsilon$

$\downarrow$                        $\downarrow$   
 $w^T \bar{x}$                $N(0, \sigma^2)$

- Gaussian noise:

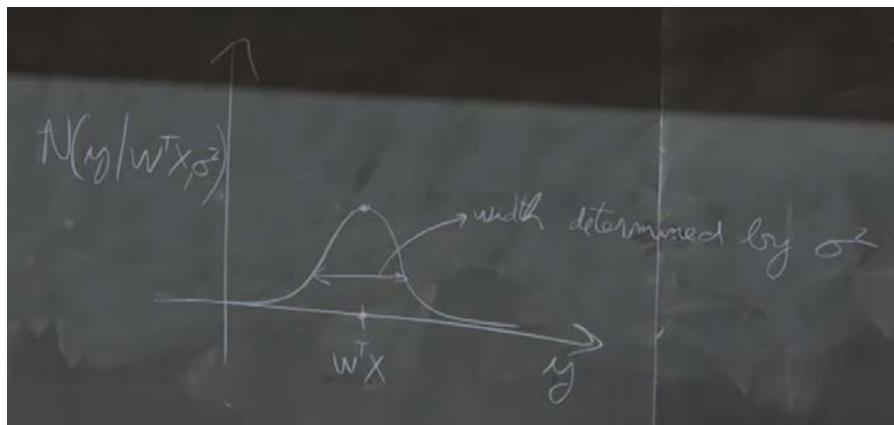
$$\Pr(y|\bar{X}, w, \sigma) = N(y|w^T \bar{X}, \sigma^2)$$

$\swarrow$                        $\nwarrow$   
 output                      input

$$= \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_n - w^T \bar{x}_n)^2}{2\sigma^2}}$$

الدكتور شرح هو ليه ال Probability هتطلع Gaussian .. عشان انت فترض ان ال function هتبقا linear .. و كمان it's deterministic .. فبكدا انت هتضيف نويز ليها جواسين ديستريبيوشن .. فالواي هيبقا ليها فاريانس زي ال noise .. و القيمه اللي هنشوفها هتبقا هي ال  $w^T x$  .. لو جيت تكتب ال pdf .. هتبقا هي ال  $\Pr(y|X, w, \sigma)$  ... تعال نرسم اللي بيحصل ده عشان يبقا عندنا mental image .. مش mental sickness ...

هنرسم curve .. الإكس أكسيز هو ال output  $y$  .. عشان هنرسم ال pdf لل Gaussian distribution بالنسبة ل  $y$  ... ال max هيبقا عند ال mean و ال variance هيحددك ال width بتاع الكيرف ...



دلوقت احنا هنعوز ن apply the maximum likelihood .. اللي شغناه المحاضره اللي فاتت .. هناخد ال expression و نقول عاوزين  $W$  اللي هتجيب ال max probabiolity لكل ال  $y$  في الداتا سبت .. انت عندك داتا .. فيها الإنبوت و الواي .. دلوقت عندك موديل بيتحدد بال  $W$  . وال  $W$  بيقلك ازاي تقدر تاخذ linear combinations من الإنبوتس اللي عندك عشان تعرف ت predict the output .. ودلوقت انت عاوز تختار أحسن  $W$  .. عشان تعمل كذا هتقول مين هي أحسن  $W$  اللي هتديك اعلي احتمالية لل  $y$  اللي احنا عاوزنها ..

## Maximum Likelihood

- Possible objective: find best  $w^*$  by maximizing the likelihood of the data

$$w^* = \operatorname{argmax}_w \Pr(y|\bar{X}, w, \sigma)$$

$$= \operatorname{argmax}_w \prod_n e^{-\frac{(y_n - w^T \bar{x}_n)^2}{2\sigma^2}}$$

$$= \operatorname{argmax}_w \sum_n -\frac{(y_n - w^T \bar{x}_n)^2}{2\sigma^2}$$

$$= \operatorname{argmin}_w \sum_n (y_n - w^T \bar{x}_n)^2$$

log  
مش جيبنا  
ال output  
انما صيغ  
Rescaling

- We arrive at the original least square problem!

هو ليه اصلا بناخد ال log .. عشان احنا عندنا ضرب و هنا ال expression is complicated و ده expression non convex ... لما بيبقي convex ببقا سهل ان الواحد يتعامل معاه انما لما بيبقي non convex مش ببقا سهل تتعامل معاه .. فبالتالي بنحتاج ال log عشان كذا .. وديه هي ال least square problem .. linear regression عنها في ال

طيب ايه الي هيحصل لو قلنا بدل ما نجيب ال  $W$  اللي بت maximize ال likelihood .. نجيب ال  $W$  اللي بت maximize the posterior .. عن طريق اننا نحسب ال posterior ...

## Maximum A Posteriori

- Alternative objective: find  $w^*$  with highest posterior probability
- Consider Gaussian prior:  $\Pr(w) = N(0, \Sigma)$

- Posterior:  $\Pr(w|X, y) \propto \underbrace{\Pr(w)}_{\text{prior}} \underbrace{\Pr(y|X, w)}_{\text{Likelihood}}$

$$= k e^{-\frac{w^T \Sigma^{-1} w}{2}} e^{-\frac{\sum_n (y_n - w^T x_n)^2}{2\sigma^2}}$$

الدكتور بعد كذا اتكلم علي وجود ال prior distribution .. ال prior انت محتاج ت express ال distribution بتاعه .. لو فكرت في ال Machine learning كأنك مش عارف ايه هي ال hypothesis .. انا عندي شوية Uncertainty في البدايه .. و ال process بتاعت ال machine learning .. هي ال process of reducing the uncertainty وبعدين نقدر ن quantify ال uncertainty بتاعتنا عن طريق ال distribution اللي هو ال prior .. وعملية ال learning هي انك تحسب ال posterior distribution .. و هي اسمها posterior عشان ال distribution ده بيطلع بعد ما نكون بصينا علي الداتا الإكس و الواي .. وده هيدينا الديرستريوشن اللي ليه less uncertainty وده هو ال bayesian learning بيتكلم عليه .. فبالنالي انا محتاج ال prior distribution اللي بيحدد ال initial uncertainty لل W ... هنا احنا هنستخدم ال domain knowledge .... فانا هخمن فين ال W و اقوم حاطط Gaussian هناك .. فهنستخدم جالوسين mean zero and covariance sigma

## Maximum A Posteriori

- Alternative objective: find  $w^*$  with highest posterior probability
- Consider Gaussian prior:  $\Pr(w) = N(0, \Sigma)$

- Posterior:  $\Pr(w|X, y) \propto \Pr(w) \Pr(y|X, w)$

$$= k e^{-\frac{w^T \Sigma^{-1} w}{2}} e^{-\frac{\sum_n (y_n - w^T x_n)^2}{2\sigma^2}}$$

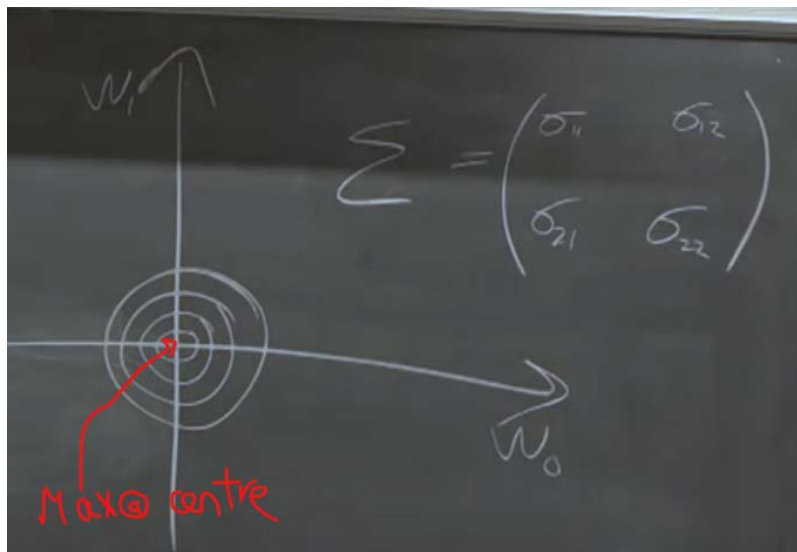
$w \rightarrow$  a vector of weights

University of Waterloo

CS480/680 Spring 2019 Pascal Poupart

4

دلوقت هنرسم مثال لل prior distribution لل W .. وهنقول يعني ان ال W ليها W0 و W1 .. ممكن يكون عندك بعد كذا أكثر من W .. بس يعني خيلنا نتكلم علي الاتنين دول دلوقت .. هنا هنرسم ال Gaussian distribution .. فهيبقا surface وكونو طالع من الصفحه .. فهيبقا عندك contour lines .. ال peak عند النص و كل contour بيغير عن spherical Gaussian .. وال Width بيتحدد عن طريق ال covariance matrix ... احنا عندنا ال covariance matrix Segma .. ال Segma هيبقا ليه 4 ارقام ... sigma11, sigma12, sigma21, sigma22 .. علي ال diagonal هيبقا عندك ال variance .. ده اللي هو زي لما اتكلمنا علي 1D .. و ال off diagonal .. هيبقا ال covariance .. وده اللي بيقلنا ازاى ال different weights بت relate لبعض .. فهي ديه الفكرة ..



دلوقت معانا خلاص ال expression بتاع ال posterior .. هنعوز نجيبو ال maximum ..

# Maximum A Posteriori

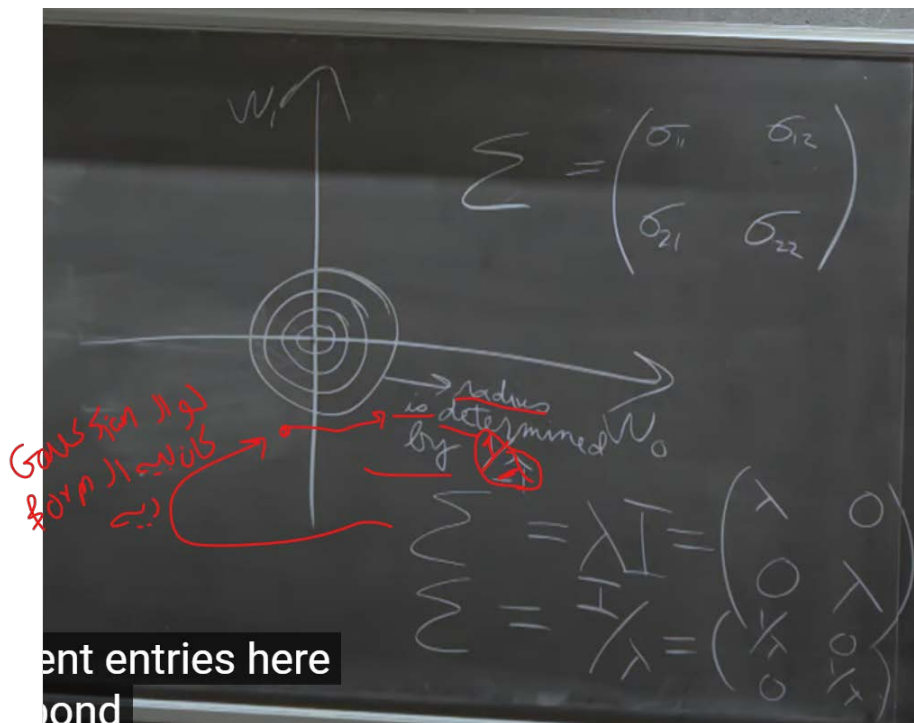
- Optimization:

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w}} \Pr(\mathbf{w} | \bar{\mathbf{X}}, \mathbf{y}) \\ &= \operatorname{argmax}_{\mathbf{w}} - \sum_n (y_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2 - \mathbf{w}^T \Sigma^{-1} \mathbf{w} \\ &= \operatorname{argmin}_{\mathbf{w}} \sum_n (y_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2 + \mathbf{w}^T \Sigma^{-1} \mathbf{w} \end{aligned}$$

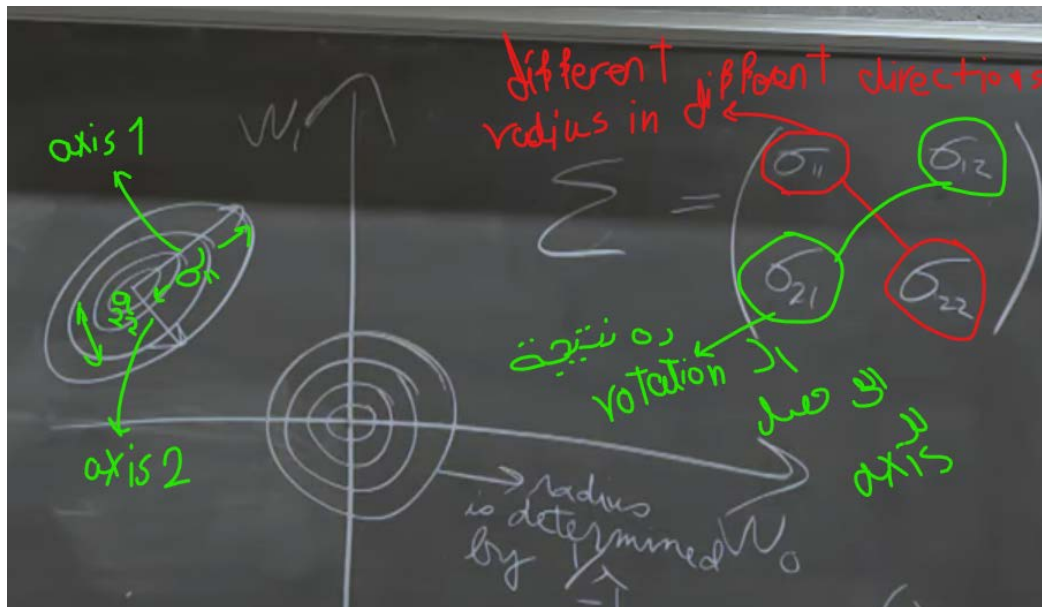
- Let  $\Sigma^{-1} = \lambda \mathbf{I}$  then *Identity Matrix*  

$$= \operatorname{argmin}_{\mathbf{w}} \sum_n (y_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$
- We arrive at the original regularized least square problem!

دلوقت انت عندك طريقه مختلفه لترجمة ال regularization term .. المحاضره اللي فاتت قلنا ان ال solution is stable عن طريق انك تحاول ت force ال entries of W انها تبقي as small as possible في هي لو قلت اني عاوز اقلل ال Norm of W .. اللي هو انا برضو عاوز أخلي ال entries بيقو قريبين من بعض ومن الصفر .. بس مكنش واضح ليه عاوزنها قريبه للصفر . ليه مثلاً مش 10 او 100 .. دده عشان ان ده بيحي نتيجة ان ال prior distribution بيديك ال knowledge اللي هو بنيقا عارفين فيه ان ال solution اللي هو أحسن W .. لازم تنقا قريبه للصفر .. وبعدين تنقا في منطقه فيها bell-shaped distribution .. فلو عندك ال prior distribution بتقولك ان الحلول القريبه من الصفر هي الأحسن ..



لو كان ال Gaussian ليه ال form التانيه اللي هي السجما الكامله



تعال نشوف ازاى ال 2 approaches بيدولنا losses مختلفه .. فدلوقت هن analysis ال loss function .. وازاي الترمات بتختلف ... بيقا انت عندك different techniques بيدولك نفس ال  $W$  .. بس ال solution بيقا impacted by our choice of  $\lambda$  لكل لمضه مختلفه هيطلعنا حلول مختلفه . لو انت given  $W$  . هنشوف ايه هو ال expected loss واحنا بنغير اللمضه ازاى ال expected loss هتتغير ... تعال نشوف ال derivations ال الي بتحصل ..

## Expected Squared Loss

- Even though we use a statistical framework, it is interesting to evaluate the expected squared loss

data  $\leftarrow$   $E[L] = \int_{x,y} \text{Pr}(x,y) (y - \mathbf{w}^T \bar{\mathbf{x}})^2 dx dy$   $\rightarrow$  Euclidean distance

from a distribution  $= \int_{x,y} \text{Pr}(x,y) (y - f(x) + f(x) - \mathbf{w}^T \bar{\mathbf{x}})^2 dx dy$

$$= \int_{x,y} \text{Pr}(x,y) \left[ (y - f(x))^2 + 2(y - f(x))(f(x) - \mathbf{w}^T \bar{\mathbf{x}}) + (f(x) - \mathbf{w}^T \bar{\mathbf{x}})^2 \right] dx dy$$

Expectation with respect to  $y$  is 0

$$E[L] = \underbrace{\int_{x,y} \text{Pr}(x,y) (y - f(x))^2 dx dy}_{\text{noise (constant)}} + \underbrace{\int_x \text{Pr}(x) (f(x) - \mathbf{w}^T \bar{\mathbf{x}})^2 dx}_{\text{error (depends on } \mathbf{w} \text{)}}$$

exp and this term

# Bias-Variance Decomposition

- Decompose squared loss

$$\begin{aligned}
 & E_S[(f(\mathbf{x}) - \mathbf{w}_S^T \bar{\mathbf{x}})^2] \\
 &= E_S[f(\mathbf{x}) - E_S[\mathbf{w}_S^T \bar{\mathbf{x}}] + E_S[\mathbf{w}_S^T \bar{\mathbf{x}}] - \mathbf{w}_S^T \bar{\mathbf{x}}]^2 \\
 &= E_S[(f(\mathbf{x}) - E_S[\mathbf{w}_S^T \bar{\mathbf{x}}])^2] \\
 &\quad + 2(f(\mathbf{x}) - E_S[\mathbf{w}_S^T \bar{\mathbf{x}}])(E_S[\mathbf{w}_S^T \bar{\mathbf{x}}] - \mathbf{w}_S^T \bar{\mathbf{x}}) \\
 &\quad + (E_S[\mathbf{w}_S^T \bar{\mathbf{x}}] - \mathbf{w}_S^T \bar{\mathbf{x}})^2 \quad \text{Expectation is 0} \\
 &= \underbrace{(f(\mathbf{x}) - E_S[\mathbf{w}_S^T \bar{\mathbf{x}}])^2}_{\text{bias}^2} + E_S[\underbrace{(E_S[\mathbf{w}_S^T \bar{\mathbf{x}}] - \mathbf{w}_S^T \bar{\mathbf{x}})^2}_{\text{variance}}]
 \end{aligned}$$

University of Waterloo

CS480/680 Spring 2019 Pascal Poupart

8

## Bias-Variance Decomposition

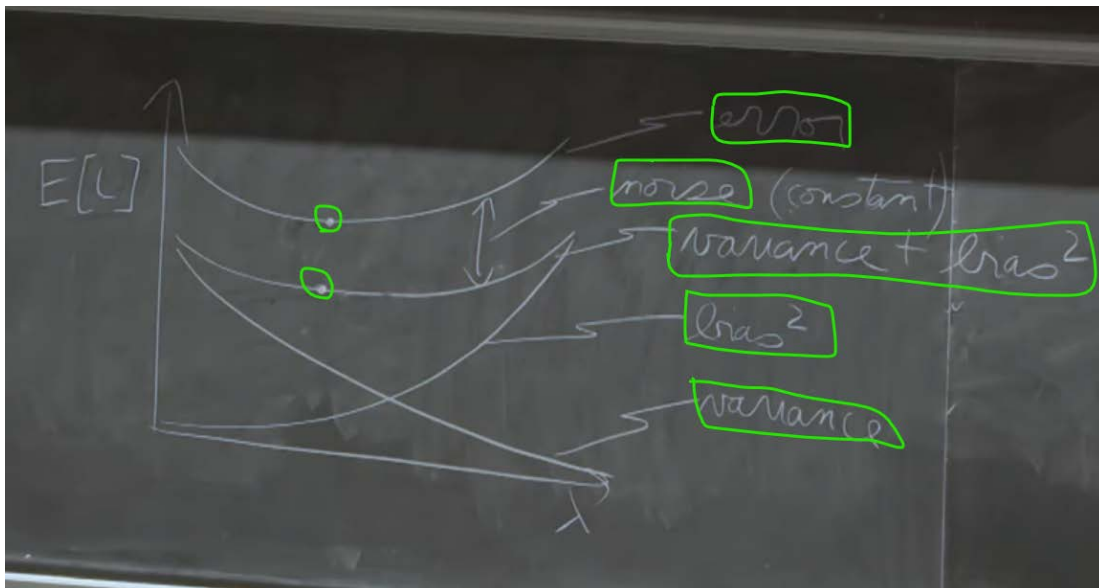
- Hence:

$$E[\text{loss}] = (\text{bias})^2 + \text{variance} + \text{noise}$$

- Picture:

تعال نرسم صوره للزرع ده عشان لما نيجي نحل الاسيمنت نفهم ايه اللي بيحصل ..

الواي أكسيس هو ال  $E[\text{loss}]$  .. والإكس أكسيس هو ال  $\lambda$  .. في الاسيمنت انت هتغير اللمضه و تحدد لمضه كويسه باستخدام ال  $\text{corss}$  validation .. عشان اللمضه هايبر بارمتر ... انت هتلاحظ ال  $E[\text{Loss}]$  زي ما الدكتور راسمها كدا .. الكيرف ده عليه 3 حاجات ... نويز و بايس تربيع و فارينس .. لو شلت النويز هتلاقي الكيرف نزل لتحت شويه لأنو بقي عباره عن  $\text{variance} + \text{bias}^2$  .. في ال practice .. انت هتعود تشوف لمضه بتجييك احسن قيمه لل error فين .. و بعدين تشيل منو النويز .. وهو ده ال sum of bias and variance



الدكتور بعد كدا دخل علي المثال اللي ف السلايد :

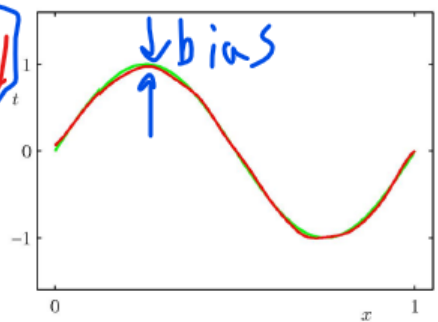
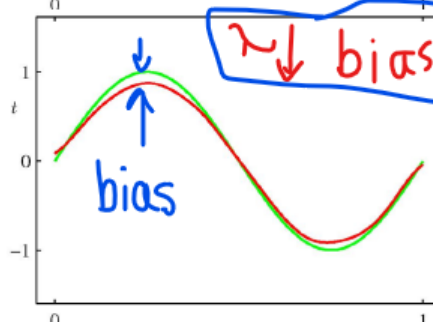
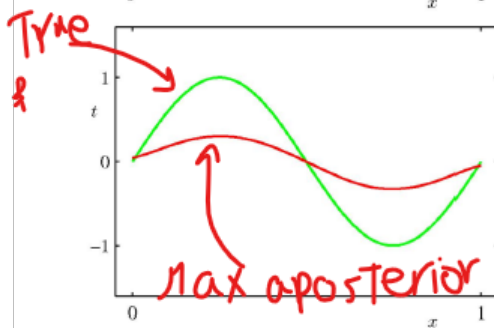
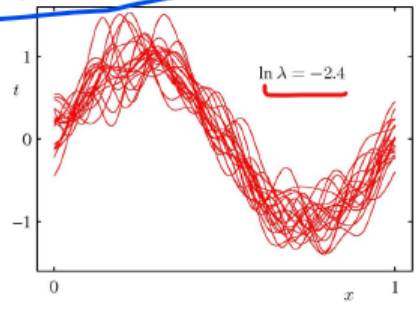
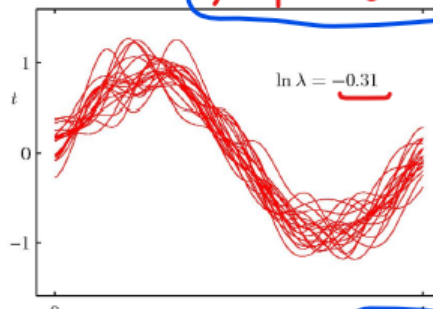
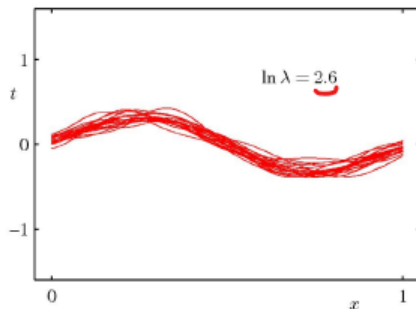


# Bias-Variance Decomposition

## • Example

non-linear regression

$\lambda \uparrow$  variance  $\uparrow$



University of Waterloo

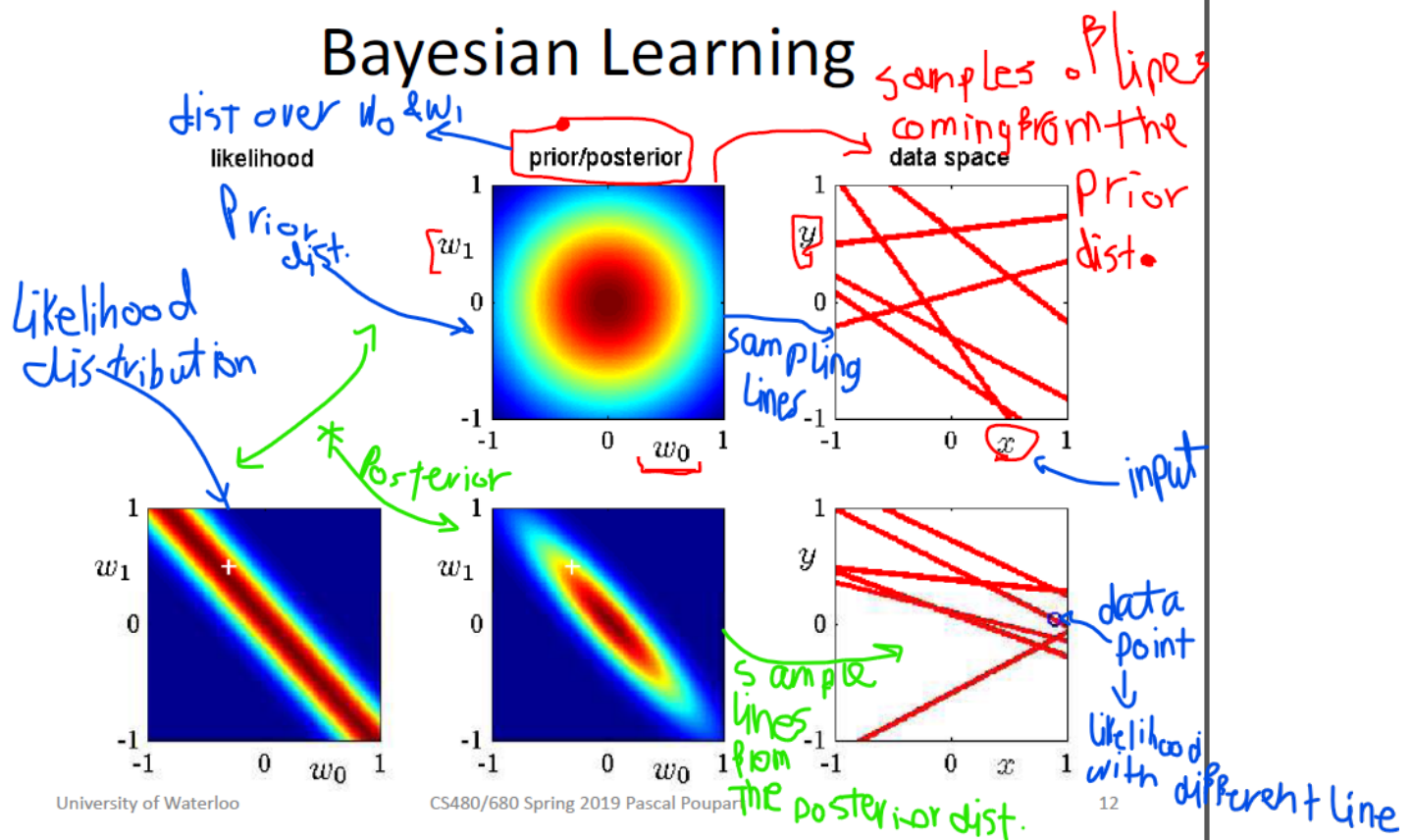
CS480/680 Spring 2019 Pascal Poupart

10

تعال نتكلم علي Bayesian linear regression .. في الحاله دييه بدل ما تاخذ ال Hypothesis  $w$  اللي بت maximize the posterior .. انت حرفياً ممكن بس تحسب ال posterior و تشتغل بيها علي طول ... انت هنا هتشتغل بال posterior .. هن rewrite ال expression بس مش أكثر ... لما بنعمل Bayesian learning .. احنا بنبدأ ب prior و نحسب ال Posterior .. فهنا انت بتحسب ال posterior distribution ..

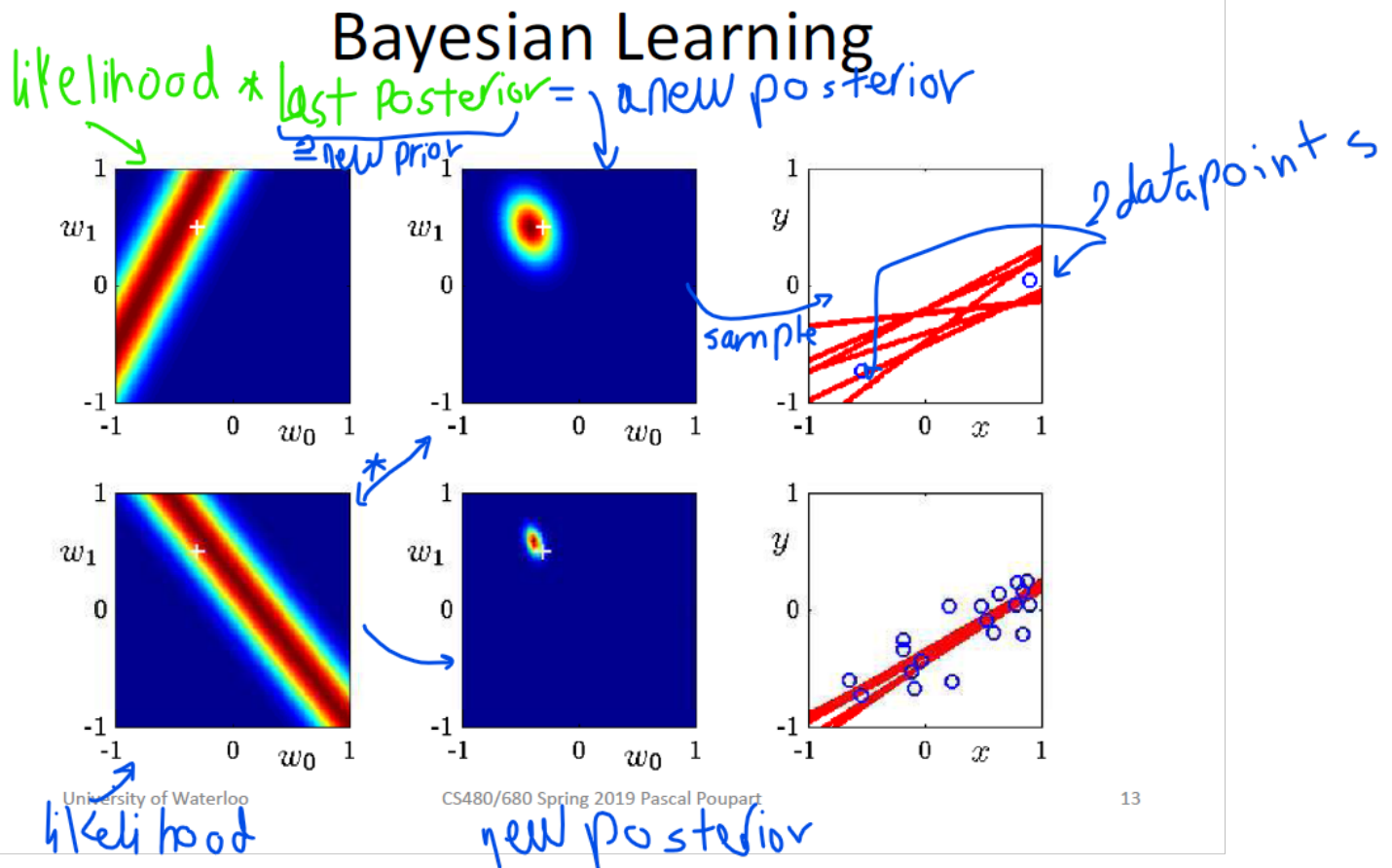
تعال نبص علي صوره

# Bayesian Learning



الصورة اللي فاتت ديه بعد خطوه واحده من ال observation for data point .. لو كملنا علي داتا بوينت تانيه .. هنلاقي الصورة الجايه ..

# Bayesian Learning





خد بالك من نقطه مهمه .. كل pair من  $W_0$  and  $W_1$  .. بيدليك خط .. ده عشان الخط ده مرسوم مع ال  $x$  and  $y$  values لما بستخدم ال  $w_0$  and  $w_1$  .

الدكتور بيقول اننا بنعمل predictions in Bayesian learning عن طريق انك تاخد weighted combinations of the predictions made by each hypothesis .. انت عندك داتا سيت  $X$  bar .. انت هتعمل prediction for a new input  $X$  star و هنتوقع ال  $y$  star .. ده نوع ال prediction اللي هنعوز نعملو ..

المشكله اني معنديش hypothesis واحده .. خد weighted combinations of the prediction of each mode where the weights are determined by posterior distribution .