# Lecture 25: Deep neural networks continued

Professor Ilias Bilionis

# Regularization through parameter penalties

*to avoid overfitting*

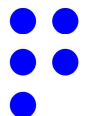# Regularization terms in loss functions

$$f(\vartheta) = \underline{L(\vartheta)} + \textcircled{$\lambda$}\ \underline{R(\vartheta)} + \mu\ R_1(\vartheta) + \dots$$

loss func.

regul. term

could have multiple

hyper-parameter
regularizing parameter.

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

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

good to default to

$$\vdots$$

# Bayesian interpretation of regularization

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

$$p(\vartheta) \implies \underbrace{p(\vartheta \mid x_{1:n}, y_{1:n})}_{\text{posterior}} \propto \underbrace{p(y_{1:n} \mid x_{1:n}, \vartheta)}_{\text{likelihood}} \underbrace{p(\vartheta)}_{\text{prior}}$$

prior over weights

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

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

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

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

$$= -\lambda \underbrace{\|\vartheta\|_2^2}_{R(\theta)} + \text{const.} \quad \text{w.r.t } \theta$$

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

Gaussian prior
↳ L2 regularization term

PREDICTIVE
SCIENCE LABORATORY