

11/08/20

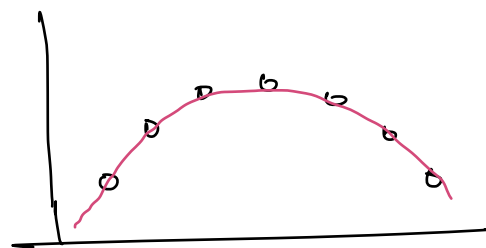
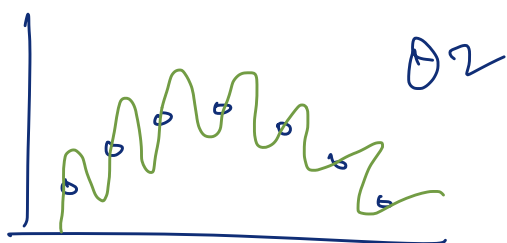
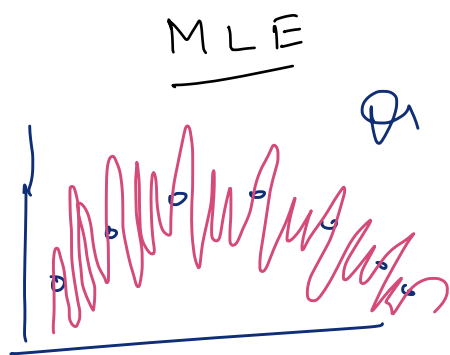
MAP (Maximum a posteriori) Estimate

- Linear regression \rightarrow MLE
- MAP \rightarrow Adding regularizer to the loss function

$$\begin{aligned} \underline{\underline{L(y, \hat{y})}} & \quad y = \text{ground truth} \quad , \quad \hat{y} = \text{prediction} \\ & \quad \theta \quad x = \text{feature vector} \\ L(\theta; x, y) & \equiv \underline{\underline{L(y, f_{\theta}(x))}} \quad y \end{aligned}$$

$$\mathcal{D} = \{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \}$$

$$\arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | x_i) = \arg \min_{\theta} \underbrace{- \sum \log p_{\theta}(y_i | x_i)}_{\text{negative log likelihood}} \quad \nearrow$$



$$\begin{aligned} & p(\theta_3) \uparrow \\ & p(\theta_2) \uparrow \\ & p(\theta_1) \end{aligned}$$

there are infinitely many θ that are the MLE

- MLE has no notion of which of these parameters might be better than others.

- prior belief about the data points.

$$p(y|x) \sim$$

$$y|x \sim \text{RV}$$

$$x \sim \text{RV}$$

- Treat the parameters θ as a RV

$p(\theta)$ - prior distribution that encodes which functions are more likely than others.

what $p(\theta)$ makes smooth f_{θ} more likely and non-smooth f_{θ} less likely?

→ $\theta, P(\theta)$

Roughly speaking, as the magnitude of θ increases, $P(\theta)$ shows decrease

$$P(\theta) = \mathcal{N}(\theta; 0, \sigma^2 I) \quad \text{--- MVG}$$

$$\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$P(\theta | \mathcal{D}) = \frac{\overbrace{P(\mathcal{D} | \theta)}^{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}}{\underbrace{P(\mathcal{D})}_N}$$

$$= \prod_{i=1}^N P(x_i) \underbrace{P(y_i | x_i, \theta)}_{P(y_i | x_i)}$$

$$\log P(\theta | \mathcal{D}) = \sum_{i=1}^N \{ \log P(x_i) + \log P(y_i | x_i, \theta) \} + \log P(\theta)$$

arg max

$$\underset{\theta}{\text{arg max}} = \sum_{i=1}^N \log P(y_i | x_i, \theta) + \log P(\theta)$$

$$\underset{\theta}{\text{arg min}} \sum -\log P(y_i | x_i, \theta) - \log P(\theta)$$

MVG

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$\mu = 0$

$$\frac{1}{(2\pi)^{d/2} (\sigma^2)^{d/2}} \exp\left(-\frac{1}{2} x^T (\sigma^2 I)^{-1} x\right) \quad \text{isotropic Gaussian}$$

$$\Sigma = \sigma^2 I_{d \times d}$$

$$= \frac{1}{(2\pi \sigma^2)^{d/2}} \exp\left(-\frac{1}{2} x^T \frac{1}{\sigma^2} I x\right)$$

$$= \frac{1}{(2\pi \sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} x^T x\right)$$

$$= \frac{1}{(2\pi \sigma^2)^{d/2}} \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right)$$

$$P(\theta) = \frac{1}{(2\pi \sigma^2)^{d/2}} \exp\left(-\frac{\|\theta\|_2^2}{2\sigma^2}\right) =$$

arg min $\theta \sum_{i=1}^N -\log p(y_i | x_i; \theta) + \left(\frac{\|\theta\|_2^2}{2\alpha^2} \right) - \text{Regularizer}$

$\star \frac{\|x\theta - y\|_2^2 + \lambda \|\theta\|_2^2}{\lambda = \text{param.}}$

$\nabla_{\theta} =$

$(\theta^T x^T x \theta - 2y^T x \theta + y^T y) + \lambda \|\theta\|_2^2$
 $= (\theta^T x^T x \theta - 2y^T x \theta + \lambda \theta^T \theta)$
 $= \theta^T (x^T x + \lambda I) \theta - 2y^T x \theta$

$\nabla_{\theta} = 2\theta (x^T x + \lambda I) - 2x^T y = 0$

$\theta = (x^T x + \lambda I)^{-1} x^T y \mid \theta_{MLE} = (x^T x)^{-1} x^T y$
 MAP

arg min $\|x\theta - y\|_2^2 + \lambda \|\theta\|_2^2$
 $\lambda = \text{hyper parameter}$

L2 regular

Ridge Reg.

LASSO Regression / L1 Regularizer

$p(\theta_i) = \frac{1}{2b} \exp\left(-\frac{|\theta_i|}{b}\right)$

$p(\theta_i) \sim \text{laplan}(\theta_i, 0, b)$

$\log p(\theta_i) = -\frac{1}{b} |\theta_i| + C$

$d = M$

arg min $\|x\theta - y\|_2^2 + \left(\frac{1}{b}\right) \sum_{i=1}^d |\theta_i|$
 $\rightarrow \|\theta\|_1$

$\|x\theta - y\|_2^2 + \lambda \|\theta\|_1$