

Model selection and averaging

David J. H. Shih

Intended learning outcomes

- ▶ Apply model selection and model averaging to statistical models.

Predictive models

Use data (\mathbf{x}, \mathbf{y}) to estimate parameters $\boldsymbol{\theta}$ of the model and to predict new outcome \hat{y} given new data point \tilde{x} .

Discriminative model with MLE

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$$

$$\hat{y} = \operatorname{argmax}_y p(y \mid \tilde{x}, \hat{\boldsymbol{\theta}})$$

Discriminative model with MAP

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{x}) = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$$

$$\hat{y} = \operatorname{argmax}_y p(y \mid \tilde{x}, \hat{\boldsymbol{\theta}})$$

Generative model with MLE

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y})$$

$$\hat{y} = \operatorname{argmax}_y p(\tilde{x} \mid \hat{\boldsymbol{\theta}}, y)$$

Generative model with MAP

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{x}, \mathbf{y}) = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) p(\boldsymbol{\theta} \mid \mathbf{y})$$

$$\hat{y} = \operatorname{argmax}_y p(y \mid \tilde{x}, \hat{\boldsymbol{\theta}}) = \operatorname{argmax}_y p(\tilde{x} \mid \hat{\boldsymbol{\theta}}, y) p(y)$$

Full Bayesian discriminative model

$$p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{x}) = \frac{p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y} \mid \mathbf{x})}$$

$$p(\tilde{y} \mid \tilde{x}) = \int_{\Theta} p(\tilde{y} \mid \tilde{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{x}) d\boldsymbol{\theta}$$

Full Bayesian generative model

$$p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{x}) = \frac{p(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) p(\boldsymbol{\theta} \mid \mathbf{y})}{p(\mathbf{x} \mid \mathbf{y})}$$

$$p(\tilde{x} \mid \tilde{y}) = \int_{\Theta} p(\tilde{x} \mid \boldsymbol{\theta}, \tilde{y}) p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{x}) d\boldsymbol{\theta}$$

$$p(\tilde{y} \mid \tilde{x}) = \frac{p(\tilde{x} \mid \tilde{y}) p(\tilde{y})}{p(\tilde{x})}$$

Summary

"All Models Are Wrong, Some are Useful." - George Box

Intended learning outcomes

- ▶ Apply model selection and model averaging to statistical models.