# Regularization

Goal: prefer "simpler" model parameters

## L0, L1, L2 Regularization:

Suppose: $\ell(\vec{\theta})$ is a likelihood fn.

Examples:

$$\ell(\vec{\theta}) \triangleq \log \prod_{i=1}^{N} p(x^{(i)}, y^{(i)} | \vec{\theta}) \longleftarrow \text{NB}$$

$$\ell(\vec{\theta}) \triangleq \log \prod_{i=1}^{N} p(y^{(i)} | x^{(i)}, \vec{\theta}) \longleftarrow \text{Log. Reg.}$$

$$\ell(\vec{\theta}) \triangleq \log \prod_{i=1}^{N} g(y^{(i)} | x^{(i)}, \vec{\theta}) \longleftarrow \text{Lin. Reg.}$$

Define:

$$\hat{\theta} = \underset{\theta}{\text{argmin}} \; J(\theta)$$
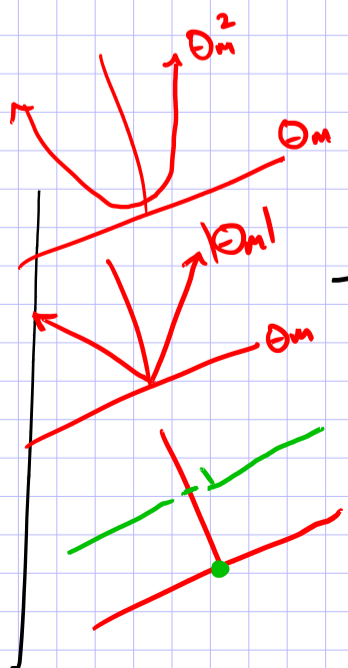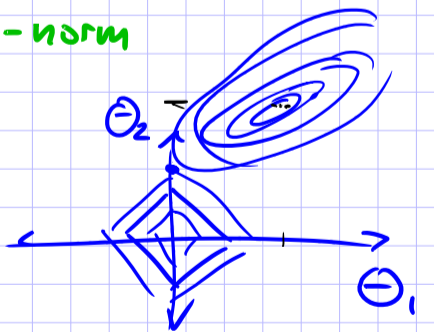
$$J(\theta) = -\ell(\theta) + \text{"model complexity"}$$

$$= \underbrace{-\ell(\theta)}_{\substack{\text{neg. log.} \\ \text{likelihood}}} + \underbrace{\lambda \, r(\theta)}_{\text{regularization}}$$

λ — tunable parameter chosen on validation

Key Idea: Define $r(\theta)$ s.t. we tradeoff between fitting the data and keeping the model simple.

## Choose form of $r(\vec{\theta})$:

Usually $r(\vec{\theta}) \triangleq \|\vec{\theta}\|_q$ — typically "p" i.e. p-norm

$$= \left[ \sum_{M=1}^{M} |\theta_m|^q \right]^{1/q}$$

| q | $r(\theta)$ | Preference for... | Notes |
|---|---|---|---|
| 2 | $\left(\|\theta\|_2\right)^2 = \sum \theta_m^2$ | Small values. | [L2 reg., differentiable] |
| 1 | $\|\theta\|_1 = \sum |\theta_m|$ | zero values. | [L1 reg., subdifferentiable] |
| 0 | $\|\theta\|_0 = \sum \mathbb{1}(\theta_m \, != 0)$ | zero values | [L0 reg., no good computational solutions] |

## Example: Linear Regression

$$r(\theta) = \left(\|\theta\|_2\right)^2 \Rightarrow L2 \text{ reg. a.k.a. "Ridge Regression"}$$

$$r(\theta) = \left(\|\theta\|_1\right) \Rightarrow L1 \text{ reg. aka "LASSO"}$$

$$r(\theta) = \left(\|\theta\|_0\right) \Rightarrow L0 \text{ reg. aka "Subset Selection"}$$

## Prob. Interp. of Regularization

Punchline: L2 reg. is MAP estimation w/ Gaussian prior
L1 reg. is MAP estimation w/ Laplace prior.

$$\ell_{MAP}(\theta) = \log\left[p(D|\theta)\, p(\theta)\right]$$

$$= \log p(D|\theta) + \log p(\theta)$$

### Ex: Zero-mean Gaussian prior on $\vec{\theta}$

Story:

$$\theta_m \sim \text{Gaussian}\left(\mu = 0, \sigma^2 = \frac{1}{2\lambda}\right) \quad \forall m$$

$$D \sim p(D|\theta)$$

$$\ell_{MAP}(\theta) = \log p(D|\theta) + \log\left[\prod_{m=1}^{M} f_{Gaussian}\left(\theta_m \mid \mu = 0, \sigma^2 = \frac{1}{2\lambda}\right)\right]$$

$$\hat{\theta} = \underset{\theta}{\arg\max}\ \ell_{MAP}(\vec{\theta})$$

$$= \underset{\theta}{\arg\max}\ \log p(D|\theta) + \log\left[\prod_{m=1}^{M} f_{Gaussian}\left(\theta_m \mid \mu = 0, \sigma^2 = \frac{1}{2\lambda}\right)\right]$$

$$= \underset{\theta}{\arg\max}\ \log p(D|\theta) + \sum_{m=1}^{M} \cancel{-\log\left(\sqrt{2\sigma^2 \pi}\right)} - \frac{1}{2\sigma^2}\left(\theta_m\right)^2$$

$$= \underset{\theta}{\arg\max}\ \log p(D|\theta) - \sum_{m=1}^{M} \frac{1}{2\sigma^2}\left(\theta_m\right)^2$$

$$= \underset{\theta}{\arg\max}\ \log p(D|\theta) - \lambda \sum_{m=1}^{M} \theta_m^2$$

L2 Reg.

Ex: Zero-mean Laplace prior on $\vec{\theta}$

Story:

$$\theta_m \sim \text{Laplace}\left(\mu=0, \; b=\frac{1}{\lambda}\right) \quad \forall m$$

$$D \sim p(D|\vec{\theta})$$

$$\Rightarrow \log p(\theta) \text{ is equiv. to } \lambda\|\theta_i\| \text{ penalty}$$