

Introduction to Machine Learning - CS725  
Instructor: Prof. Ganesh Ramakrishnan  
Lecture 07 - Bayesian Linear Regression  
(Gaussian and Laplacian priors), Regularized  
Linear Regression (Ridge Regression and Lasso)

# Recap: Summary for MAP estimation with Normal Distribution

- Univariate: With  $\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$  and  $x \sim \mathcal{N}(\mu, \sigma^2)$ ,  
 $p(\mu|\mathcal{D}) \sim \mathcal{N}(\mu_m, \sigma_m^2)$

$$\frac{1}{\sigma_m^2} = \frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}$$
$$\frac{\mu_m}{\sigma_m^2} = \frac{m}{\sigma^2} \hat{\mu}_{mle} + \frac{\mu_0}{\sigma_0^2}$$

- Multivariate: By **extrapolation** (Bayesian setting for fixed  $\Sigma$ )  
 $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$ ,  $\mu \sim \mathcal{N}(\mu_0, \Sigma_0)$  &  $p(\mu|\mathcal{D}) \sim \mathcal{N}(\mu_m, \Sigma_m)$

$$\Sigma_m^{-1} = m\Sigma^{-1} + \Sigma_0^{-1}$$
$$\Sigma_m^{-1} \mu_m = m\Sigma^{-1} \hat{\mu}_{mle} + \Sigma_0^{-1} \mu_0$$

# Different Estimators

Recap: Mean and Mode coincide for (Multivariate) Gaussian ==>

Bayes Estimate = MAP estimate

	Point?	$p(x D)$
MLE	$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} LL(D \theta)$	$p(x \theta_{MLE})$
Bayes Estimator	$\hat{\theta}_B = E_{p(\theta D)} E[\theta]$	$p(x \theta_B)$
MAP	$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} p(\theta D)$	$p(x \theta_{MAP})$
Pure Bayesian	$p(\theta D) \sim N(u_m, \Sigma_m)$ $p(x D) \sim N(u_m + \dots, \Sigma_m + \dots)$	$p(\theta D) = \frac{p(D \theta)p(\theta)}{\int_m p(D \theta)p(\theta)d\theta}$ $p(D \theta) = \prod_{i=1} p(x_i \theta)$ $p(x D) = \int_{\theta} p(x \theta)p(\theta D)$

# Recap: Back to Linear Regression: Why Bayesian?

- The Bayesian interpretation of probabilistic estimation is a logical extension that enables reasoning with uncertainty **but in the light of some background belief**
- **Bayesian linear regression**: A Bayesian alternative to **Maximum Likelihood** least squares regression to address **overfitting**
- Continue with Normally distributed errors
- Model the  $\mathbf{w}$  using a prior distribution and use the posterior over  $\mathbf{w}$  as the result
- Intuitive Prior: Components of  $\mathbf{w}$  should not become too large!

# Recap: Prior Distribution for $\mathbf{w}$

$$y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- Maximum (log)-likelihood estimate is  $\hat{\mathbf{w}}_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T y$
- We can use a Prior distribution on  $\mathbf{w}$  to avoid over-fitting

$u=0$ , since  $\phi_i$  should have  $w_i=0$  unless  $\phi_i$  provides "signal"

$w_i \sim \mathcal{N}(0, \frac{1}{\lambda}) \rightarrow \frac{1}{\lambda} = \sigma^2 = \frac{1}{\text{precision}}$

(that is, each component  $w_i$  is approximately bounded within  $\pm \frac{3}{\sqrt{\lambda}}$  by the 3 -  $\sigma$  rule).  $\lambda$  is also called the precision of the Gaussian

- Q: Bayesian Estimation?

# Recap: Multivariate Normal Distribution and MAP estimate

- 1 If  $w_i \sim \mathcal{N}(0, \frac{1}{\lambda})$  then  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\lambda}I)$  where  $I$  is an  $n \times n$  identity matrix
- 2  $\Rightarrow$  That is,  $\mathbf{w}$  has a multivariate Gaussian distribution  $\Pr(\mathbf{w}) = \frac{1}{(\frac{2\pi}{\lambda})^{\frac{n}{2}}} e^{-\frac{\lambda}{2}\|\mathbf{w}\|_2^2}$  with  $\mu_0 = \mathbf{0}$ .  $\Sigma_0 = \frac{1}{\lambda}I$
- 3 Consider Bayesian Estimation for multivariate Gaussian on  $\mathbf{w}$

# Posterior Distribution for $\mathbf{w}$ for Linear Regression

- Given  $y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2) \Rightarrow$   $\mu_0 = 0$   
 $\Sigma_0 = \frac{1}{\lambda} \mathbf{I}$   
 $y \sim \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}), \sigma^2)$ ,  $\mathbf{w} \sim \mathcal{N}(\mu_0, \Sigma_0)$  where,  $w_i \sim \mathcal{N}(0, \frac{1}{\lambda})$
- We want to find  $P(\mathbf{w}|D) = \mathcal{N}(\mu_m, \Sigma_m)$

Invoking the Bayes Estimation results from before (homework):

$$\begin{aligned}
 P(\mathbf{w}|D) &\propto P(D|\mathbf{w}) P(\mathbf{w}) \propto \prod_{i=1}^m \exp\left(-\frac{(y_i - \mathbf{w}^T \phi(x_i))^2}{2\sigma^2}\right) \times \exp\left[-\frac{1}{2}(\mathbf{w} - \mu_0)^T \Sigma_0^{-1} (\mathbf{w} - \mu_0)\right] \\
 &= \exp\left[-\mathbf{w}^T \left(\Sigma_0^{-1} + \underbrace{\sum_{i=1}^m \frac{\phi^T(x_i) \phi(x_i)}{2\sigma^2}}_{\text{green wavy line}}\right) \mathbf{w} + \underbrace{2 \sum_{i=1}^m \mathbf{w}^T \left(\frac{y_i \phi(x_i)}{2\sigma^2} + \mu_0^T \Sigma_0^{-1}\right)}_{\text{green wavy line}} - \dots\right]
 \end{aligned}$$

# Posterior Distribution for $\mathbf{w}$ for Linear Regression

- Given  $y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2) \Rightarrow y \sim \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}), \sigma^2)$ ,  $\mathbf{w} \sim \mathcal{N}(\mu_0, \Sigma_0)$  where,  $w_i \sim \mathcal{N}(0, \frac{1}{\lambda})$
- We want to find  $P(\mathbf{w}|D) = \mathcal{N}(\mu_m, \Sigma_m)$   
Invoking the Bayes Estimation results from before (homework):

$$\mu_0 = 0$$

$$\Sigma_0^{-1} = \lambda \mathbf{I}$$

Substitute!

$$\Sigma_m^{-1} = \frac{1}{\sigma^2} \Phi^T \Phi + \Sigma_0^{-1}$$
$$\Sigma_m^{-1} \mu_m = \Phi^T y / \sigma^2 + \cancel{\Sigma_0^{-1} \mu_0}$$

$$\Sigma_m^{-1} = \frac{\Phi^T \Phi}{\sigma^2} + \lambda \mathbf{I}$$

$$\mu_m = \left( \frac{\Phi^T \Phi}{\sigma^2} + \lambda \mathbf{I} \right)^{-1} \frac{\Phi^T y}{\sigma^2}$$



# Finding $\mu_m$ & $\Sigma_m$ for $\mathbf{w}$

Setting  $\Sigma_0 = \frac{1}{\lambda}I$  and  $\mu_0 = \mathbf{0}$

$$\Sigma_m^{-1} \mu_m = \Phi^T \mathbf{y} / \sigma^2$$

$$\Sigma_m^{-1} = \lambda I + \Phi^T \Phi / \sigma^2$$

$$\mu_m = \frac{(\lambda I + \Phi^T \Phi / \sigma^2)^{-1} \Phi^T \mathbf{y}}{\sigma^2}$$

or

$$\mu_m = \frac{(\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}}{\sigma^2}$$

# MAP and Bayes Estimates

- $\Pr(\mathbf{w} \mid \mathcal{D}) = \mathcal{N}(\mathbf{w} \mid \underline{\mu_m}, \underline{\Sigma_m})$
- The **MAP estimate** or mode under the Gaussian posterior is the mode of the posterior  $\Rightarrow$

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \mathcal{N}(\mathbf{w} \mid \mu_m, \Sigma_m) = \underline{\mu_m}$$

- Similarly, the **Bayes Estimate**, or the expected value under the Gaussian posterior is the mean  $\Rightarrow$

$$\hat{\mathbf{w}}_{Bayes} = E_{\Pr(\mathbf{w} \mid \mathcal{D})}[\mathbf{w}] = E_{\mathcal{N}(\mu_m, \Sigma_m)}[\mathbf{w}] = \underline{\mu_m}$$

- Summarily:

Recall:  $\hat{\mathbf{w}}_{MLE}$  had no  $\sigma^2$  (of course no  $\lambda$ )

$$\mu_{MAP} = \mu_{Bayes} = \mu_m = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

# Predictive distribution for linear Regression

Pure Bayesian:  $\Pr(y|x, D)$

- $\hat{\mathbf{w}}_{MAP}$  helps avoid overfitting as it takes regularization into account
- But we miss the modeling of uncertainty when we consider only  $\hat{\mathbf{w}}_{MAP}$
- **Eg:** While predicting diagnostic results on a new patient  $x$ , along with the value  $y$ , we would also like to know the uncertainty of the prediction  $\Pr(y | x, D)$ . Recall that  $y = \mathbf{w}^T \phi(x) + \varepsilon$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

$$\Pr(y | \mathbf{x}, D) = \Pr(y | \mathbf{x}, \langle \mathbf{x}_1, y_1 \rangle \dots \langle \mathbf{x}_m, y_m \rangle)$$

Expect:  $\Pr(y|x, D) \sim \mathcal{N}(\mu_n^T \phi(x), \phi^T(x) \Sigma_n \phi(x) + \sigma^2)$

# Pure Bayesian Regression Summarized (optional)

By definition, regression is about finding  $(y \mid \mathbf{x}, \mathcal{D})$ . By Bayes Rule

$$\begin{aligned}\Pr(y \mid \mathbf{x}, \mathcal{D}) &= \Pr(y \mid \mathbf{x}, \langle \mathbf{x}_1, y_1 \rangle \dots \langle \mathbf{x}_m, y_m \rangle) \\ &= \int_{\mathbf{w}} \Pr(y \mid \mathbf{w}; \mathbf{x}) \Pr(\mathbf{w} \mid \mathcal{D}) d\mathbf{w} \\ &\sim \mathcal{N}(\mu_m^T \phi(\mathbf{x}), \underbrace{\sigma^2 + \phi^T(\mathbf{x}) \Sigma_m \phi(\mathbf{x})})\end{aligned}$$

where

$$y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon \text{ and } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\mathbf{w} \sim \mathcal{N}(0, \alpha I) \text{ and } \mathbf{w} \mid \mathcal{D} \sim \mathcal{N}(\mu_m, \Sigma_m)$$

$$\mu_m = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \text{ and } \Sigma_m^{-1} = \lambda I + \Phi^T \Phi / \sigma^2$$

$$\text{Finally } y \sim \mathcal{N}(\underbrace{\mu_m^T \phi(\mathbf{x})}, \underbrace{\phi^T(\mathbf{x}) \Sigma_m \phi(\mathbf{x})})$$

# MAP (and Bayes) Inference *(Rewriting differently)*

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w} \mid \mathcal{D}) = \underset{\mathbf{w}}{\operatorname{argmax}} \log \Pr(\mathbf{w} \mid \mathcal{D}), \text{ where,}$$

$$-\log \Pr(\mathbf{w} \mid \mathcal{D}) = \frac{n}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_m| + \frac{1}{2} (\mathbf{w} - \mu_m)^T \Sigma_m^{-1} (\mathbf{w} - \mu_m)$$

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} -\log \Pr(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \Sigma_m^{-1} \mathbf{w} - \mathbf{w}^T \Sigma_m^{-1} \mu_m$$

(expanding/canceling redundant terms & completing squares: Tut 3)

*By substituting for  $\Sigma_m$  &  $\mu_m$*

*Recall: log is monotonically increasing*

$$\operatorname{argmax}_{\mathbf{w}} f(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} -f(\mathbf{w})$$

# MAP (and Bayes) Inference

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w} \mid \mathcal{D}) = \underset{\mathbf{w}}{\operatorname{argmax}} \log \Pr(\mathbf{w} \mid \mathcal{D}), \text{ where,}$$

$$-\log \Pr(\mathbf{w} \mid \mathcal{D}) = \frac{n}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_m| + \frac{1}{2} (\mathbf{w} - \mu_m)^T \Sigma_m^{-1} (\mathbf{w} - \mu_m)$$

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} -\log \Pr(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \Sigma_m^{-1} \mathbf{w} - \mathbf{w}^T \Sigma_m^{-1} \mu_m$$

(expanding/canceling redundant terms & completing squares: Tut 3)

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \underbrace{\|\phi \mathbf{w} - \mathbf{y}\|^2}_{\text{Least squares loss}} + \underbrace{\sigma^2 \lambda \|\mathbf{w}\|^2}_{\text{Regularizer}} = \mathbf{w}_{Ridge}$$

is the same as that of **Penalized** *Regularized Regression*.

$$w_{\text{MAP}} = \underset{w}{\text{argmin}} \quad \frac{1}{2} \|\Phi w - y\|^2 + \underbrace{\lambda \sigma^2 \|w\|^2}_{\text{independently } \|w\|^2}$$

Independently  $\|w\|^2$

is minimized only  
when each  $w_i = 0$

THIS IS COMMON SENSE! ALL WE HAVE DONE IS GIVEN SOME  
PROBABILISTIC INTERPRETATION TO COMMON SENSE!

# Penalized Regularized Least Squares Regression

- The Bayes and MAP estimates for Linear Regression coincide with *Regularized Ridge Regression*

$$\mathbf{w}_{Ridge} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Eg: If  $\phi_i = \phi_j$ ,  $w_i \neq 0$  I want  $w_j = 0$

- Intuition:** To discourage redundancy and/or stop coefficients of  $\mathbf{w}$  from becoming too large in magnitude, add a penalty to the error term used to estimate parameters of the model.
- The general **Penalized Regularized L.S Problem:**

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \Omega(\mathbf{w})$$

Can also be generic Error( $\phi, w, y$ )



# Penalized Regularized Least Squares: Examples

- The general **Penalized Regularized L.S Problem**:

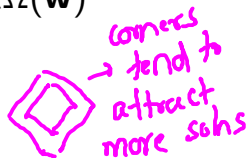
$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \Omega(\mathbf{w})$$

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2 \Rightarrow$  Ridge Regression

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 \Rightarrow$  Lasso

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_0 \Rightarrow$  Support-based penalty

(A compromise)



corners  
→ tend to  
attract  
more solns

# non-zeros..  
not (to) differentiable

- Some  $\Omega(\mathbf{w})$  correspond to priors that can be expressed in close form. Some give good working solutions. Some norms are mathematically easier to handle

$$\|\mathbf{w}\|_2^2 \sim N(\cdot) \sim \exp(-\|\cdot\|^2) \cdot \|\mathbf{w}\|_1 \sim \text{Lap}(\cdot) \sim \exp(-|\cdot|)$$

# Constrained Regularized Least Squares Regression

- **Intuition:** To discourage redundancy and/or stop coefficients of  $\mathbf{w}$  from becoming too large in magnitude, constrain the error minimizing estimate using a penalty
- The general **Constrained Regularized L.S. Problem:**

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$$

*such that  $\Omega(\mathbf{w}) \leq \theta$*

- Claim: For any **Penalized** formulation with a particular  $\lambda$ , there exists a corresponding **Constrained** formulation with a corresponding  $\theta$
- **Proof of Equivalence:** Requires tools of Optimization/duality

# Constrained Regularized Least Squares: Examples

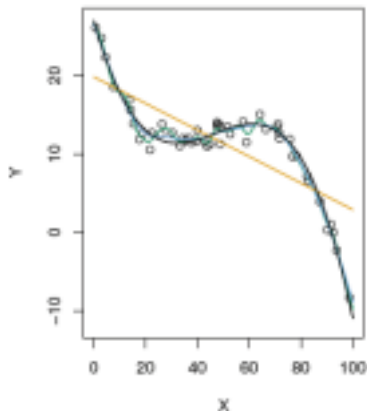
- The general **Constrained Regularized L.S. Problem**:

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$$

*such that  $\Omega(\mathbf{w}) \leq \theta$*

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2 \Rightarrow$  **Ridge Regression**
- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 \Rightarrow$  **Lasso**
- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_0 \Rightarrow$  **Support-based penalty**

# Polynomial regression



- Consider a degree 3 polynomial regression model as shown in the figure
- Each bend in the curve corresponds to increase in  $\|w\|$
- Eigen values of  $(\Phi^T \Phi + \lambda I)$  are indicative of curvature. Increasing  $\lambda$  reduces the curvature

# Do Closed-form solutions Always Exist?

- Linear regression and Ridge regression both have closed-form solutions

- For linear regression,

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

- For ridge regression,

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

(for linear regression,  $\lambda = 0$ )

- What about optimizing the formulations (constrained/penalized) of Lasso ( $L_1$  norm)? And support-based penalty ( $L_0$  norm)? **Also requires tools of Optimization/duality**

# Lasso Regularized Least Squares Regression

- The general **Penalized Regularized L.S Problem**:

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \Omega(\mathbf{w})$$

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 \Rightarrow$  **Lasso**
- Lasso Regression*

$$\mathbf{w}_{lasso} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

- Lasso is the MAP estimate of Linear Regression subject to Laplace Prior on  $\mathbf{w}$   $\sim Laplace(0, \theta)$

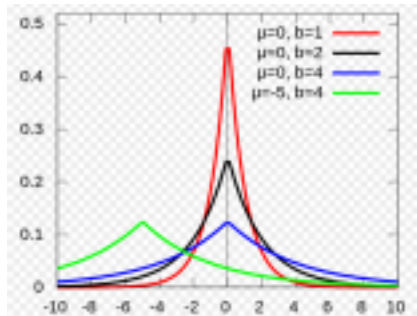
$$Laplace(w_i \mid \mu, b) = \frac{1}{2b} \exp \left( -\frac{|w_i - \mu|}{b} \right)$$

# Gaussian Hare vs. Laplacian Tortoise



- Gaussian easier to estimate

$$w_{MAP}^{Gaussian} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$



- Laplacian yields more sparsity

No closed form for  $w_{MAP}^{Lap.}$   
Iterative algos...

























