# Lecture 5: Linear Regression by Maximum Likelihood, Maximum A Posteriori and Bayesian Learning CS480/680 Intro to Machine Learning

2023-1-24

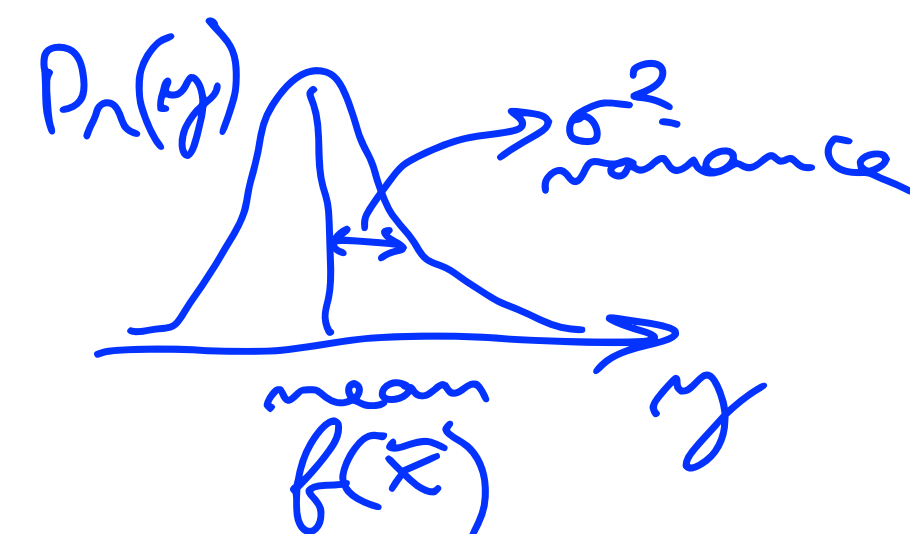Pascal Poupart
David R. Cheriton School of Computer Science

UNIVERSITY OF
WATERLOO

# Noisy Linear Regression

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

$$y = f(\overline{x}) + \epsilon$$

$$w^T\overline{x} \qquad N(0, \sigma^2)$$

- Gaussian noise:

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

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

*(handwritten annotations: "dataset", "normalization constant", "mean", "variance", $\Pr(y)$, $\sigma^2$ variance, mean $f(\overline{x})$, $y$)*

UNIVERSITY OF
WATERLOO

# Maximum Likelihood

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

$$\boldsymbol{w}^* = argmax_{\boldsymbol{w}} \Pr(\boldsymbol{y}|\overline{\boldsymbol{X}}, \boldsymbol{w}, \sigma)$$

$$= argmax_{\boldsymbol{w}} \prod_n e^{-\frac{\left(y_n - \boldsymbol{w}^T \overline{x}_n\right)^2}{2\sigma^2}}$$

*take log*

$$= argmax_{\boldsymbol{w}} \sum_n -\frac{\left(y_n - \boldsymbol{w}^T \overline{x}_n\right)^2}{2\sigma^2}$$

$$= argmin_{\boldsymbol{w}} \sum_n \left(y_n - \boldsymbol{w}^T \overline{\boldsymbol{x}}_n\right)^2$$

- We arrive at the original least square problem!

UNIVERSITY OF
WATERLOO

# Maximum A Posteriori

- Alternative objective: find $\boldsymbol{w}^*$ with highest posterior probability

- Consider Gaussian prior: $\Pr(\boldsymbol{w}) = N(\boldsymbol{0}, \boldsymbol{\Sigma})$

  *proportional* $\left(\begin{smallmatrix} 0 \\ \vdots \\ 0 \end{smallmatrix}\right)$ (*matrix*)

- Posterior:

  $$\Pr(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}) \propto \Pr(\boldsymbol{w}) \Pr(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w})$$

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

UNIVERSITY OF
WATERLOO

# Maximum A Posteriori

- Optimization:

$$\boldsymbol{w}^* = argmax_{\boldsymbol{w}} \; \Pr(\boldsymbol{w}|\overline{\boldsymbol{X}}, \boldsymbol{y})$$

$$= argmax_{\boldsymbol{w}} \; -\sum_n\left(y_n - \boldsymbol{w}^T\overline{\boldsymbol{x}}_n\right)^2 - \boldsymbol{w}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{w}$$

$$= argmin_{\boldsymbol{w}} \; \sum_n\left(y_n - \boldsymbol{w}^T\overline{\boldsymbol{x}}_n\right)^2 + \boldsymbol{w}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{w}$$

Let $\boldsymbol{\Sigma}^{-1} = \lambda\boldsymbol{I}$ then

$$= argmin_{\boldsymbol{w}} \; \sum_n\left(y_n - \boldsymbol{w}^T\overline{\boldsymbol{x}}_n\right)^2 + \lambda\left\|\boldsymbol{w}\right\|_2^2$$

- We arrive at the original **regularized** least square problem!

*Handwritten annotations:*

Take log and drop constants

$$\Sigma^{-1} = \begin{pmatrix} \lambda & 0 \\ 0 & \ddots \\ & & \lambda \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 1/\lambda & 0 \\ 1/\lambda & \ddots \\ 0 & & 1/\lambda \end{pmatrix}$$

UNIVERSITY OF
WATERLOO

# Expected Squared Loss

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

$$E[L] = \int_{\boldsymbol{x},y} \Pr(\boldsymbol{x}, y)\left(y - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right)^2 d\boldsymbol{x}dy$$

$$= \int_{\boldsymbol{x},y} \Pr(\boldsymbol{x}, y)\left(y - f(\boldsymbol{x}) + f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right)^2 d\boldsymbol{x}dy$$

$$= \int_{\boldsymbol{x},y} \Pr(\boldsymbol{x}, y)\left[\left(y - f(\boldsymbol{x})\right)^2 + 2\underbrace{\left(y - f(\boldsymbol{x})\right)}\left(f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right) + \left(f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right)^2\right] d\boldsymbol{x}dy$$

Expectation with respect to $y$ is 0

$$E[L] = \underbrace{\int_{\boldsymbol{x},y} \Pr(\boldsymbol{x}, y)\left(y - f(\boldsymbol{x})\right)^2 d\boldsymbol{x}dy}_{\text{noise (constant)}} + \underbrace{\int_{\boldsymbol{x}} \Pr(\boldsymbol{x})\left(f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right)^2 d\boldsymbol{x}}_{\text{error (depends on } \boldsymbol{w})}$$

UNIVERSITY OF
WATERLOO

# Expected Squared Loss

- Let's focus on the error part, which depends on $\boldsymbol{w}$

$$E_{\boldsymbol{x}}[(f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}})^2] = \int_{\boldsymbol{x}} \Pr(\boldsymbol{x})\left(f(\boldsymbol{x}) - \boldsymbol{w}^T\overline{\boldsymbol{x}}\right)^2 d\boldsymbol{x}$$

- But the choice of $\boldsymbol{w}$ depends on the dataset $S$

- Instead consider expectation with respect to $S$

$$E_S\left[(f(\boldsymbol{x}) - \boldsymbol{w}_S^T\overline{\boldsymbol{x}})^2\right]$$

where $\boldsymbol{w}_S$ is the weight vector obtained based on $S$

UNIVERSITY OF
**WATERLOO**

# Bias-Variance Decomposition

- Decompose squared loss

$$E_S\big[\big(f(\pmb{x}) - \pmb{w}_S^T\overline{\pmb{x}}\big)^2\big]$$

$$= E_S\big[f(\pmb{x}) - E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big] + E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big] - \pmb{w}_S^T\overline{\pmb{x}}\big]^2$$

$$= E_S\left[\big(f(\pmb{x}) - E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big]\big)^2 + 2\big(f(\pmb{x}) - E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big]\big)\underbrace{\big(E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big] - \pmb{w}_S^T\overline{\pmb{x}}\big)}_{\text{Expectation is 0}} + \big(E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big] - \pmb{w}_S^T\overline{\pmb{x}}\big)^2\right]$$

$$= \underbrace{\big(f(\pmb{x}) - E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big]\big)^2}_{\text{bias}^2} + \underbrace{E_S\left[\big(E_S\big[\pmb{w}_S^T\overline{\pmb{x}}\big] - \pmb{w}_S^T\overline{\pmb{x}}\big)^2\right]}_{\text{variance}}$$

UNIVERSITY OF
WATERLOO

# Bias-Variance Decomposition
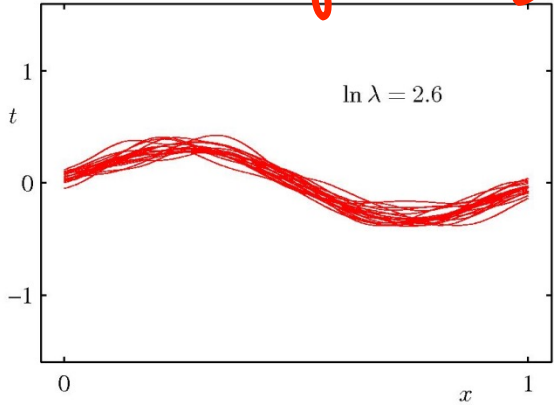
- Hence:

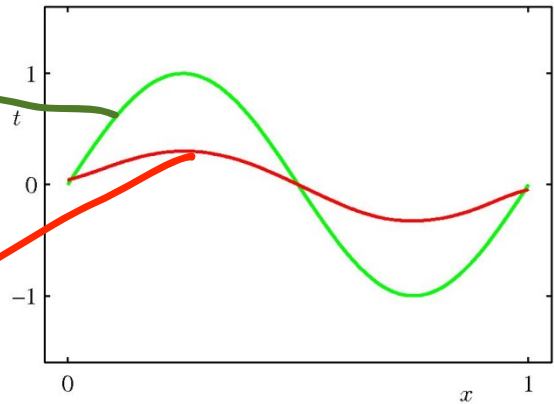$$E[loss] = (bias)^2 + variance + noise$$

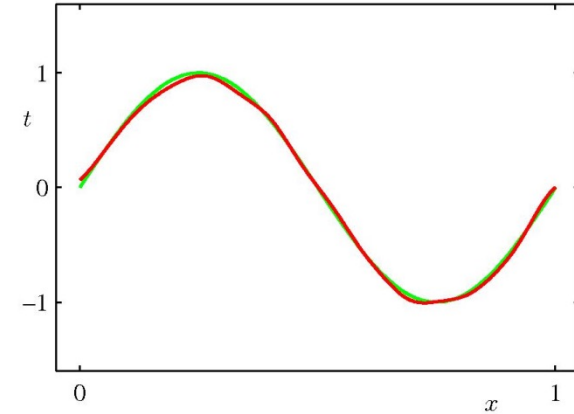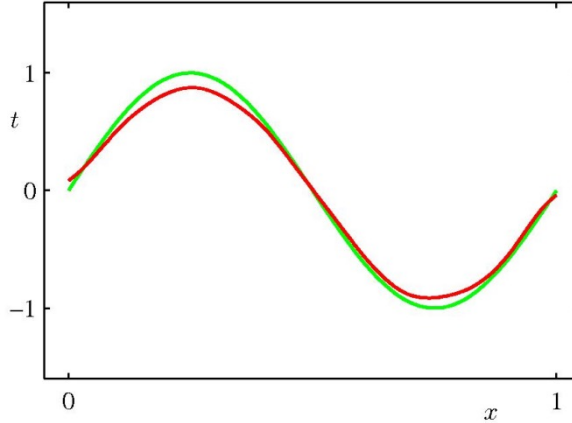- Picture:

# Bias-Variance Decomposition

- Example

underfitting

overfitting

variance



$\ln \lambda = 2.6$

$\ln \lambda = -0.31$

$\ln \lambda = -2.4$

$f(x)$ true fn

estimate $E(w_s^T x)$

bias

UNIVERSITY OF
WATERLOO

# Bayesian Linear Regression

- We don't know if $\boldsymbol{w}^*$ is the true underlying $\boldsymbol{w}$

- Instead of making predictions according to $\boldsymbol{w}^*$, compute the weighted average prediction according to $\Pr(\boldsymbol{w}|\bar{\boldsymbol{X}}, \boldsymbol{y})$

$$\Pr(\boldsymbol{w}|\bar{\boldsymbol{X}}, \boldsymbol{y}) = ke^{-\frac{\boldsymbol{w}^T\Sigma^{-1}\boldsymbol{w}}{2}}e^{-\frac{\Sigma_n\left(y_n - \boldsymbol{w}^T\bar{\boldsymbol{x}}_n\right)^2}{2\sigma^2}}$$

$$= ke^{-\frac{1}{2}(\boldsymbol{w}-\bar{\boldsymbol{w}})^T A(\boldsymbol{w}-\bar{\boldsymbol{w}})} = N(\bar{\boldsymbol{w}}, A^{-1})$$

*mean* (→ $\bar{\boldsymbol{w}}$)
*covariance* (→ $A^{-1}$)

where $\bar{\boldsymbol{w}} = \sigma^{-2}A^{-1}\bar{\boldsymbol{X}}\boldsymbol{y}$

$$A = \sigma^{-2}\bar{\boldsymbol{X}}\bar{\boldsymbol{X}}^T + \Sigma^{-1}$$

UNIVERSITY OF
WATERLOO

# Bayesian Learning

# Bayesian Learning

true weights

UNIVERSITY OF
WATERLOO

# Bayesian Prediction

- Let $\boldsymbol{x}_*$ be the input for which we want a prediction and $y_*$ be the corresponding prediction

$$\Pr(y_*|\overline{\boldsymbol{x}}_*, \overline{\boldsymbol{X}}, \boldsymbol{y}) = \int_{\boldsymbol{w}} \Pr(y_*|\overline{\boldsymbol{x}}_*, \boldsymbol{w}) \Pr(\boldsymbol{w}|\overline{\boldsymbol{X}}, y) \, d\boldsymbol{w}$$

*query point*

$$= k \int_{\boldsymbol{w}} e^{-\frac{\left(y_* - \overline{\boldsymbol{x}}_*^T \boldsymbol{w}\right)^2}{2\sigma^2}} e^{-\frac{1}{2}(\boldsymbol{w}-\overline{\boldsymbol{w}})^T A (\boldsymbol{w}-\overline{\boldsymbol{w}})} \, d\boldsymbol{w}$$

$$= N(\sigma^{-2}\overline{\boldsymbol{x}}_*^T A^{-1} \overline{\boldsymbol{X}} \boldsymbol{y}, \sigma^2 + \overline{\boldsymbol{x}}_*^T A^{-1} \overline{\boldsymbol{x}}_*)$$

*mean*  *variance*

UNIVERSITY OF
WATERLOO