# Lecture 25: Deep neural networks continued

Professor Ilias Bilionis

# Regularization through parameter penalties
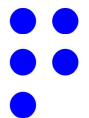
**PREDICTIVE SCIENCE LABORATORY**

# Regularization terms in loss functions

$$\mathcal{L}(\vartheta) = L(\vartheta) + \lambda R(\vartheta) + \mu R_L(\vartheta) + \ldots$$

loss funct.

regul. term

hyper-parameter
regularizing parameter.

$$+ \quad R(\vartheta) = \|\vartheta\|_2^2 = \sum_i \vartheta_i^2$$

$$+ \quad R(\vartheta) = \|\vartheta\|_1 = \sum_i |\vartheta_i|$$

$$\vdots$$

# Bayesian interpretation of regularization

$$\max_{\vartheta} p(y_{1:n} \mid x_{1:n}, \vartheta) \implies \min_{\vartheta} L(\vartheta)$$

prior over weights $p(\vartheta) \implies p(\vartheta \mid x_{1:n}, y_{1:n}) \propto p(y_{1:n} \mid x_{1:n}, \vartheta) p(\vartheta)$

MAP of $\vartheta$: $\max_{\vartheta} \log p(\vartheta \mid x_{1:n}, y_{1:n})$

$$J(\vartheta) = -\log p(\vartheta \mid x_{1:n}, y_{1:n}) = \underbrace{-\log p(y_{1:n} \mid x_{1:n}, \vartheta)}_{L(\vartheta)} \underbrace{-\log p(\vartheta)}_{\lambda R(\vartheta)}$$

Gaussian prior: $\vartheta \sim \mathcal{N}(0, \lambda^{-1})$

$$-\log p(\vartheta) = -\log \mathcal{N}(\vartheta \mid 0, \lambda^{-1})$$

$$= -\lambda \|\vartheta\|_2^2 + \text{const.}$$

$$\implies R(\vartheta) = \|\vartheta\|_2^2$$