

Bias/Variance and Regularization

Fundamentals of Data Science
12, 15 December 2022
Prof. Fabio Galasso

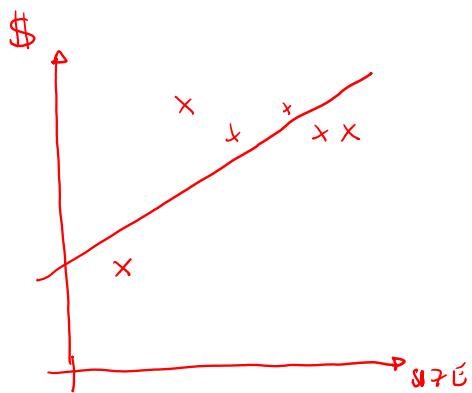


Outline

- Bias/Variance
- Regularization
- Hold-out Cross Validation

Bias/Variance

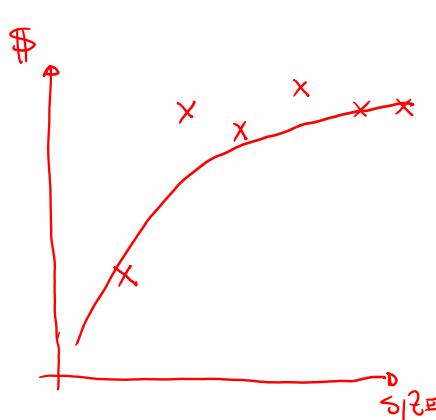
Bias/Variance Regression



$$\theta_0 + \theta_1 x$$

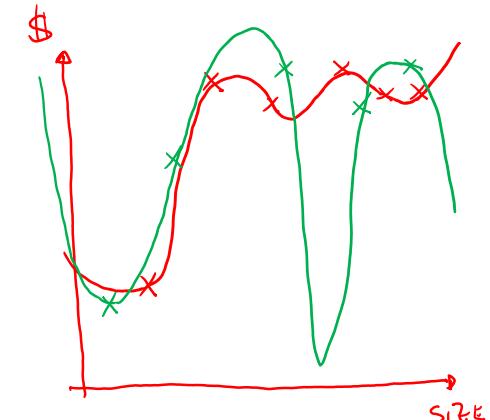
UNDERFIT

HIGH BIAS



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"JUST RIGHT"

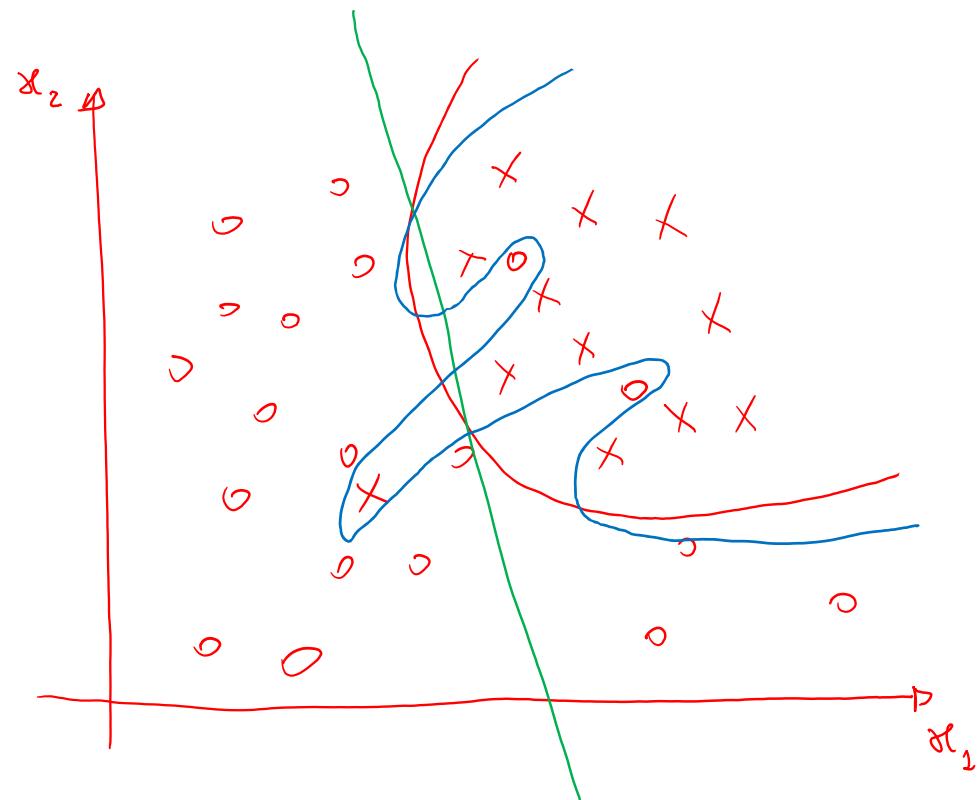


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_5 x^5$$

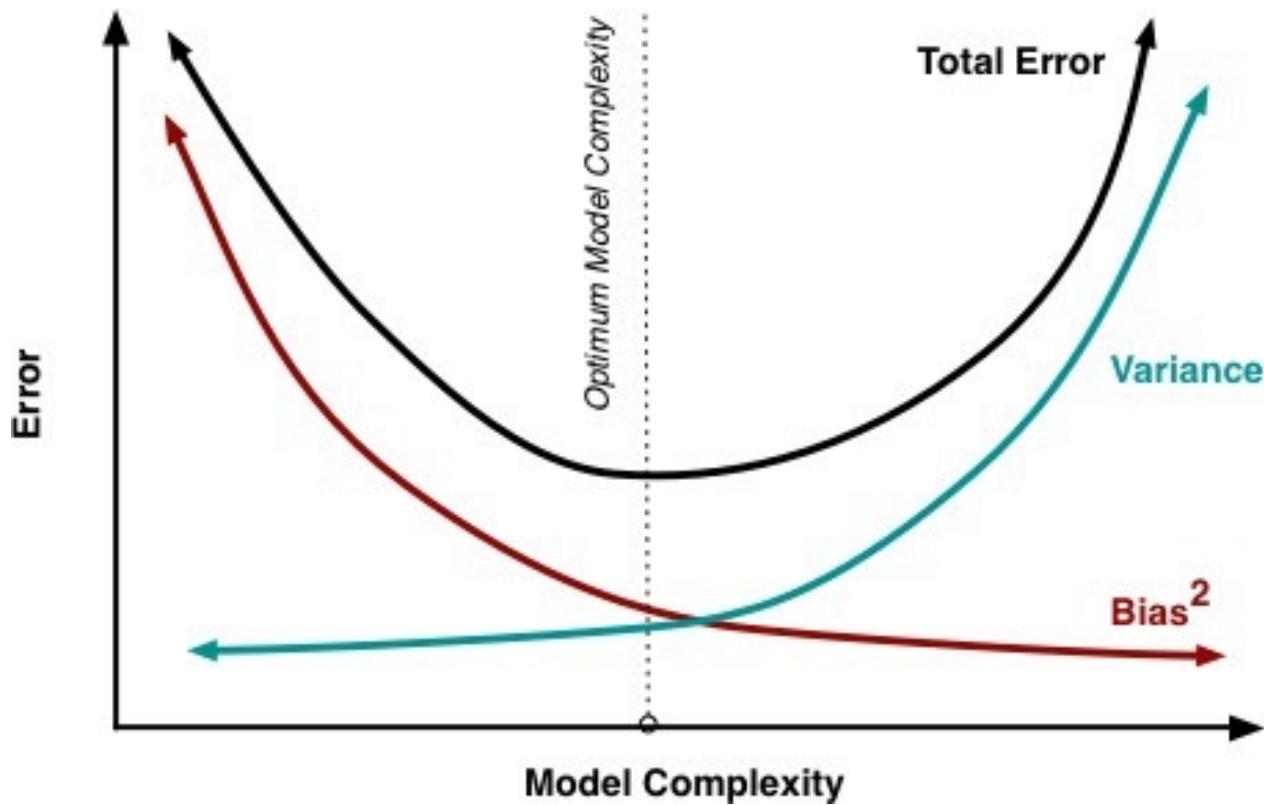
OVERFIT

HIGH VARIANCE

Bias/Variance Classification



Bias-Variance Tradeoff



<http://scott.fortmann-roe.com/docs/BiasVariance.html>

Setup/Assumptions

1) DATA DISTRIBUTION

$$(\mathbf{x}, y) \sim D \quad \begin{matrix} \text{TRAIN} \\ \text{TEST} \end{matrix}$$

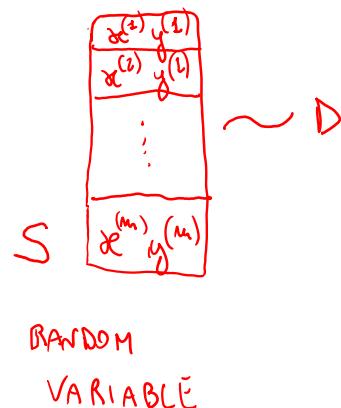
2) INDEPENDENT SAMPLES

Θ^* or h^* "TRUE"

PARAMETERS

(NOT RANDOM)

PROCESS OF
LEARNING



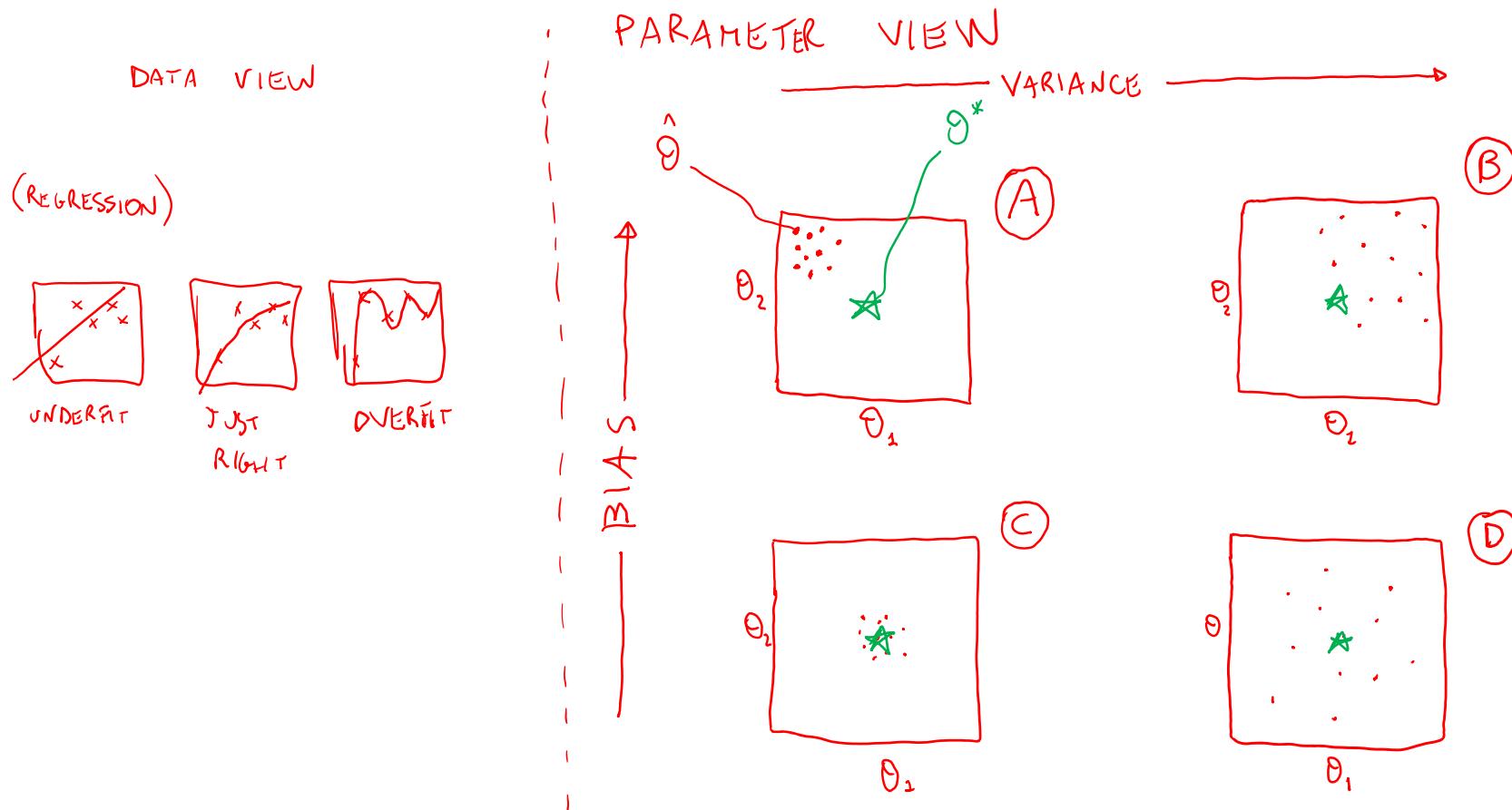
\hat{h} or \hat{g} ~ SAMPLING DISTN.

DETERMINISTIC
FUNCTION

RANDOM
VARIABLE

Bias/Variance

Data and Parameter View



Consistency, Statistical Efficiency, Bias

$$m \rightarrow \infty$$

$$\text{VAR}[\hat{\theta}] \rightarrow 0$$

E.g. $\mathcal{O}\left(\frac{1}{m^2}\right) \Theta(e^{-m})$

"STATISTICAL EFFICIENCY": RATE $\text{VAR}[\hat{\theta}] \rightarrow 0$ AS $m \rightarrow \infty$

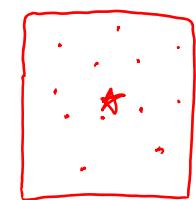
$$\hat{\theta} \rightarrow \theta^* \quad m \rightarrow \infty : \text{CONSISTENT}$$

$$E[\hat{\theta}] = \theta^* \quad \text{for all } m ; \text{ UNBIASED}$$

Fighting High Variance

1) $m \rightarrow \infty$

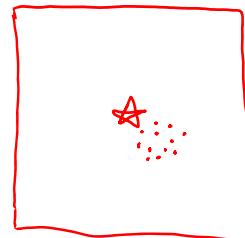
2) REGULARIZATION



LOW BIAS

HIGH VARIANCE

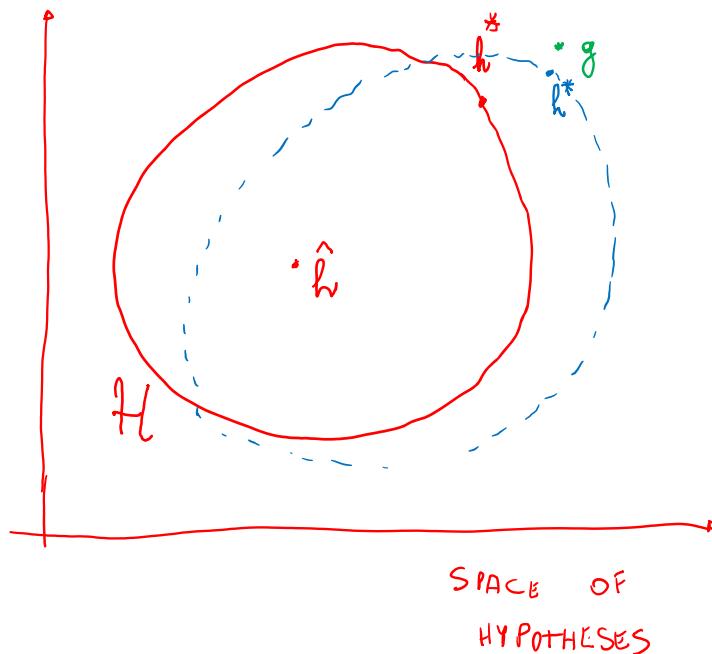
Reg



SMALL BIAS

LOW VARIANCE

Generalization and Empirical Risk



g BEST POSSIBLE HYPOTHESIS

h^* BEST IN CLASS H

\hat{h} LEARNT FROM FINITE DATA

$$E(h) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [\mathbb{1}\{h(\mathbf{x}) \neq y\}]$$

RISK
GENERALIZATION ERROR

$$\hat{E}_s(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(\mathbf{x}^{(i)}) \neq y^{(i)}\}$$

EMPIRICAL
RISK

Approximation and Estimation Error

$$\mathcal{E}(g) = \text{BAYES ERROR} / \text{IRREDUCIBLE ERROR}$$

$$\mathcal{E}(h^*) - \mathcal{E}(g), \text{ APPROXIMATION ERROR} \quad \} \text{ CLASS}$$

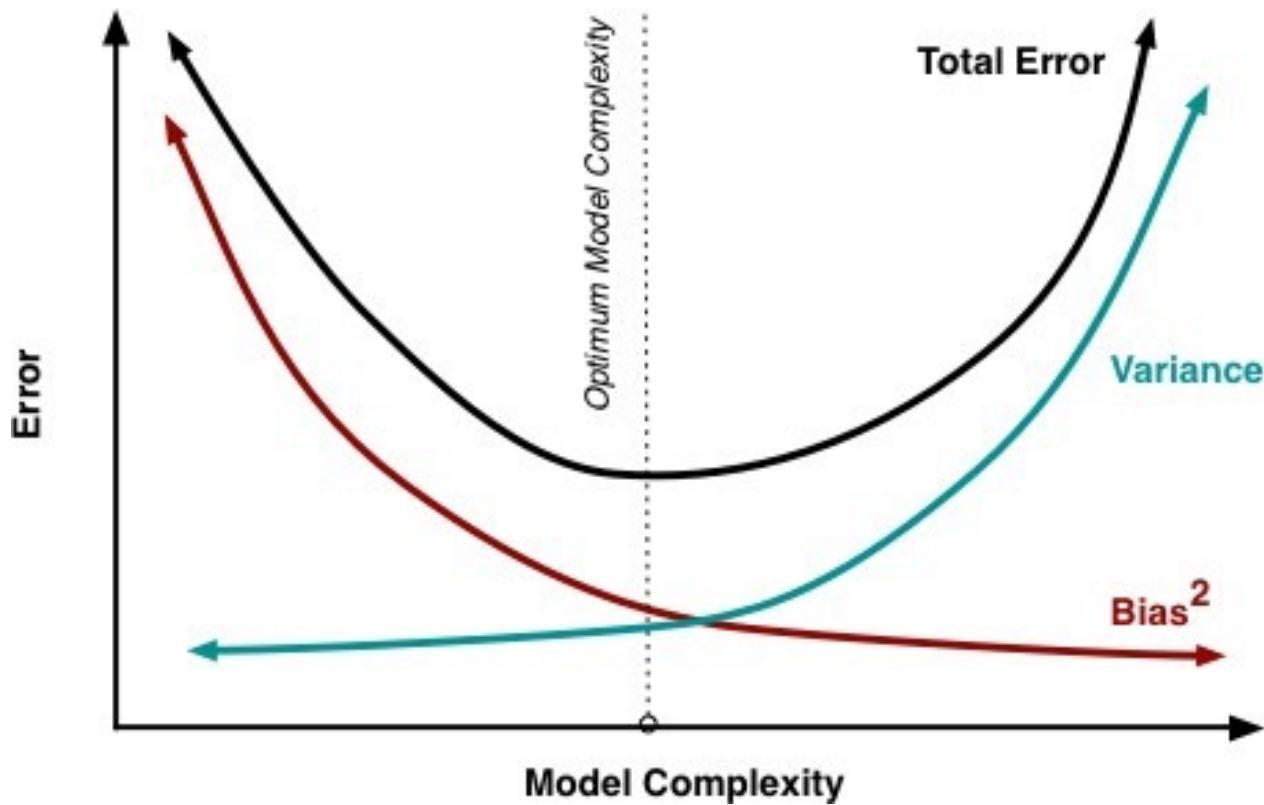
$$\mathcal{E}(\hat{h}) - \mathcal{E}(h^*), \text{ ESTIMATION ERROR} \quad \} \text{ DATA}$$

$$\mathcal{E}(\hat{h}) = \text{EST. ERROR} + \text{APPROX. ERROR} + \text{IRREDUCIBLE ERROR}$$

$$= \underbrace{\text{EST. VAR}}_{\text{VARIANCE}} + \underbrace{\text{EST. BIAS} + \text{APPROX. ERROR}}_{\text{BIAS}} + \underbrace{\text{IRREDUCIBLE}}_{\text{ERROR}}$$

$$\mathcal{E}(\hat{h}) \approx \text{VARIANCE} + \text{BIAS} + \text{IRREDUCIBLE}$$

Bias-Variance Tradeoff



<http://scott.fortmann-roe.com/docs/BiasVariance.html>

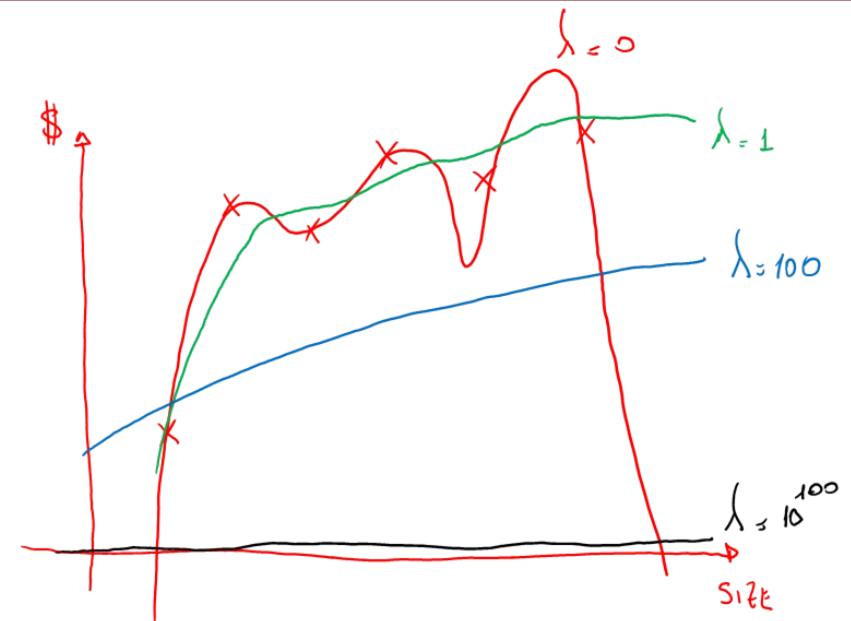
Outline

- Bias/Variance
- Regularization
- Hold-out Cross Validation

Regularization

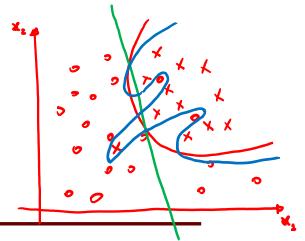
Regularization

$$\underset{\theta}{\text{MIN}} \quad \frac{1}{2} \sum_{i=1}^m \|y^{(i)} - \theta^T x^{(i)}\|^2 + \underbrace{\frac{\lambda}{2} \|\theta\|^2}_{\text{REGULARIZATION}}$$



$$\underset{\theta}{\text{ARGMAX}} \quad \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}, \theta) - \lambda \|\theta\|^2 \quad h_{\theta}(x) \approx \phi$$

Regularization and Gradient Descent



$$\text{ARGMAX}_{\theta} \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) - \frac{\lambda}{2} \|\theta\|^2 \rightsquigarrow -\frac{\lambda}{2} \sum_{j=1}^m \theta_j^2$$

NO REGULARIZATION
FOR $j=0$

$$l(\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) - \frac{\lambda}{2} \sum_{j=1}^m \theta_j^2$$

$$\theta_j = \theta_j + \alpha \left[\sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)} - \lambda \theta_j \right] \quad j \in [1, 2, \dots, m]$$

$\underbrace{\qquad\qquad\qquad}_{\frac{\partial}{\partial \theta_j} l(\theta)}$

FOR LOGISTIC REGRESSION: $h_{\theta}(x^{(i)}) = g(\theta^T x^{(i)})$

$$\theta_0 = \theta_0 + \alpha \sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right)$$

SINCE $x_0^{(i)} = 1$ AND NO REG FOR θ_0

ALSO APPLICABLE TO LINEAR REG $h_{\theta}(x^{(i)}) = \theta^T x^{(i)}$

Application to Text Classification

$M = 100$ EXAMPLES

$n = d = 10'000$

$$x_i = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{adam} \quad \left. \begin{array}{c} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{array} \right\} d$$

! torro

Regularization: Probabilistic Interpretation

$$S = \left\{ \mathbf{x}^{(i)}, y^{(i)} \right\}_{i=1}^m$$

$$P(\theta | s) = \frac{P(s|\theta) P(\theta)}{P(s)}$$

$$\underset{\theta}{\text{ARGMAX}} \quad P(\theta | s) = \underset{\theta}{\text{ARGMAX}} \quad P(s|\theta) P(\theta) = \underset{\theta}{\text{ARGMAX}} \left(\prod_{i=1}^m P(y^{(i)} | \mathbf{x}^{(i)}; \theta) \right) P(\theta)$$

$$P(\theta): \quad \theta \sim \mathcal{N}(0, \sigma^2 I) \quad P(\theta) = \frac{1}{\sqrt{(2\pi)^m |\sigma^2 I|^{\frac{m}{2}}}} \exp\left(-\frac{1}{2} \theta^\top (\sigma^2 I)^{-1} \theta\right)$$

LOGISTIC REGRESSION

FREQUENTIST: $\underset{\theta}{\text{ARGMAX}} \quad P(s|\theta)$ MLE

BAYESIAN: PRIOR DISTN. $P(\theta)$ $\underset{\theta}{\text{ARGMAX}} \quad P(\theta | s)$

MAP MAXIMUM A
POSTERIORI

Generative vs. Discriminative; Bayesian vs. Frequentist

$y \rightarrow$ CLASS

$x \rightarrow$ INPUT FEATURES

$\theta \rightarrow$ PARAMETERS

FREQUENTIST

BAYESIAN

DISCRIMINATIVE

$$P(y; x, \theta)$$

$$P(y, \theta; x) = P(y | \theta; x) P(\theta)$$

GENERATIVE

$$P(y, x; \theta)$$

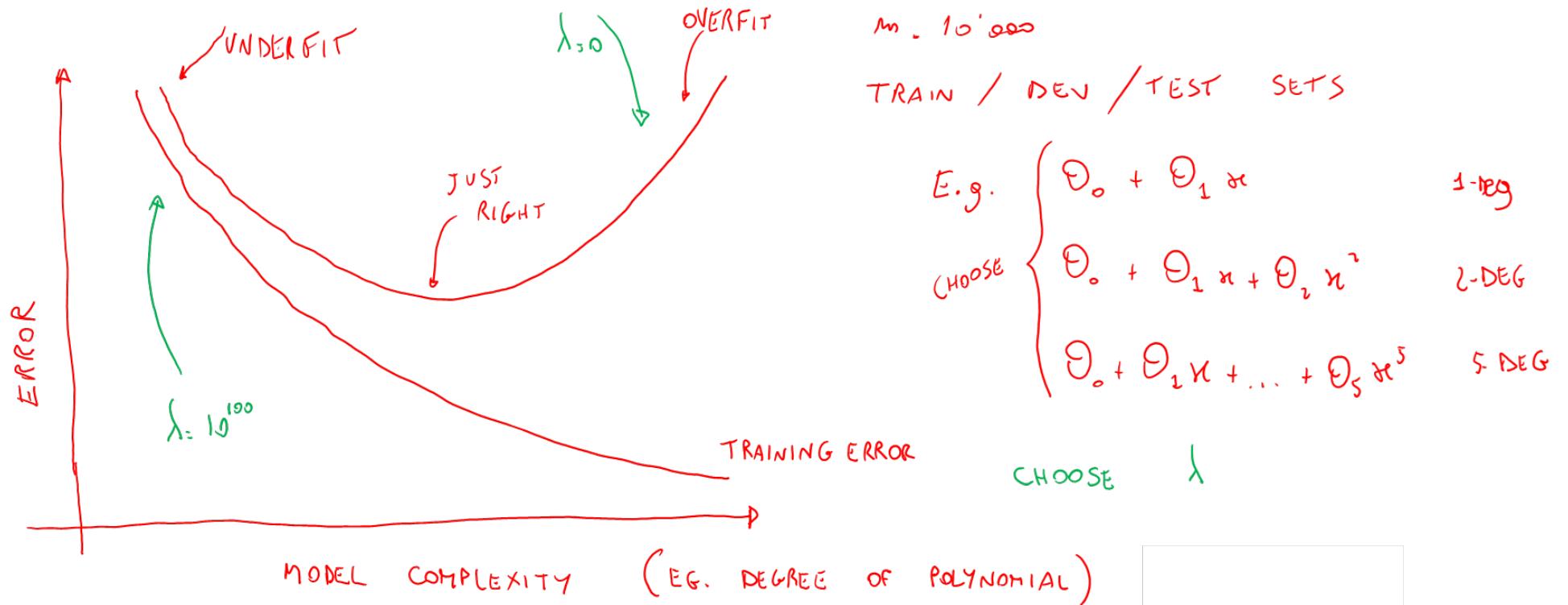
$$P(y, x, \theta) = P(y, x | \theta) P(\theta)$$

Outline

- Bias/Variance
- Regularization
- Hold-out Cross Validation

Hold-out Cross Validation

Model Complexity Vs Error



Train/Dev/Test

1. Split S into S_{train} , S_{dev} , S_{test}

DEV = DEVELOPMENT / CROSS VALIDATION

2. Train each model S_{train}

GET HYPOTHESES h_θ 'S E.G. VARIOUS DEGREES OF POLYNOMIAL

3. Choose the model with lowest error on S_{dev}

NOT ON S_{train}

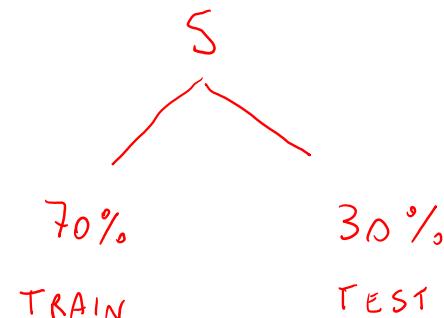
- Optional: evaluate on S_{test} and report performance

NOT ON S_{dev}

(SIMPLE) HOLD-OUT CROSS VALIDATION

DEGREE	S_{dev} ERROR	S_{test} ERROR
1	10	10
2	5.1	5.0
3	5.0	5.0
4	4.9	5.0
5	?	?
6	10	10

Hold-out Cross Validation



60% 20% 20
TRAIN DEV TEST

10'000'000

TRAIN 6,000,000
DEV 2,000,000
TEST 2,000,000

Setting Hyperparameters

Idea #1: Choose hyperparameters
that work best on the data

Your Dataset

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: would always choose most complex model

Your Dataset

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: would always choose most complex model

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: would always choose most complex model

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: would always choose most complex model

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Idea #3 (hold-out cross validation): Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!

train

validation

test

Setting Hyperparameters

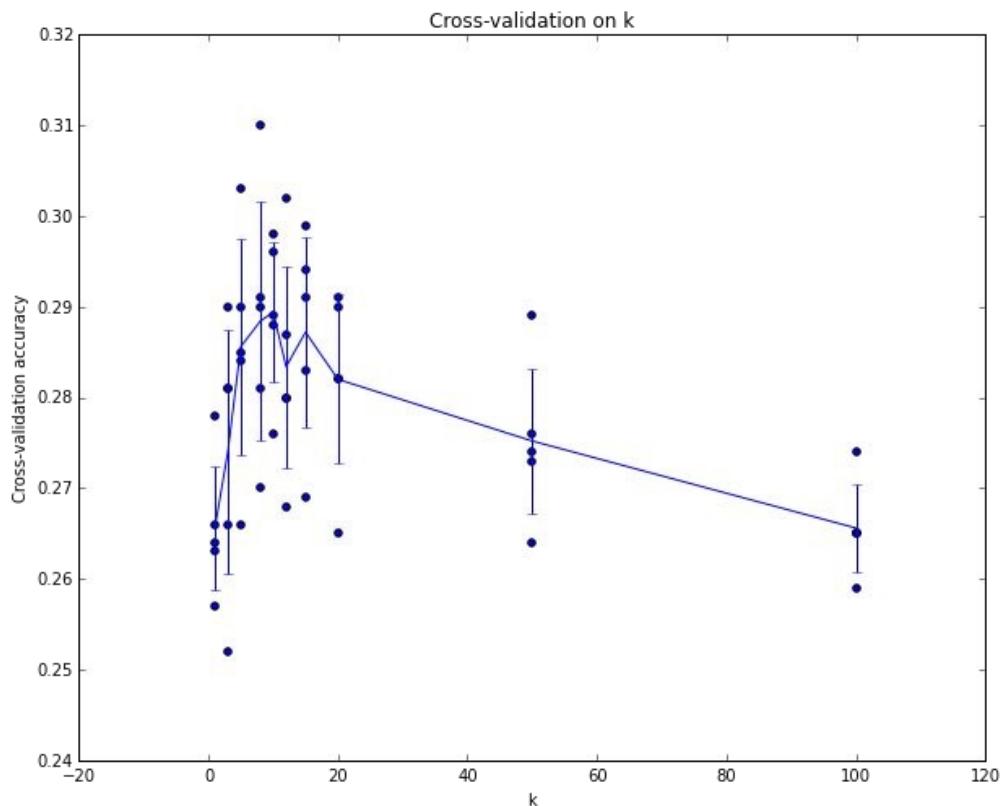
Your Dataset

Idea #4: k-fold cross-validation: Split data into **folds**,
try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but not used too frequently in deep learning

Setting Hyperparameters



Example of
5-fold cross-validation
for the value of k .

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)

References

- Sections 1.2.3, 1.3 and 3.2 in [Bishop, 2006. Pattern Recognition and Machine Learning]
- Sections 5-12, 20-32 in [Andrew Ng, 2019. Machine Learning Yearning]

Thank you

Acknowledges: slides and material from Andrew Ng, Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau

