Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 6 - Support Vector Regression and Optimization Basics

|  | **Point?** | $p(x\|D)$ |
|---|---|---|
| MLE | $\hat{\theta}_{MLE} = \text{argmax}_\theta \, LL(D\|\theta)$ | $p(x\|\theta_{MLE})$ |
| Bayes Estimator | $\hat{\theta}_B = E_{p(\theta\|D)}E[\theta]$ | $p(x\|\theta_B)$ |
| MAP | $\hat{\theta}_{MAP} = \text{argmax}_\theta \, p(\theta\|D)$ | $p(x\|\theta_{MAP})$ |
| Pure Bayesian |  | $p(\theta\|D) = \dfrac{p(D\|\theta)p(\theta)}{\int_m p(D\|\theta)p(\theta)d\theta}$ $p(D\|\theta) = \prod_{i=1}^{m} p(x_i\|\theta)$ $p(x\|D) = \int_\theta p(x\|\theta)p(\theta\|D)d\theta$ |

where $\theta$ is the parameter

- $\hat{\mathbf{w}}_{MAP}$ 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 $\Pr(y \mid x, D)$. Recall that $y = \mathbf{w}^T \phi(x) + \varepsilon$ and $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

$$\Pr(y \mid \mathbf{x}, \mathcal{D}) = \Pr(y \mid \mathbf{x}, <\mathbf{x}_1, y_1> ... <\mathbf{x}_m, y_m>)$$

# Pure Bayesian Regression Summarized

- By definition, regression is about finding $(y \mid \mathbf{x}, <\mathbf{x}_1, y_1> \ldots <\mathbf{x}_m, y_m>)$
- By Bayes Rule

$$\Pr(y \mid \mathbf{x}, \mathcal{D}) = \Pr(y \mid \mathbf{x}, <\mathbf{x}_1, y_1> \ldots <\mathbf{x}_m, y_m>)$$
$$= \int_{\mathbf{w}} \Pr(y \mid \mathbf{w}; \mathbf{x}) \Pr(\mathbf{w} \mid \mathcal{D}) d\mathbf{w}$$
$$\sim \mathcal{N}\left(\mu_m^T \phi(\mathbf{x}), \sigma^2 + \phi^T(\mathbf{x}) \Sigma_m \phi(\mathbf{x})\right)$$

*where*

$$y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon \text{ and } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$
$$\mathbf{w} \sim \mathcal{N}(0, \alpha I) \text{ and } \mathbf{w} \mid \mathcal{D} \sim \mathcal{N}(\mu_m, \Sigma_m)$$
$$\mu_m = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \text{ and } \Sigma_m^{-1} = \lambda I + \Phi^T \Phi / \sigma^2$$
$$\text{Finally } y \sim \mathcal{N}(\mu_m^T \phi(\mathbf{x}), \phi^T(\mathbf{x}) \Sigma_m \phi(\mathbf{x}))$$

## MAP (and Bayes) Inference

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

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

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

..... (expanding & canceling out redundant terms & completing squares: Tutorial 3)

## MAP (and Bayes) Inference

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

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

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

..... (expanding & canceling out redundant terms & completing squares: Tutorial 3)

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \frac{1}{2\sigma^2}\mathbf{w}^T \left(\phi^T\phi\mathbf{w} - 2\phi^T\mathbf{y}\right) + \lambda\mathbf{w}^T\mathbf{w} = \underset{\mathbf{w}}{\operatorname{argmax}} \frac{1}{2}||\phi\mathbf{w}-\mathbf{y}||^2 + \sigma^2\lambda||\mathbf{w}||^2 = \mathbf{w}_{Ridge}$$

is the same as that of *Regularized Regression*.

$$\mathbf{w}_{Ridge} = \underset{\mathbf{w}}{\operatorname{argmin}} \ ||\phi\mathbf{w} - \mathbf{y}||_2^2 + \lambda\sigma^2||\mathbf{w}||_2^2$$

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

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

- **Intuition:** To discourage redundancy and/or stop coefficients of **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:

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

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

- 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

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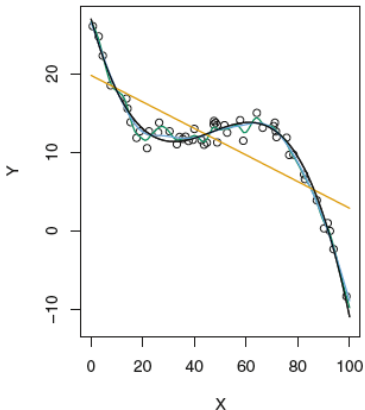- The general **Constrained** **Regularized L.S. Problem**:

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

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

- Claim: For any Penalized formulation with a particular $\lambda$, there exists a corresponding Constrained formulation with a corresponding $\theta$
  - $\Omega(\mathbf{w}) = ||\mathbf{w}||_2^2 \Rightarrow$ **Ridge Regression**
  - $\Omega(\mathbf{w}) = ||\mathbf{w}||_1 \Rightarrow$ **Lasso**
  - $\Omega(\mathbf{w}) = ||\mathbf{w}||_0 \Rightarrow$ **Support**-based penalty
- **Proof of Equivalence:** Requires tools of Optimization/duality

# 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 $(\Phi^\top \Phi + \lambda I)$ 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,
  $$w^* = (\Phi^\top \Phi)^{-1} \Phi^\top y$$
  - For ridge regression,
  $$w^* = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top y$$

  (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

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

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

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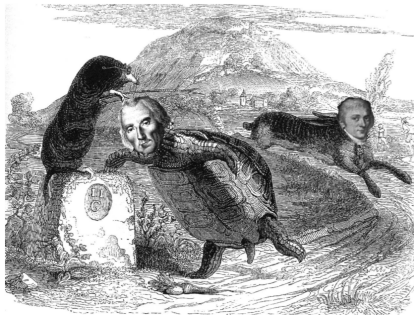  - $\Omega(\mathbf{w}) = ||\mathbf{w}||_0 \Rightarrow$ **Support-based penalty**

- *Lasso* Regression

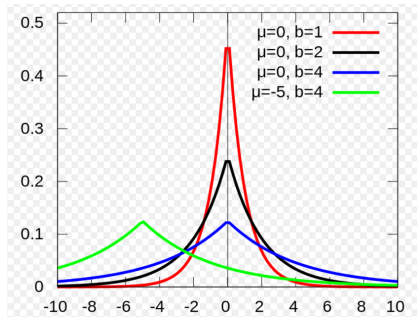$$\mathbf{w}_{lasso} = \underset{\mathbf{w}}{\arg\min} \ ||\Phi\mathbf{w} - \mathbf{y}||_2^2 + \lambda||\mathbf{w}||_1^2$$

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

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

- Gaussian easier to estimate



- Laplacian yields more sparsity

# Support Vector Regression

One more formulation before we look at Tools of Optimization/duality

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate

2. Addressing overfitting
   - Bayesian and Maximum Aposteriori Estimates, Regularization, **Support Vector Regression**

3. How to minimize the resultant and more complex error functions?
   - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

$$f = w \cdot \phi(x) + b$$

- Any point in the band (of $\epsilon$) is not penalized. Thus the loss function is known as *$\epsilon$-insensitive loss*
- Any point outside the band is penalized, and has slackness $\xi_i$ or $\xi_i^*$
- The SVR model curve may not pass through any training point

- The tolerance $\epsilon$ is fixed
- It is desirable that $\forall i$:

- The tolerance $\epsilon$ is fixed
- It is desirable that $\forall i$:
    - $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i$
    - $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

- 1-norm Error, and $L_2$ regularized:

- 1-norm Error, and $L_2$ regularized:

  - $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$
    s.t. $\forall i$,
    $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i$,
    $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$,
    $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and $L_2$ regularized:

- 1-norm Error, and $L_2$ regularized:

  - $\min_{\mathbf{w},b,\xi_i,\xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$
    s.t. $\forall i$,
    $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \le \epsilon + \xi_i$,
    $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \le \epsilon + \xi_i^*$,
    $\xi_i, \xi_i^* \ge 0$

- 2-norm Error, and $L_2$ regularized:

  - $\min_{\mathbf{w},b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$
    s.t. $\forall i$,
    $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \le \epsilon + \xi_i$,
    $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \le \epsilon + \xi_i^*$
  - Here, the constraints $\xi_i, \xi_i^* \ge 0$ are not necessary

- **Unconstrained (Penalized) Optimization:**

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

- **Constrained Optimization 1:**

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

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

- **Constrained Optimization 2 ($t = 1$ or $2$):**

$$\arg\min_{\mathbf{w},b,\xi_i,\xi_i^*} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_i(\xi_i^t + \xi_i^{*t})$$

s.t. $\forall i,\ y_i - w^\top\phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i;\ b + \mathbf{w}^\top\phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

- **Equivalence**: $\lambda$ (Penalized) $\equiv \theta$ (Constrained)
- **Duality**: Dual of Support Vector Regression

## Solving Unconstrained Minimization Problem

- Intuitively: Minimize by setting derivative (gradient) to 0 and hoping to find **closed form** solution.
- When is such a solution a global minimum?
- For most optimization problems, finding closed form solutions is difficult. Even for linear regression (for which closed form solution exists), are there alternative methods?
  - Eg: Consider, $\mathbf{y} = \Phi\mathbf{w}$, where $\Phi$ is a matrix with full column rank, the least squares solution, $\mathbf{w}^* = (\Phi^T\Phi)^{-1}\Phi^T\mathbf{y}$ . Now, imagine that $\Phi$ is a very large matrix. with say, 100,000 columns and 1,000,000 rows. Computation of closed form solution might be challenging.
- How about iterative methods?