

Boston, Sep 19th, 2015

ACTING ON DATA

Rahul Dave (rahuldave@gmail.com)

@rahuldave

\mathcal{L}_X Prior

\mathcal{L}_X Prior

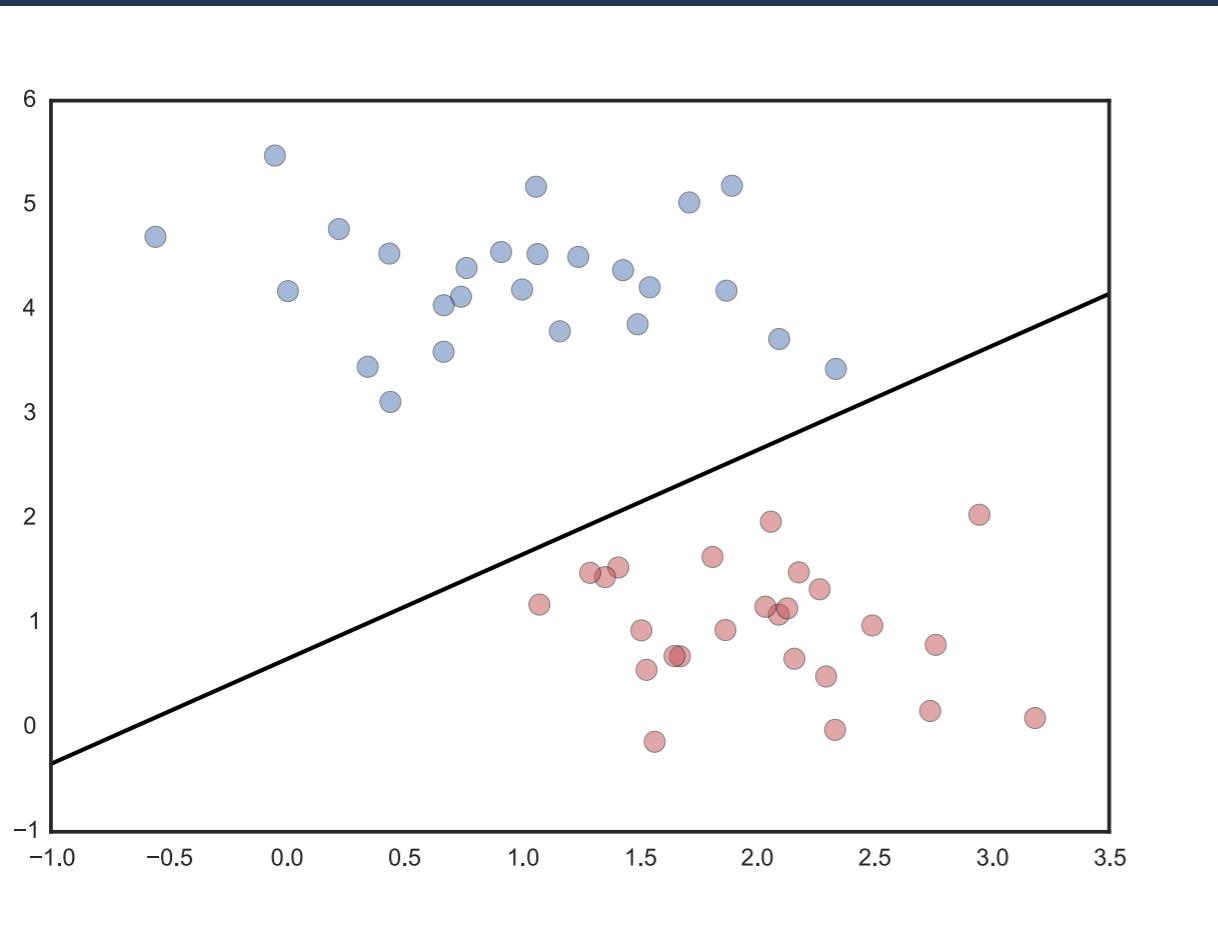
Data Science.
Simulation.
Software.
Social Good.
Interesting Problems.

machine learning, complex systems, stochastic methods, viz, extreme computing

DEGREE PROGRAMS:

- Master's of science- one year
- Master's of engineering - two year with thesis/research project

CLASSIFICATION

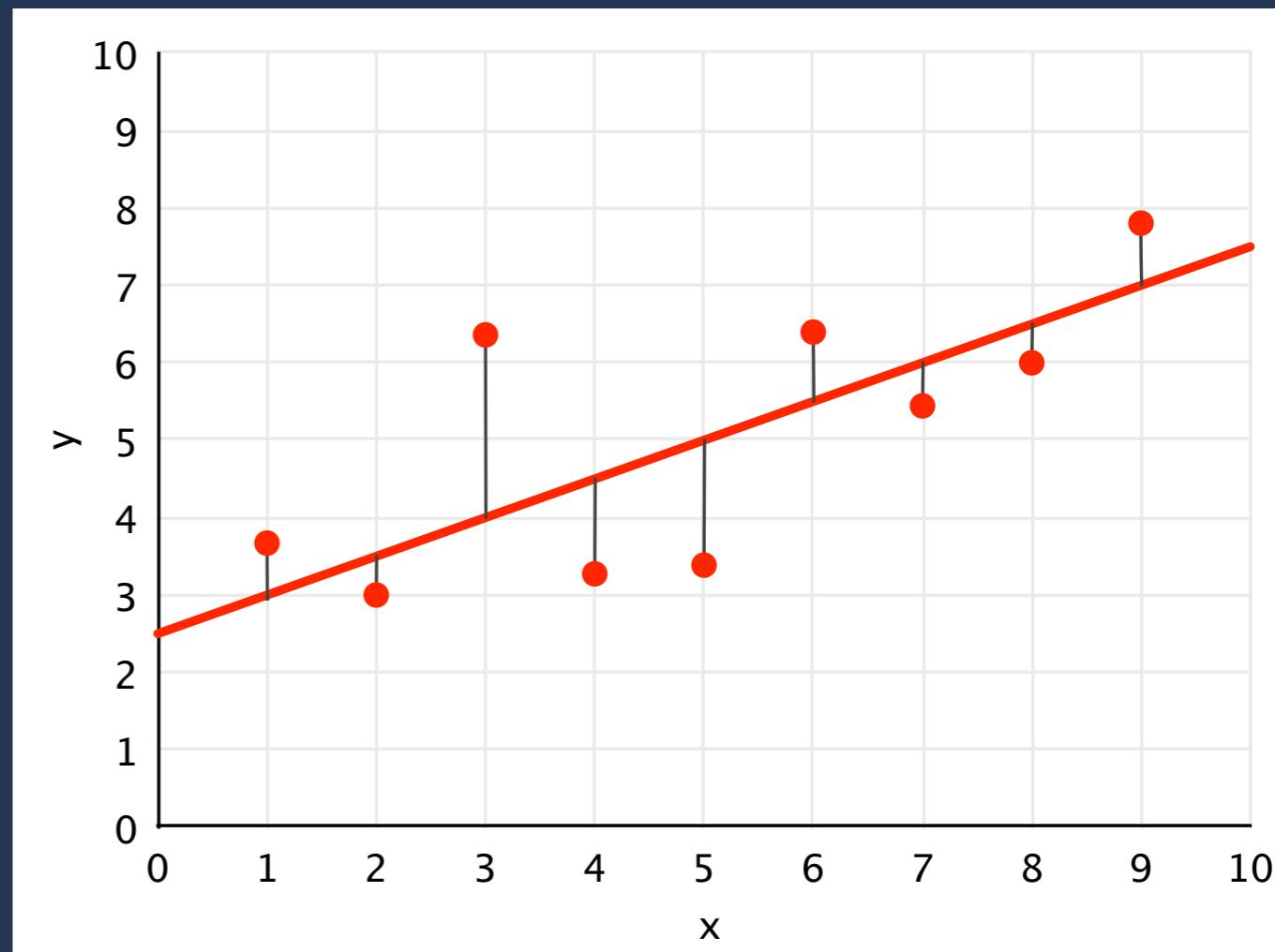


- will a customer churn?
- is this a check? For how much?
- a man or a woman?
- will this customer buy?
- do you have cancer?
- is this spam?
- whose picture is this?
- what is this text about?^j

^j image from code in <http://bit.ly/1Azg29G>

REGRESSION

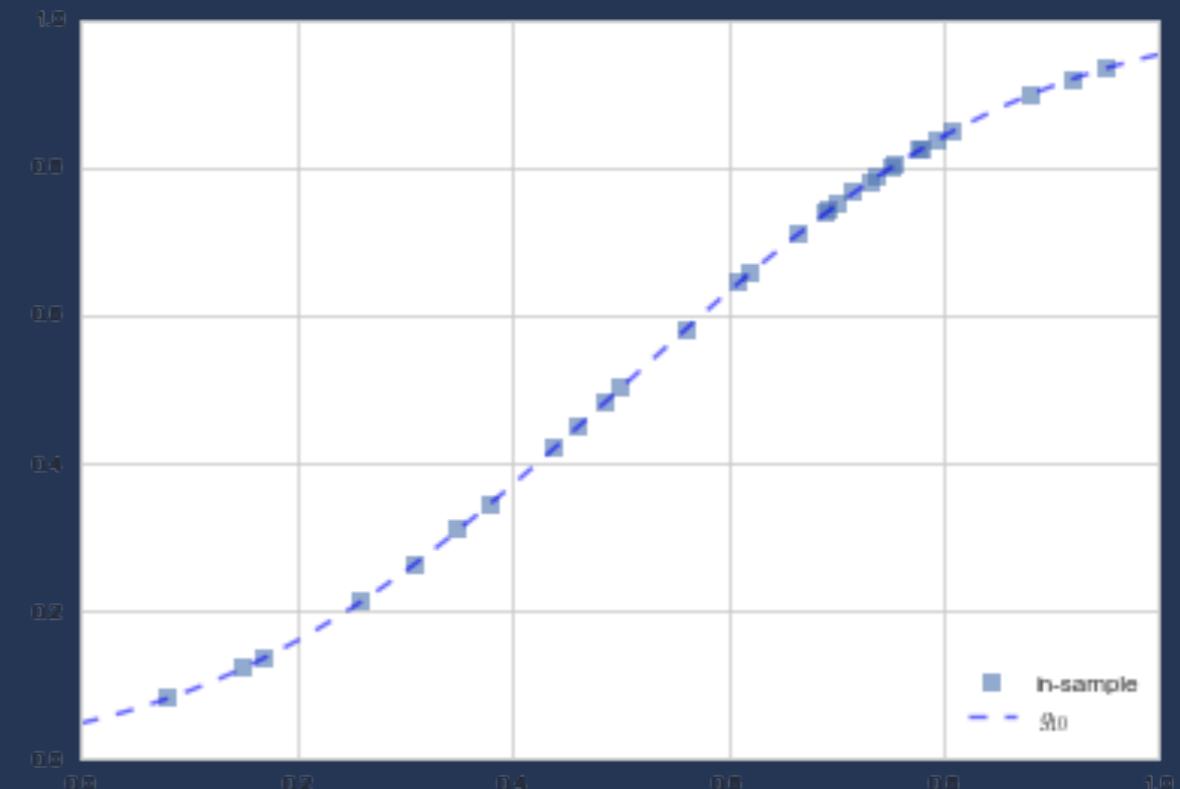
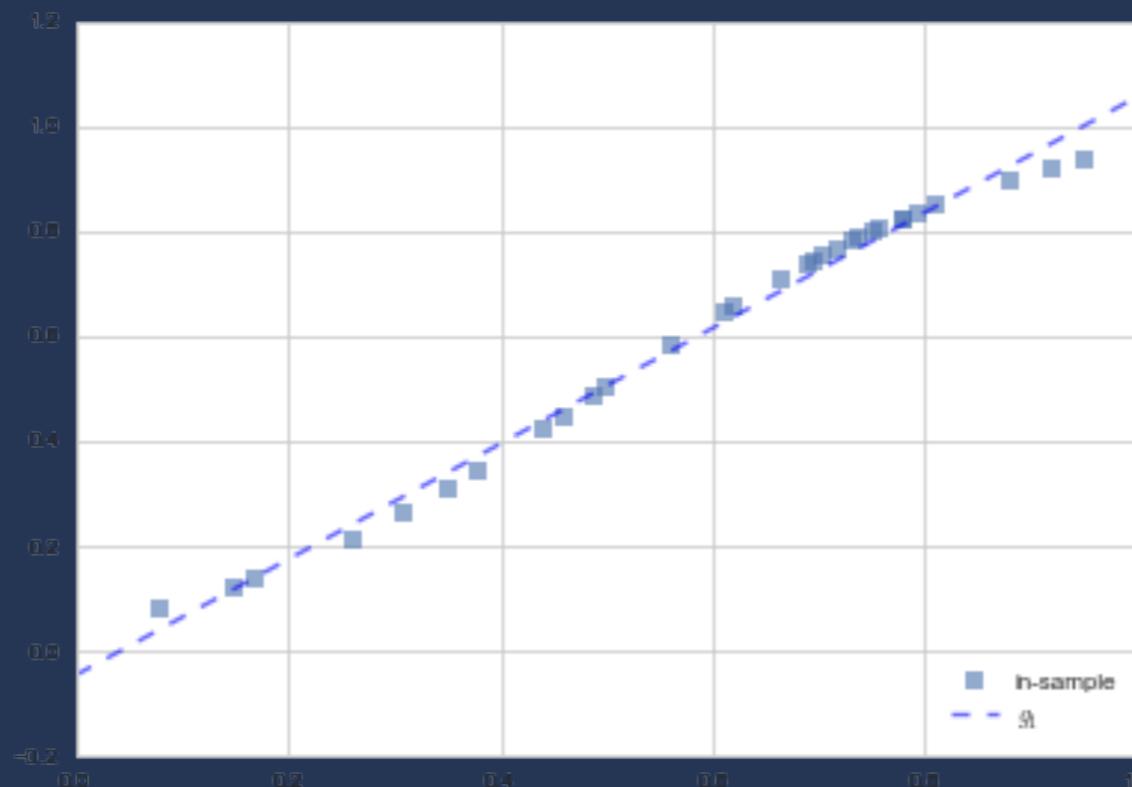
- how many dollars will you spend?
- what is your creditworthiness
- how many people will vote for Bernie t days before election
- use to predict probabilities for classification
- causal modeling in econometrics



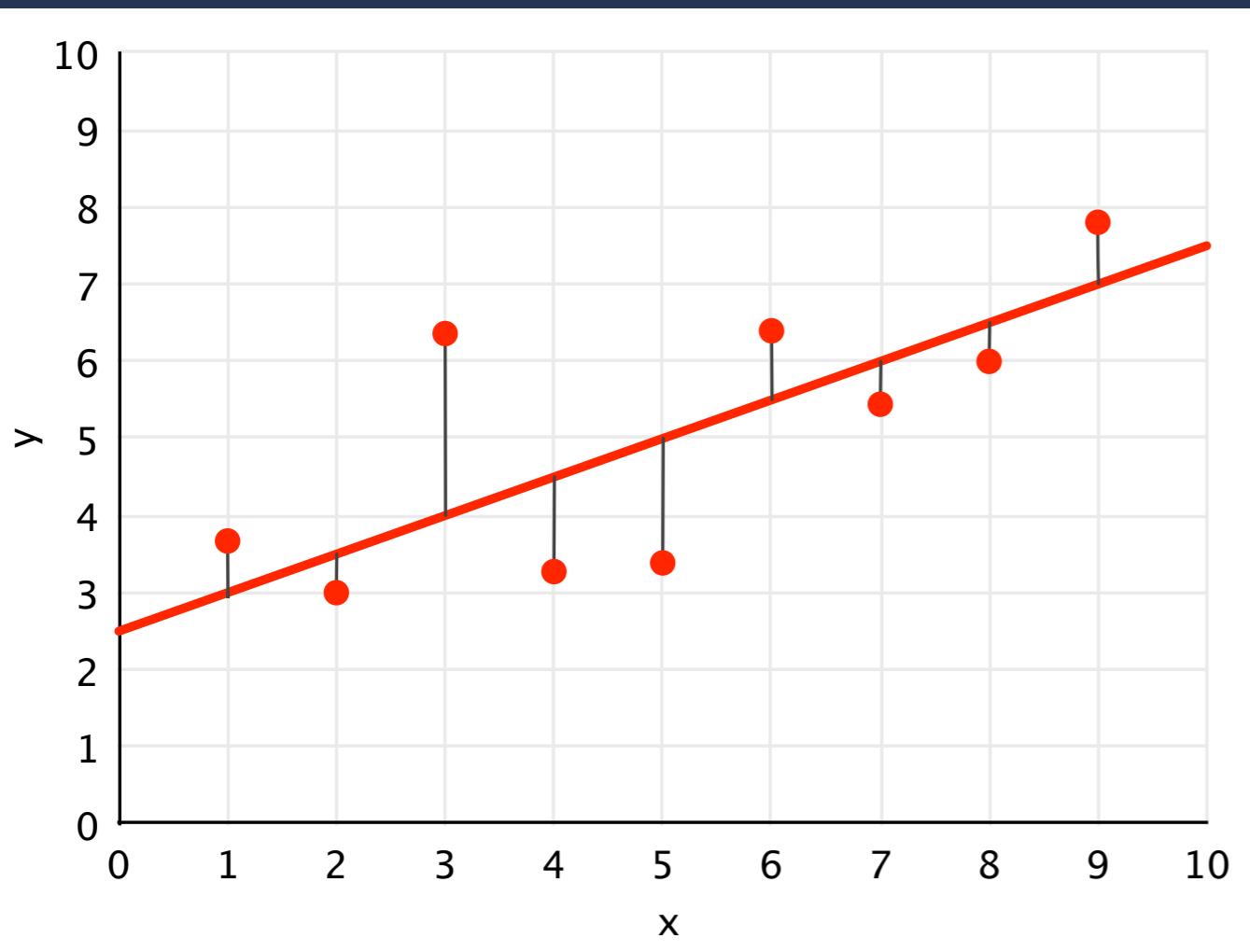
From Bayesian Reasoning and Machine Learning, David Barber:

"A father decides to teach his young son what a sports car is. Finding it difficult to explain in words, he decides to give some examples. They stand on a motorway bridge and ... the father cries out 'that's a sports car!' when a sports car passes by. After ten minutes, the father asks his son if he's understood what a sports car is. The son says, 'sure, it's easy'. An old red VW Beetle passes by, and the son shouts - 'that's a sports car!'. Dejected, the father asks - 'why do you say that?'. 'Because all sports cars are red!', replies the son."

30 points of data. Which fit is better? Line in \mathcal{H}_1 or curve in \mathcal{H}_{20} ?



What does it mean to FIT?



Minimize distance from the line?

$$R_{\mathcal{D}}(h_1(x)) = \frac{1}{N} \sum_{y_i \in \mathcal{D}} (y_i - h_1(x_i))^2$$

Minimize squared distance from the line.

$$g_1(x) = \arg \min_{h_1(x) \in \mathcal{H}} R_{\mathcal{D}}(h_1(x)).$$

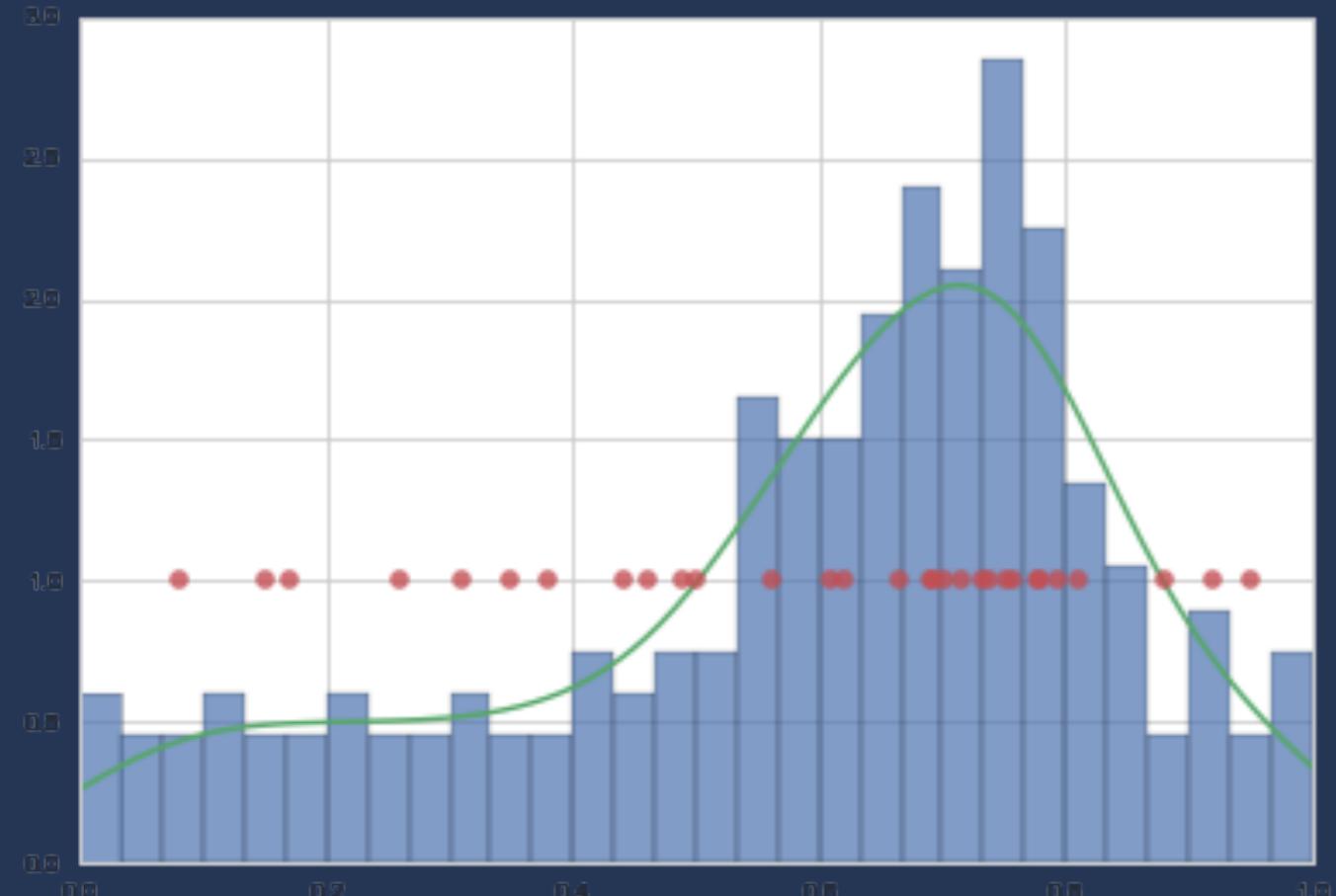
Get intercept w_0 and slope w_1 .

EMPIRICAL RISK MINIMIZATION

The sample must be representative of the population!

$$A : R_{\mathcal{D}}(g) \text{ smallest on } \mathcal{H}$$
$$B : R_{\text{out of sample}}(g) \approx R_{\mathcal{D}}(g)$$

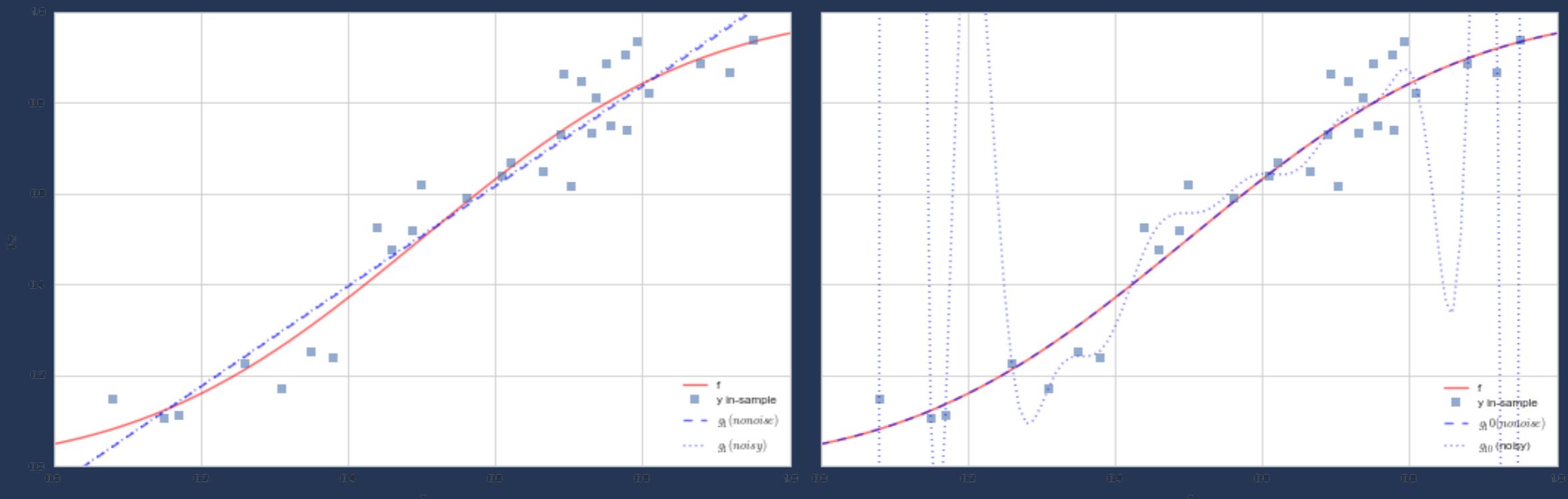
- A: Empirical risk estimates out of sample risk.
B: Thus the out of sample risk is also small.



THE REAL WORLD HAS NOISE

Which fit is better now?

The line or the curve?



10

9

8

7

6

5

4

3

2

1

0

y

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

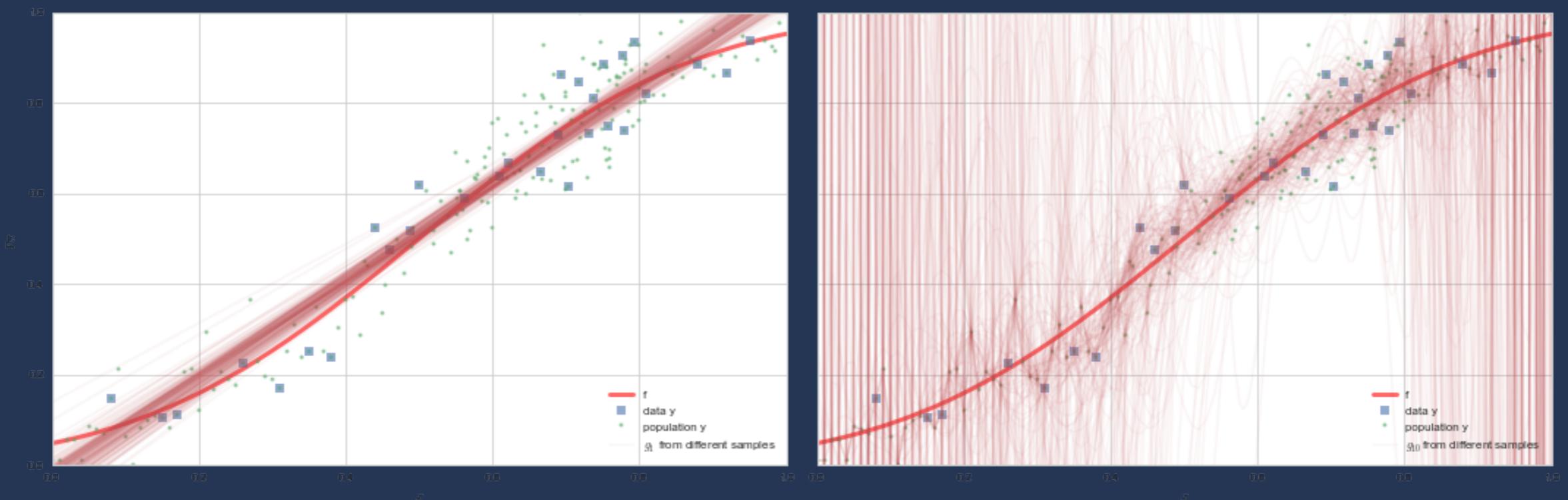
257

258

259

260

UNDERFITTING (Bias) vs OVERFITTING (Variance)



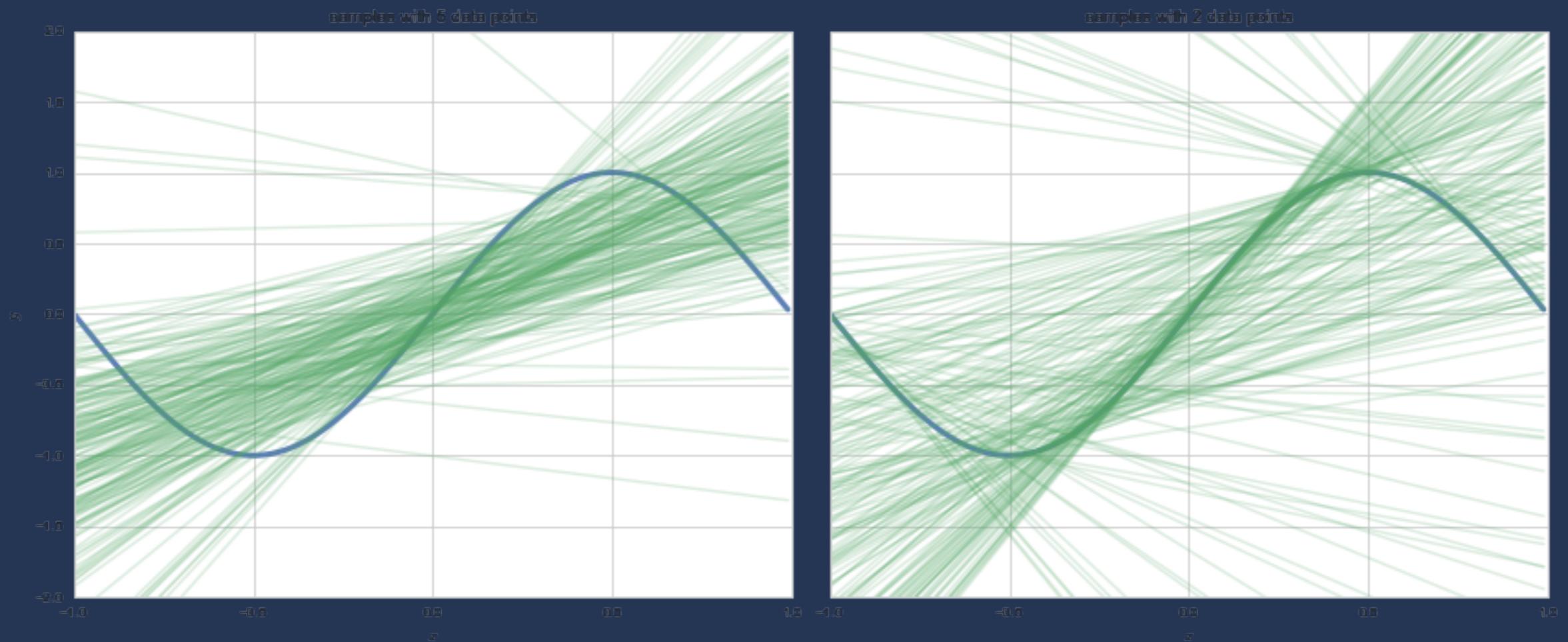
A background scatter plot featuring two classes of data points: red circles and blue circles. Overlaid on the plot are several green and black wavy lines of varying complexity, representing different machine learning models. Some lines closely follow the general trend of the data points, while others are extremely complex, oscillating wildly between the points, which is a visual representation of overfitting.

If you are having problems in machine learning the reason is almost always

OVERFITTING

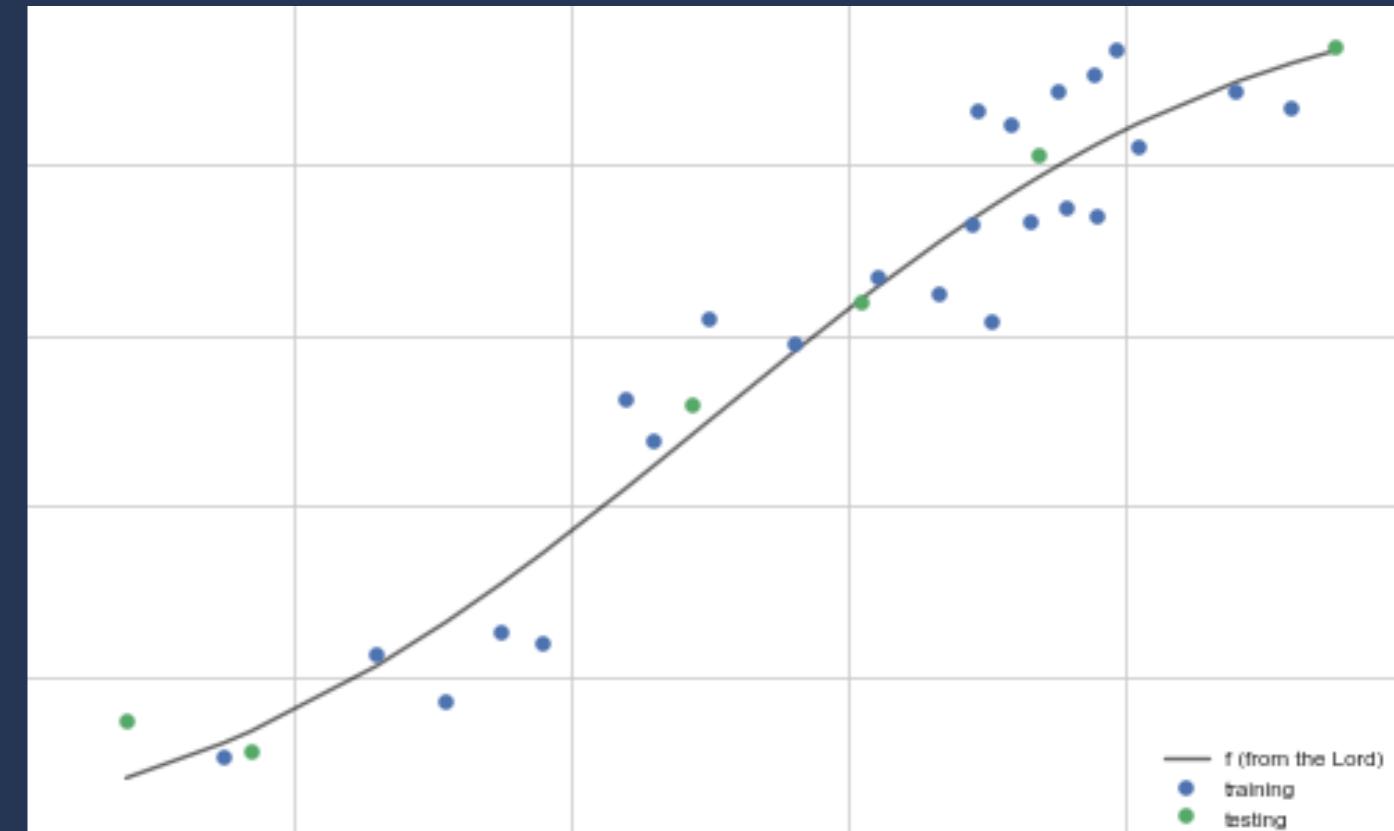
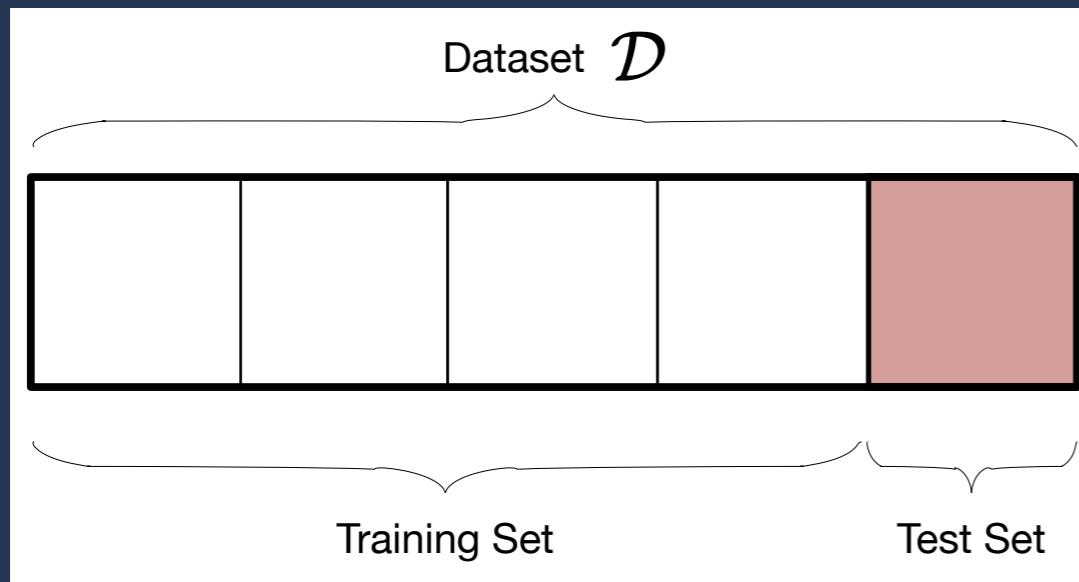
- Background from <http://commons.wikimedia.org/wiki/File:Overfitting.svg>

DATA SIZE MATTERS: straight line fits to a sine curve

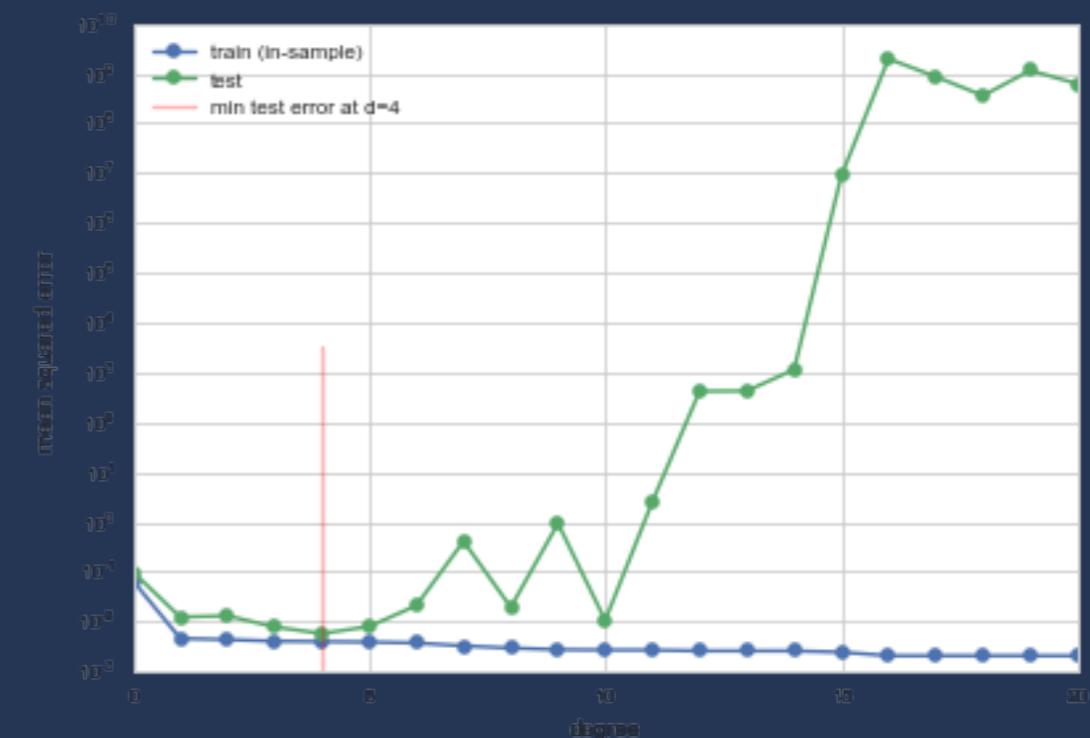
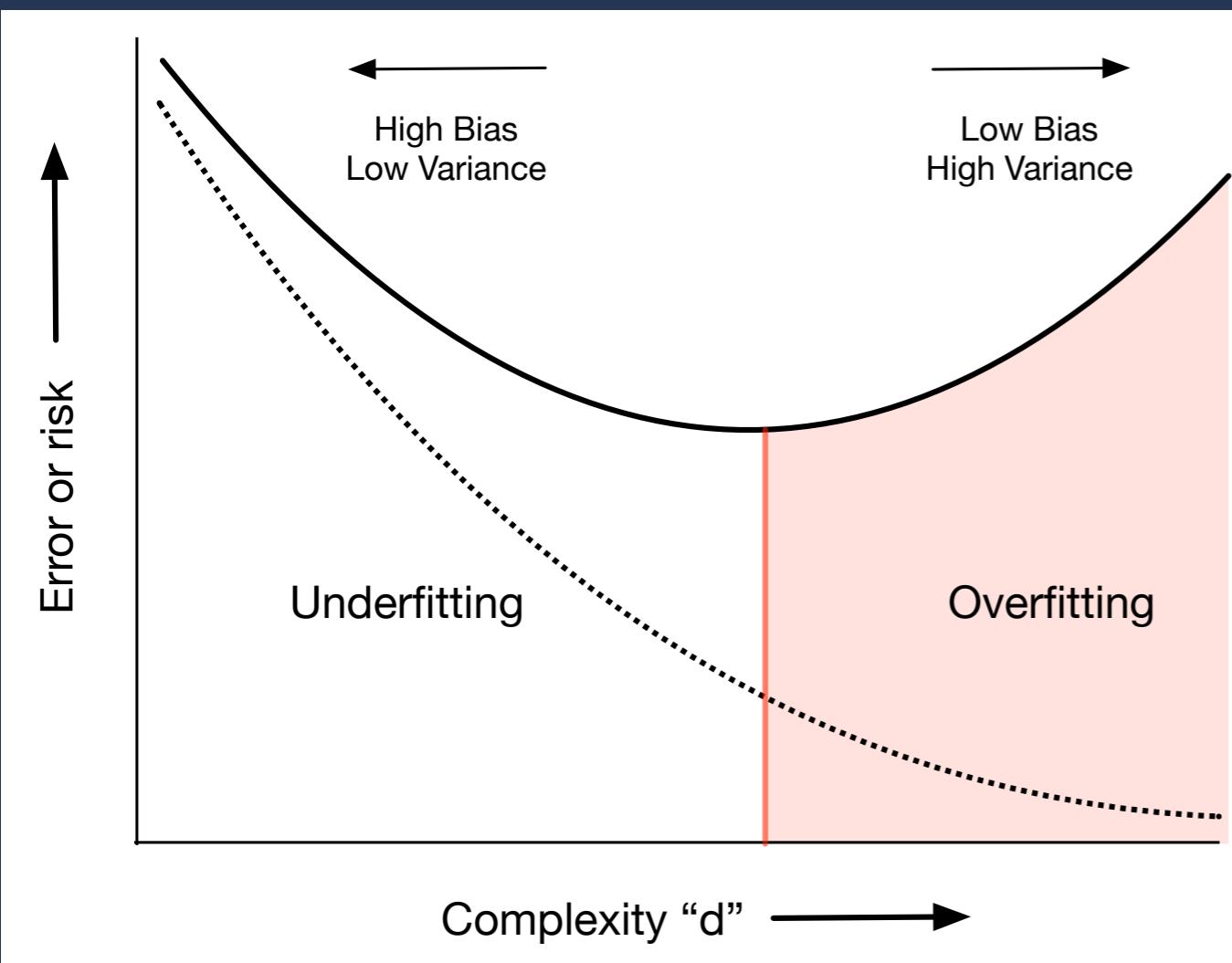


Corollary: Must fit simpler models to less data!

HOW DO WE LEARN?

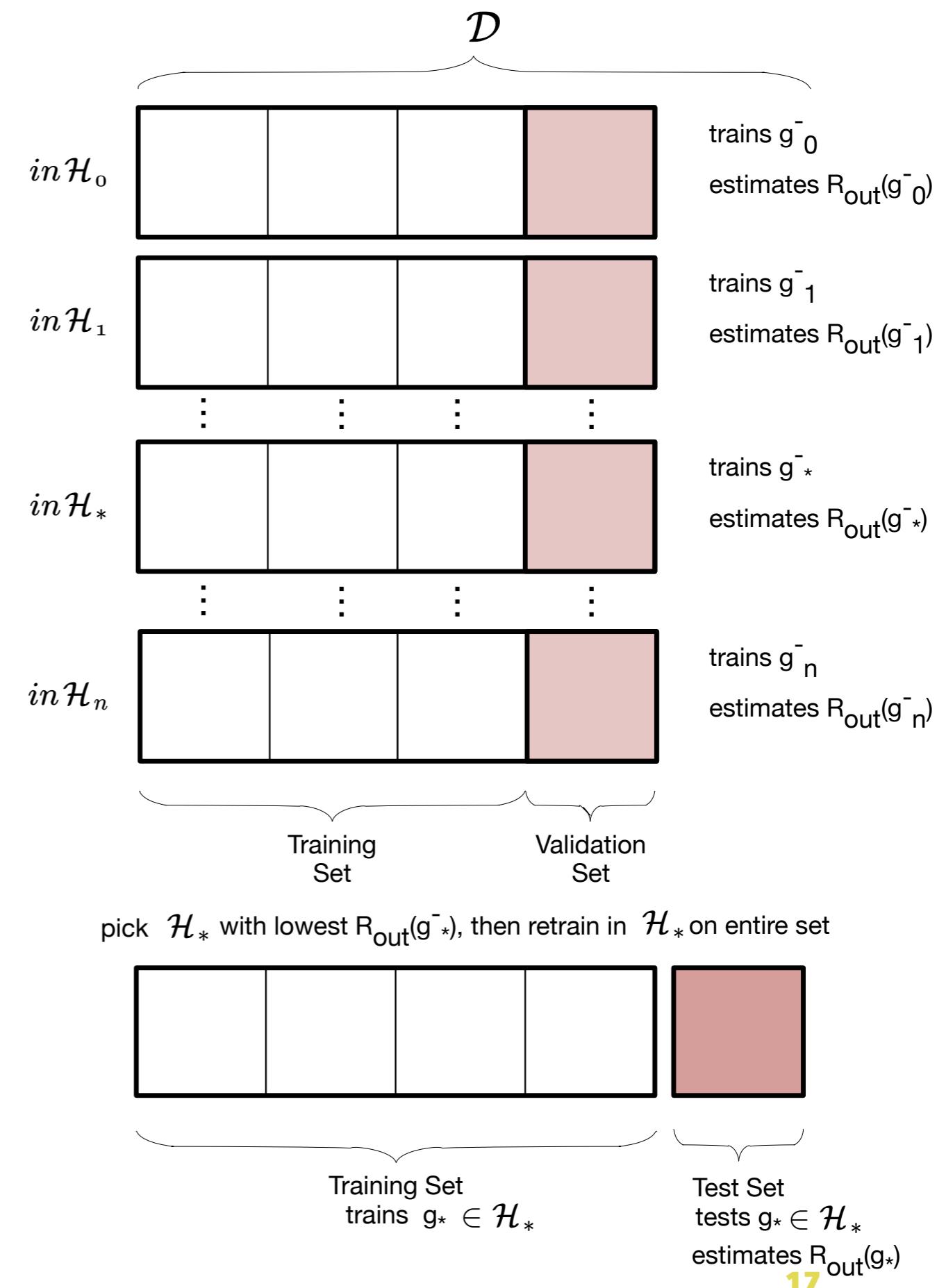
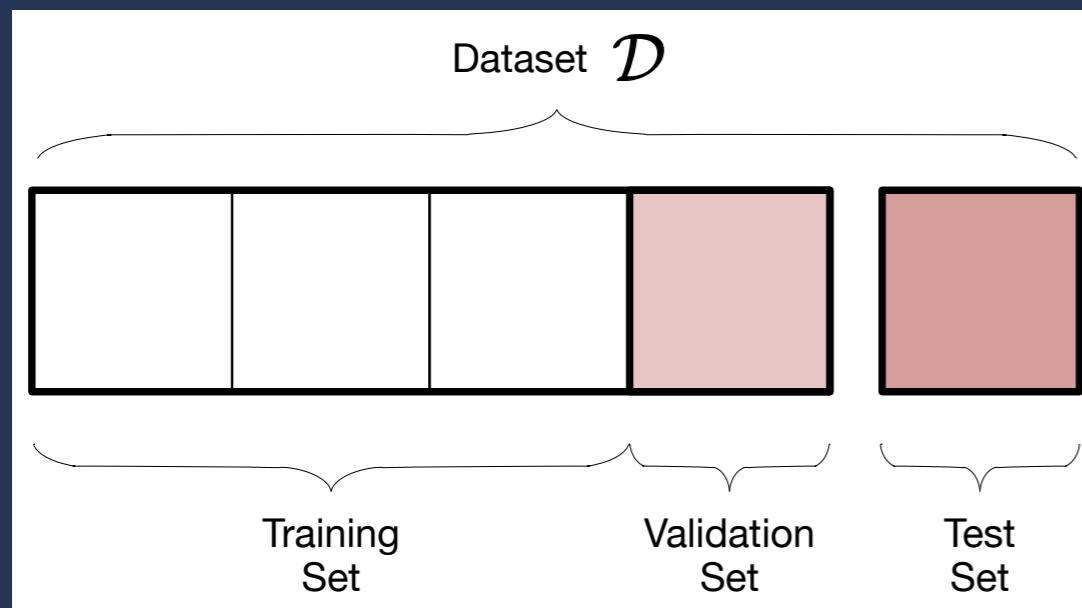


BALANCE THE COMPLEXITY

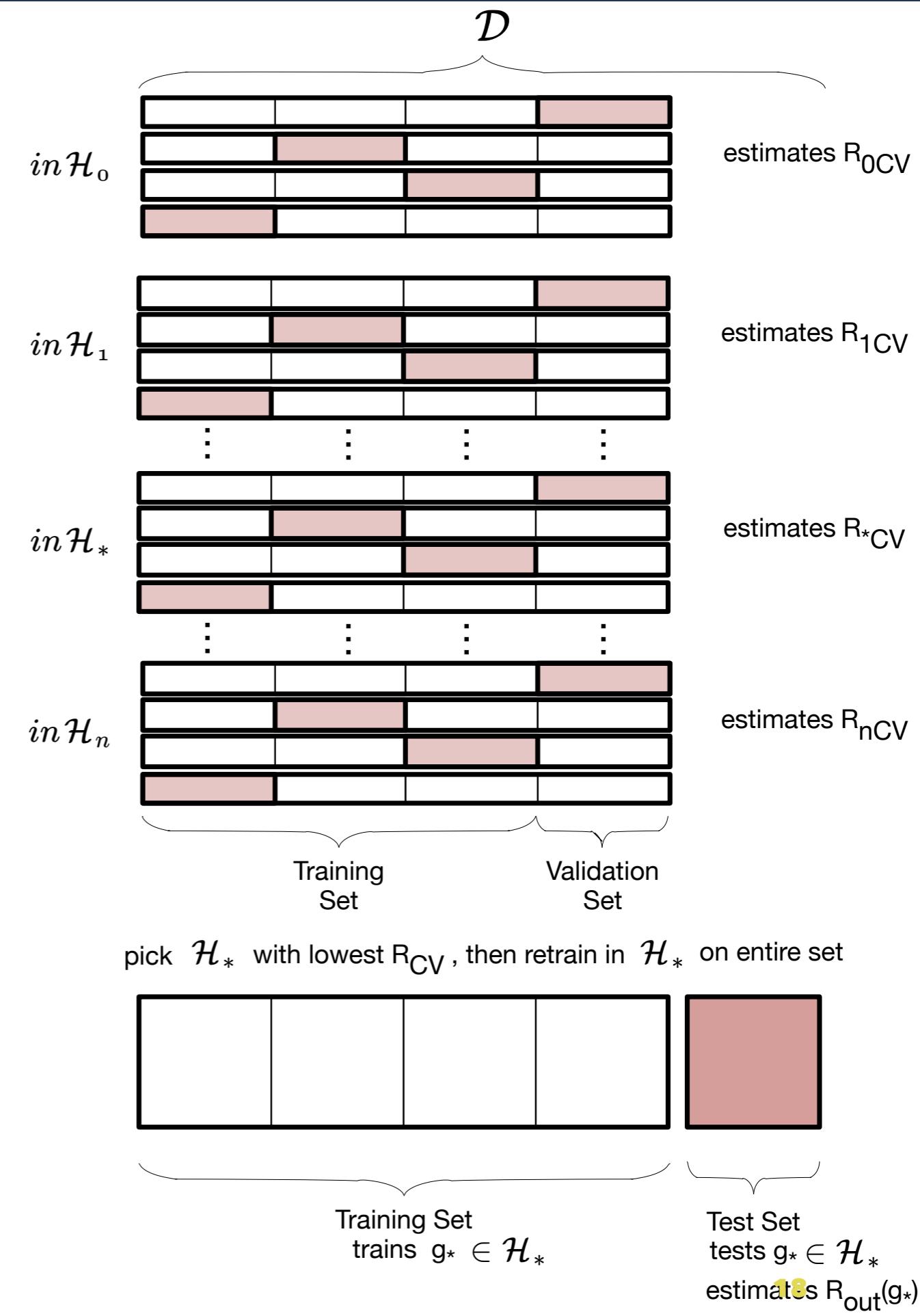
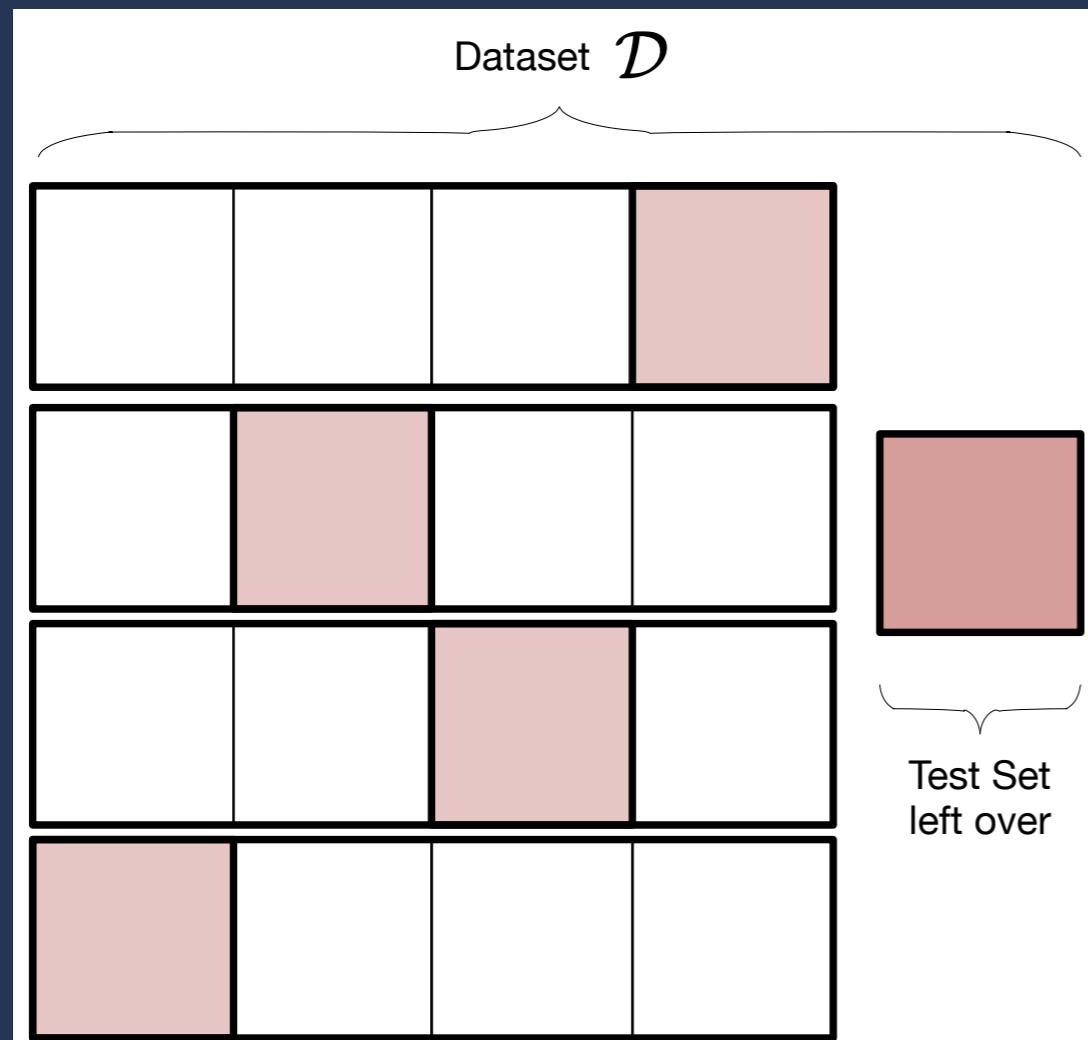


VALIDATION

- train-test not enough as we *fit* for d on test set and contaminate it
- thus do train-validate-test



CROSS-VALIDATION



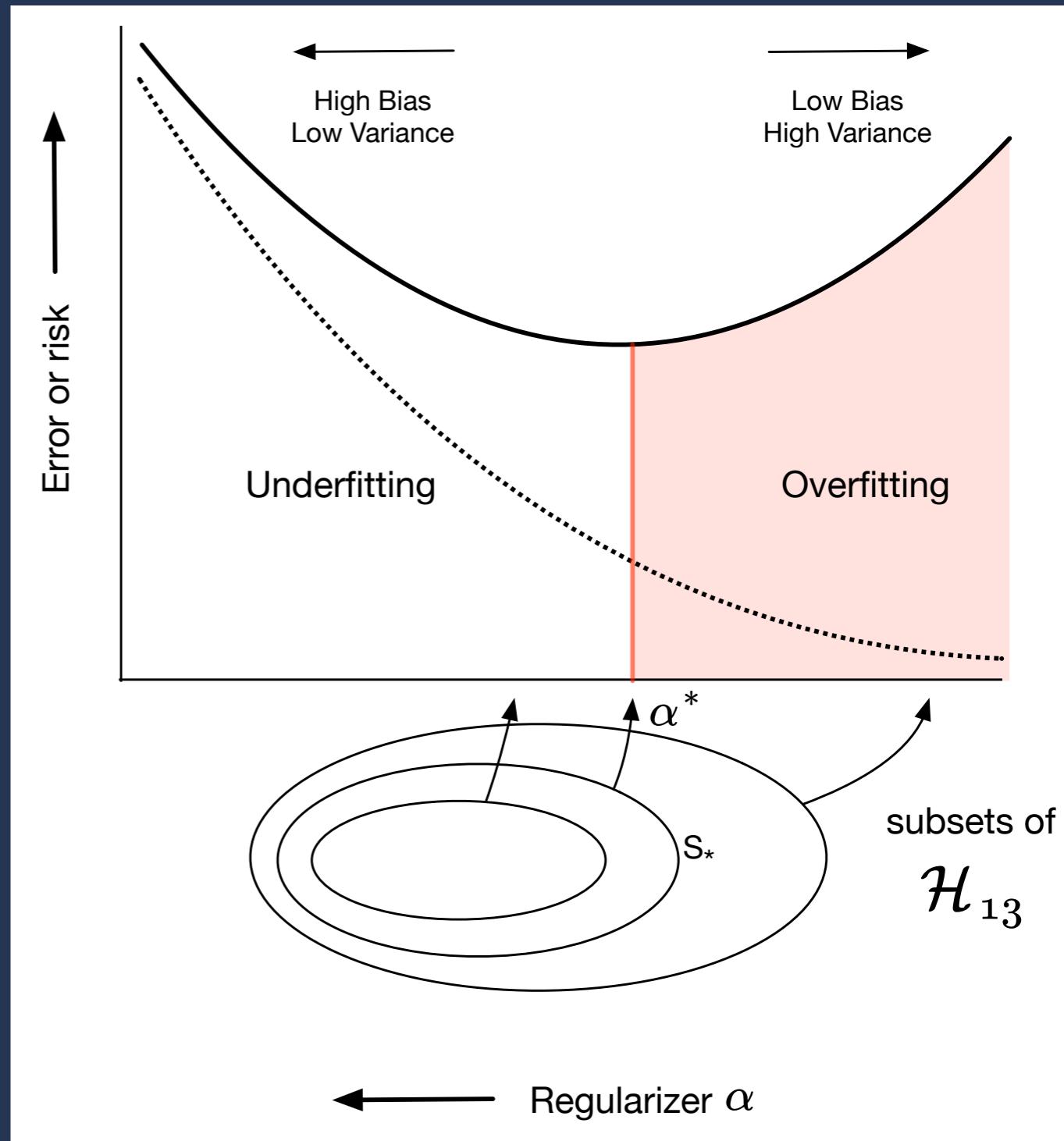
REGULARIZATION

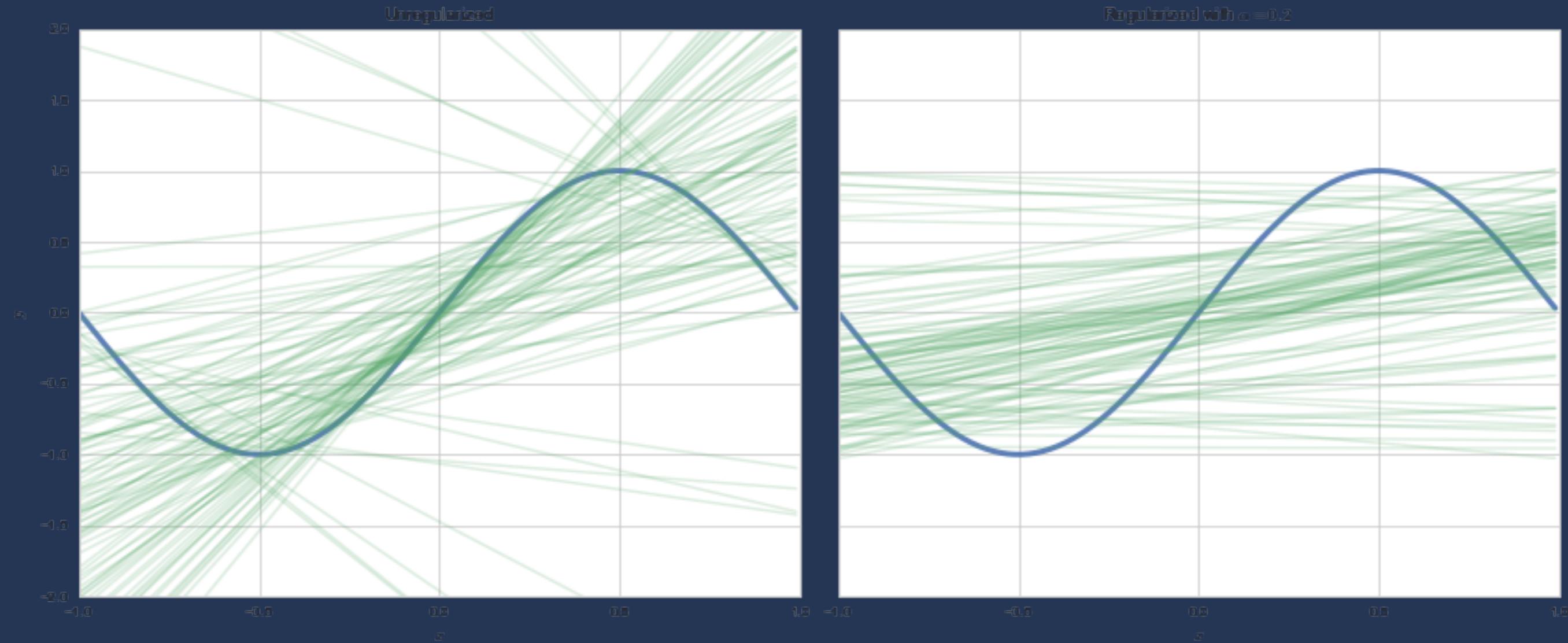
Keep higher a-priori complexity and impose a

complexity penalty

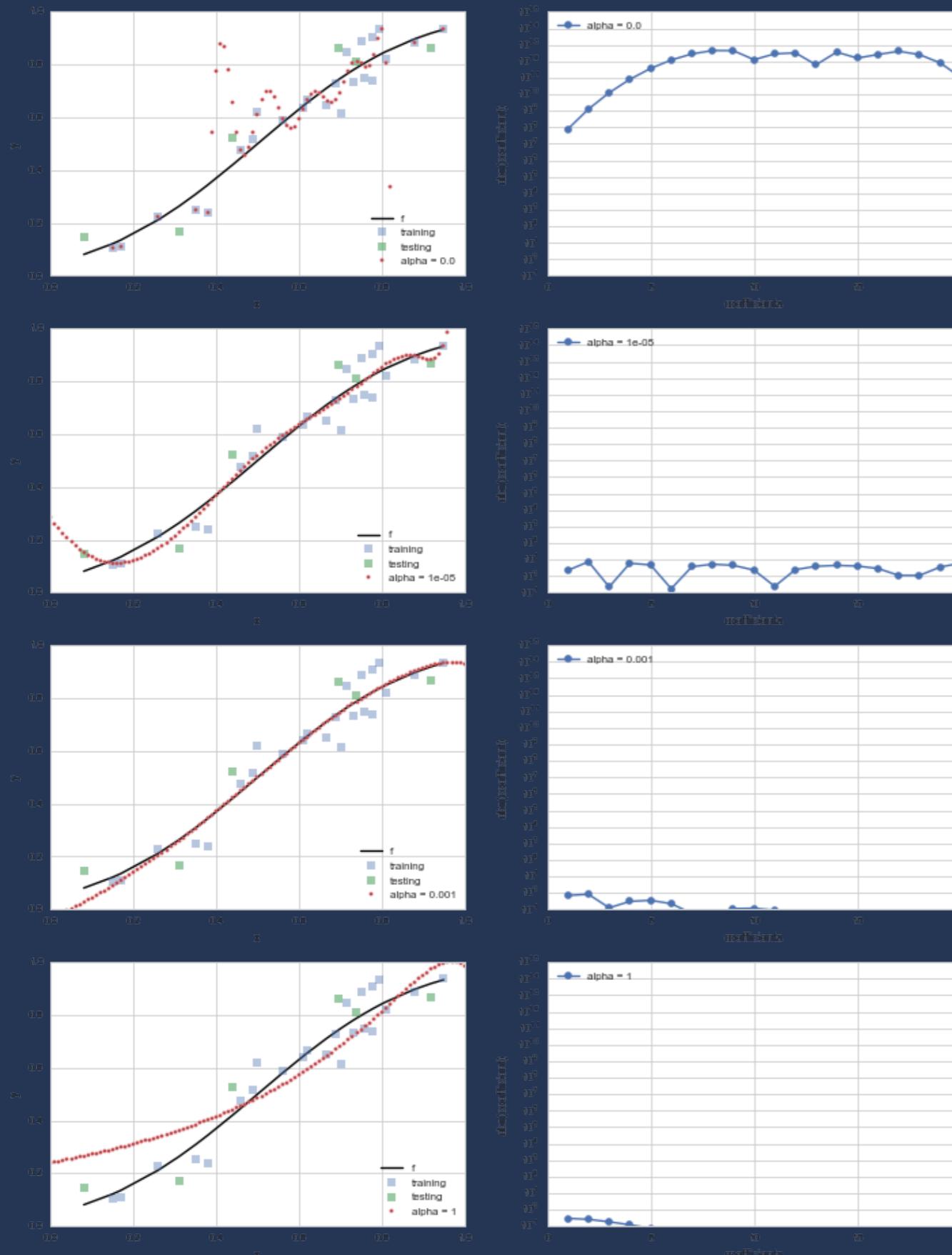
on risk instead, to choose a SUBSET of \mathcal{H}_{big} .
We'll make the coefficients small:

$$\sum_{i=0}^j a_i^2 < C.$$





REGULARIZATION



$$\mathcal{R}(h_j) = \sum_{y_i \in \mathcal{D}} (y_i - h_j(x_i))^2 + \alpha \sum_{i=0}^j a_i^2.$$

As we increase α , coefficients go towards 0.

Lasso uses $\alpha \sum_{i=0}^j |a_i|$, sets coefficients to exactly 0.

Thus regularization automates:

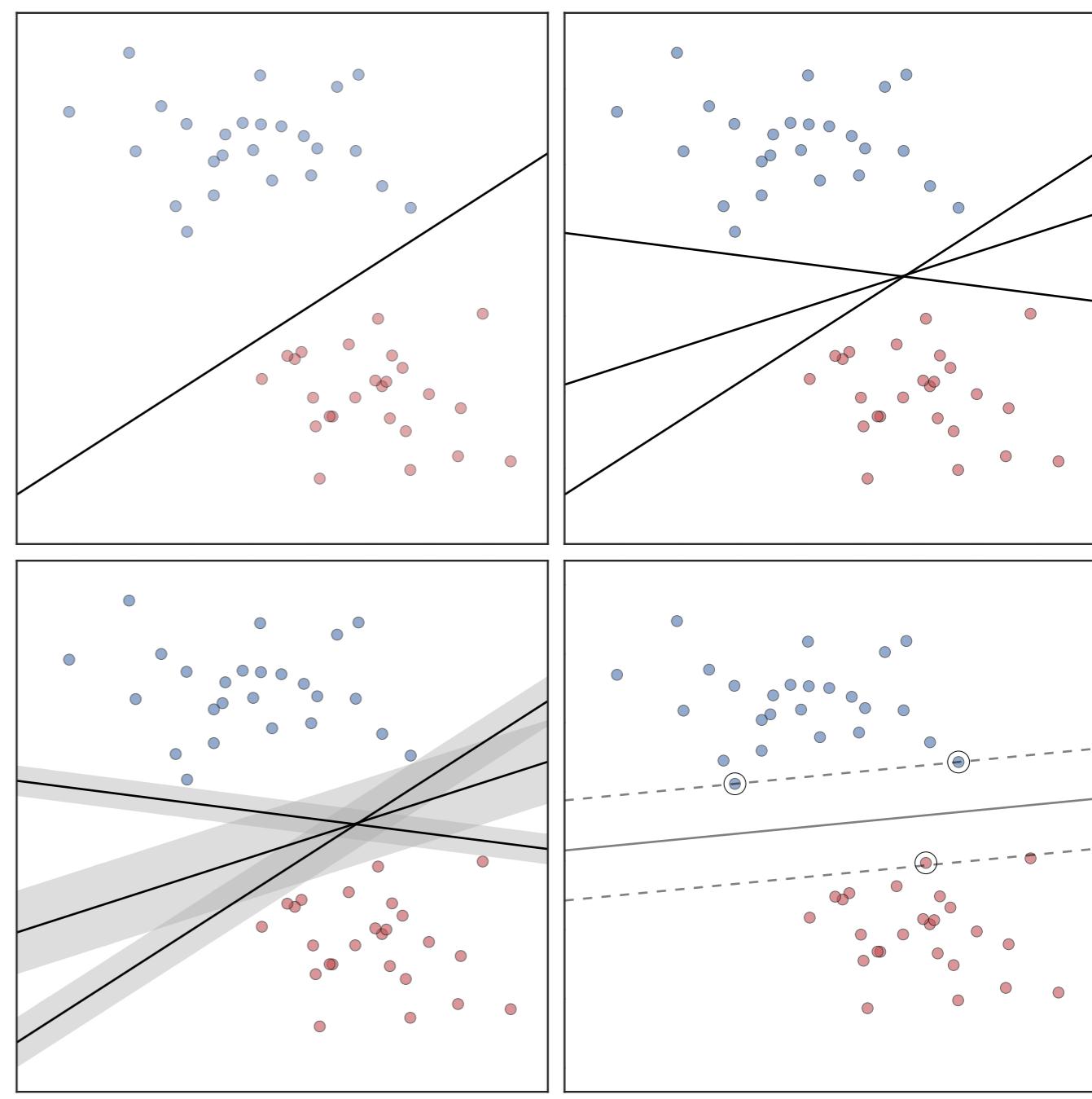
FEATURE ENGINEERING

CLASSIFICATION

BY LINEAR SEPARATION

Which line?

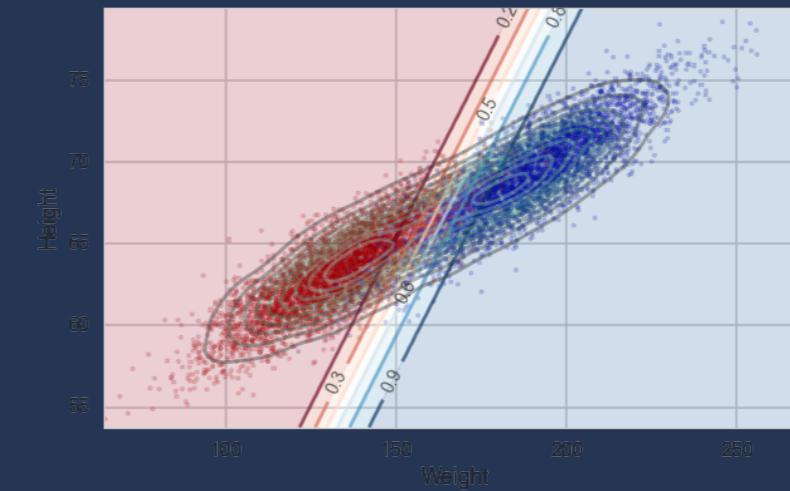
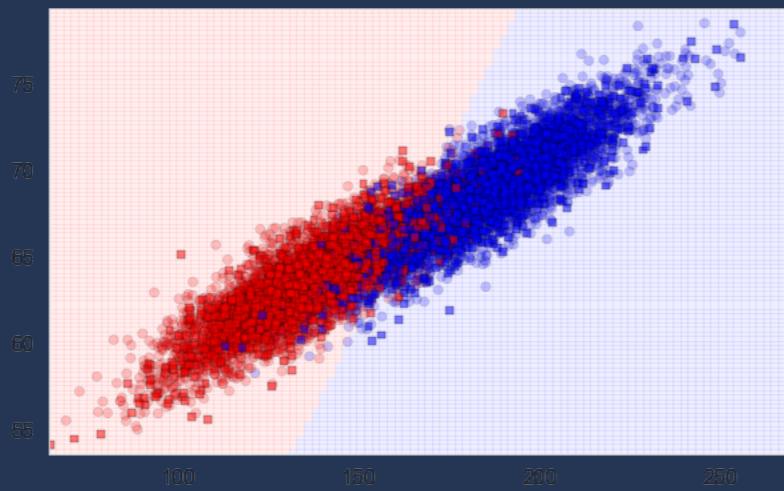
- Different Algorithms, different lines.
- SVM uses max-margin^j



^j image from code in <http://bit.ly/1Azg29G>

DISCRIMINATIVE CLASSIFIER

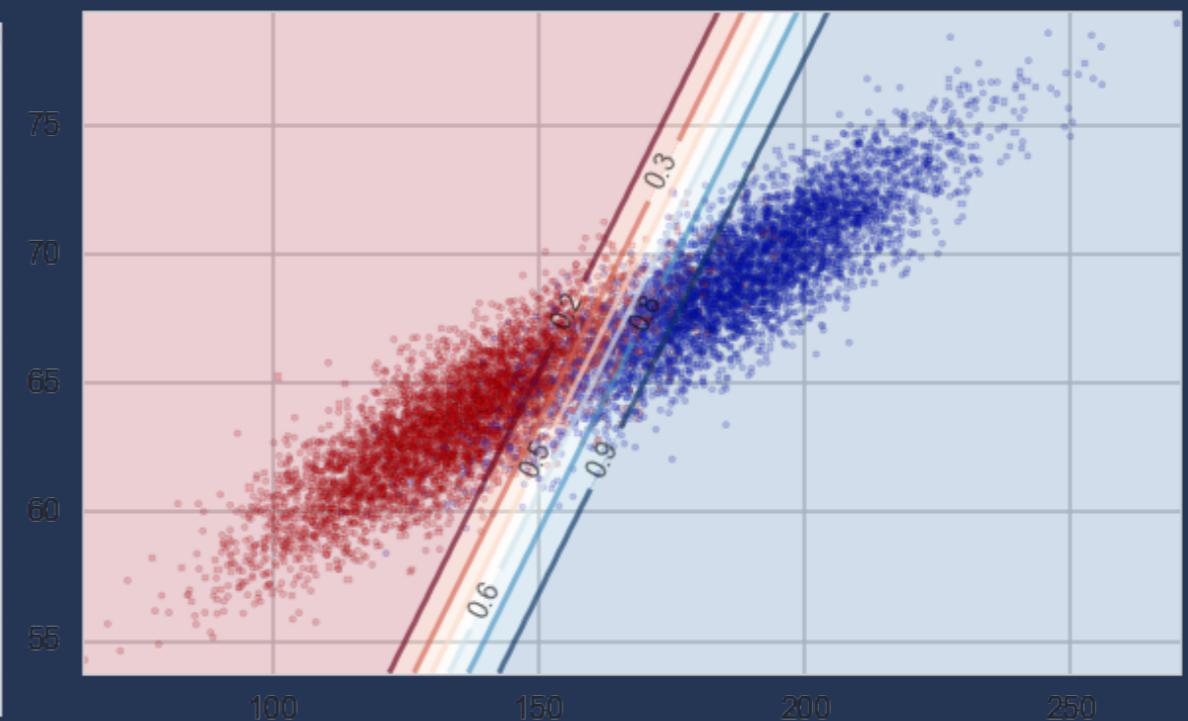
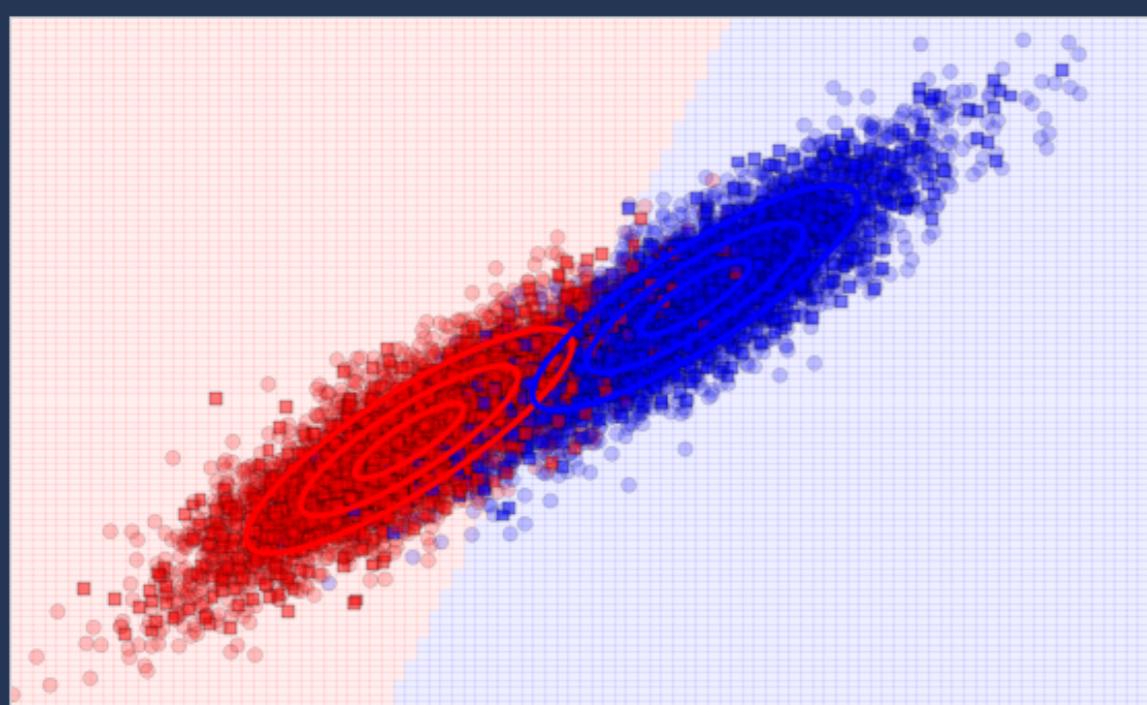
$$P(y|x) : P(\text{male}|\text{height}, \text{weight})$$



VS. DISCRIMINANT

GENERATIVE CLASSIFIER

$$P(y|x) \propto P(x|y)P(x) : P(\text{height}, \text{weight}|\text{male}) \times P(\text{male})$$



The virtual ATM:

: Inspired by <http://blog.yhathq.com/posts/image-classification-in-Python.html>

PCA

unsupervised dim reduction from 332x137x3 to 50

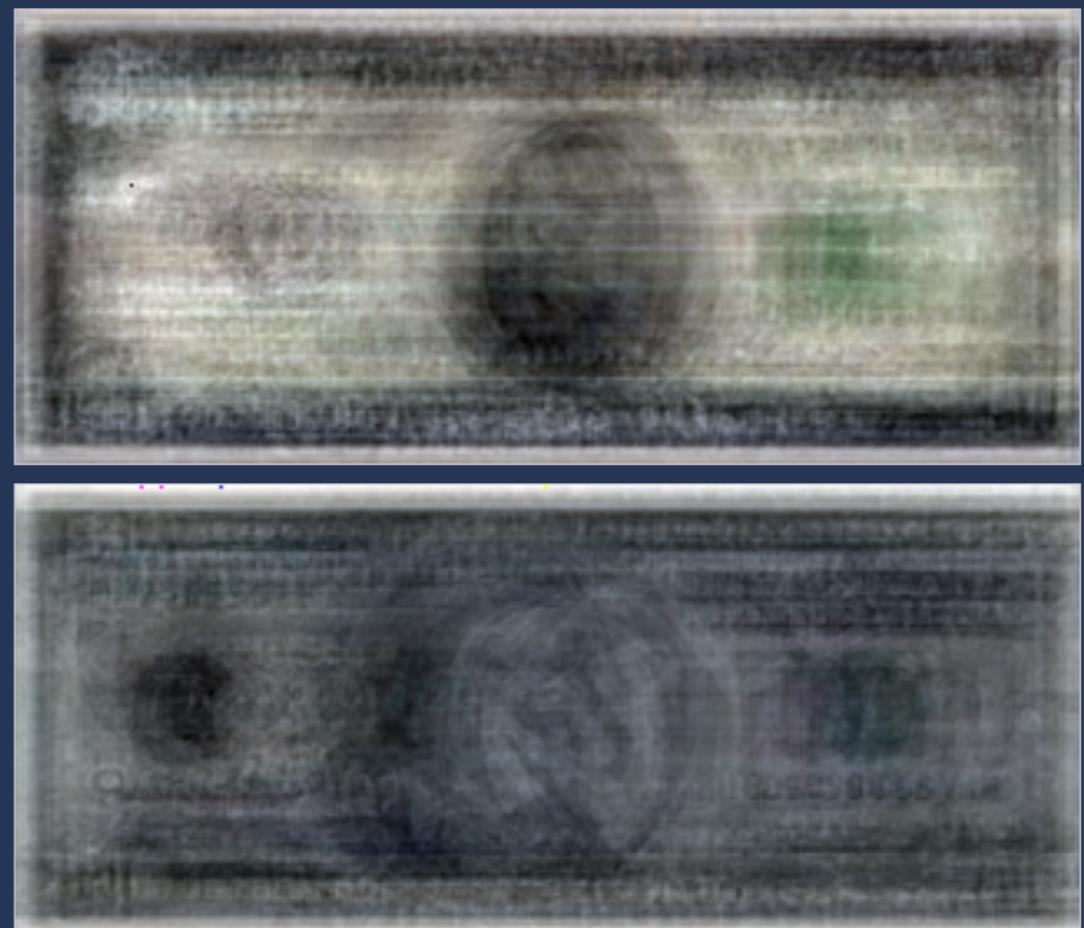
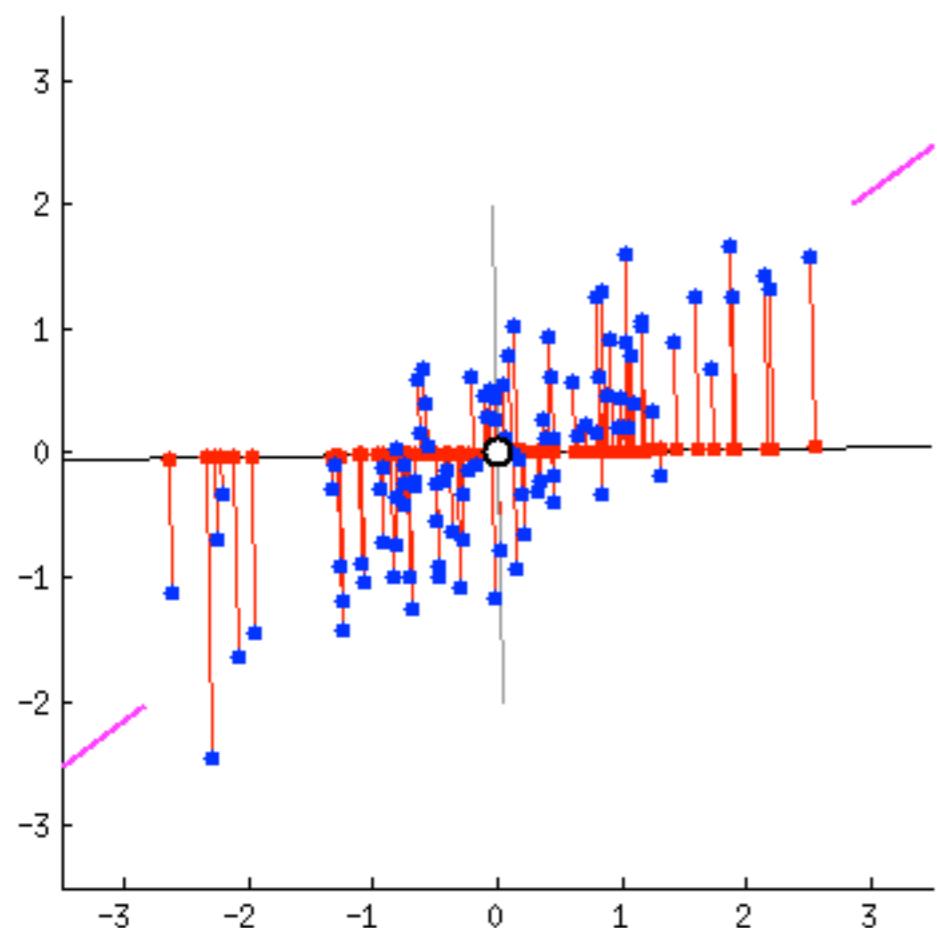
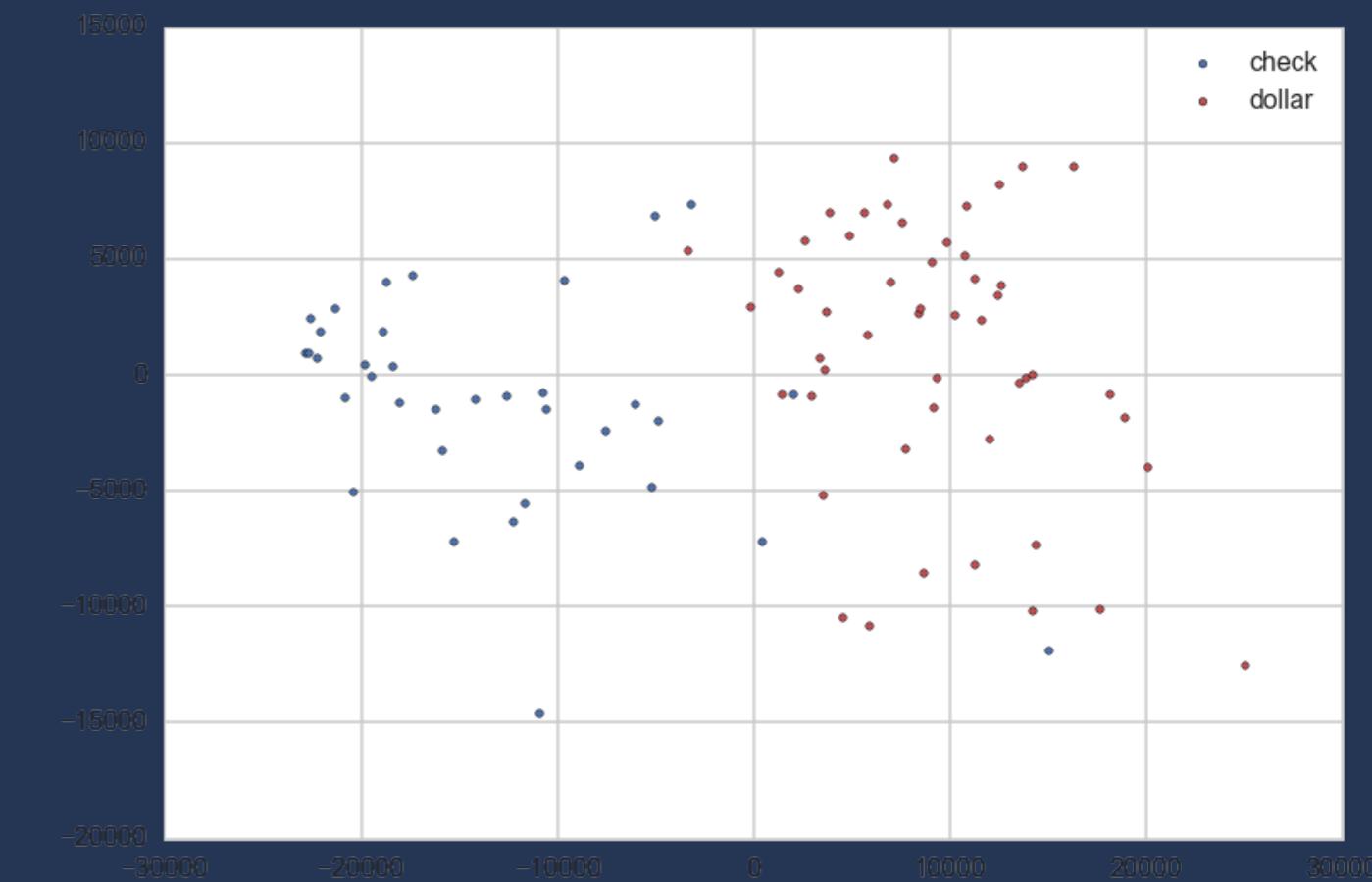
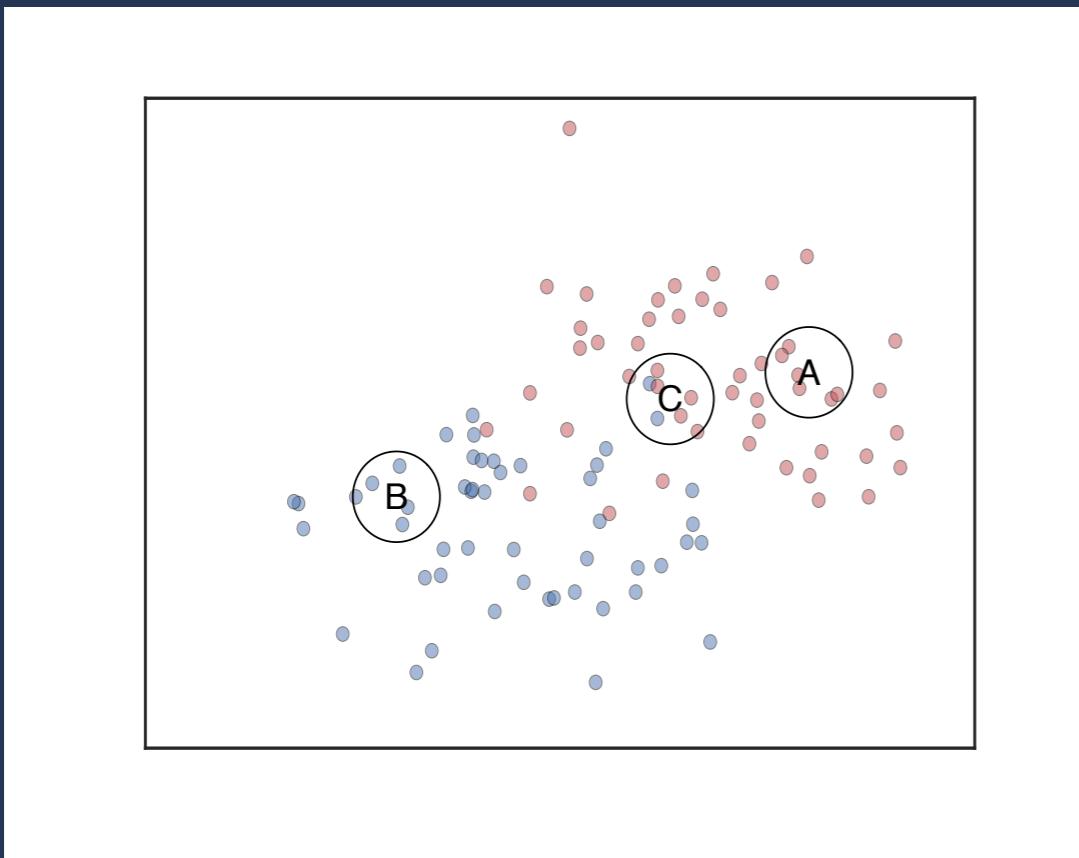
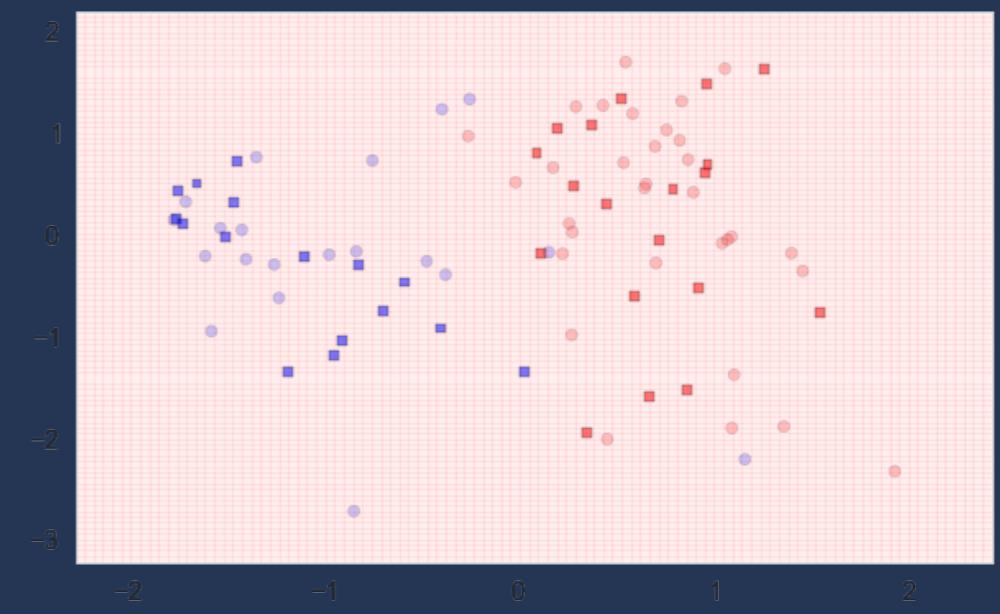
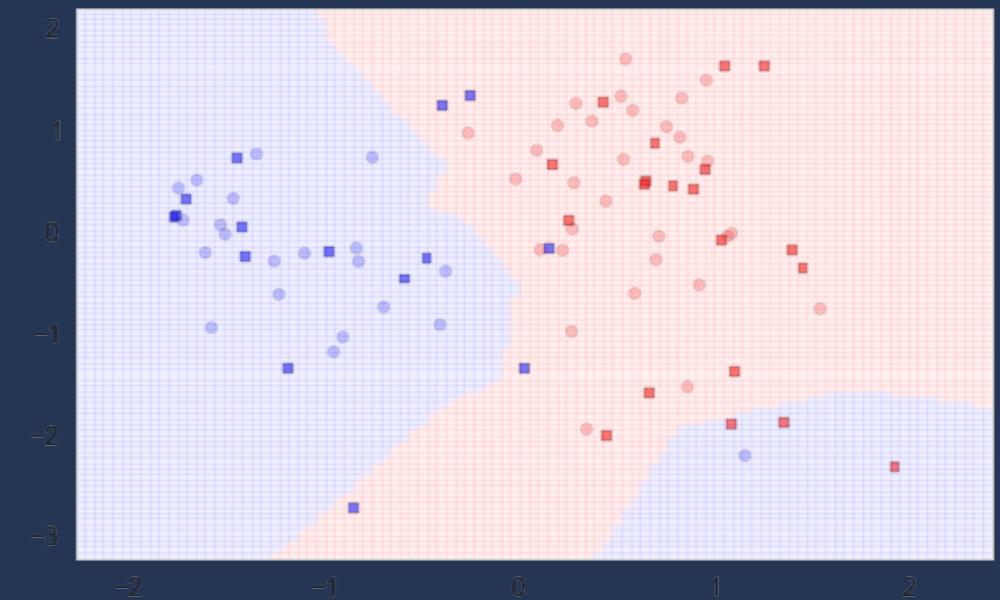
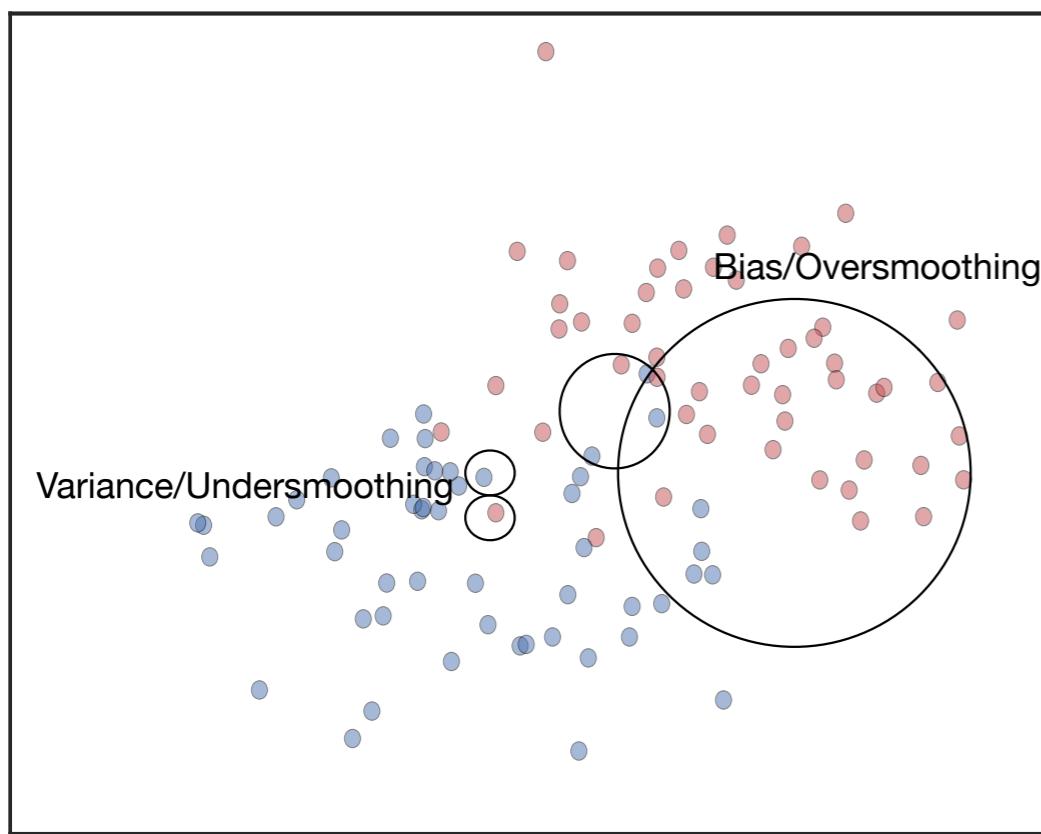


Diagram from <http://stats.stackexchange.com/a/140579>

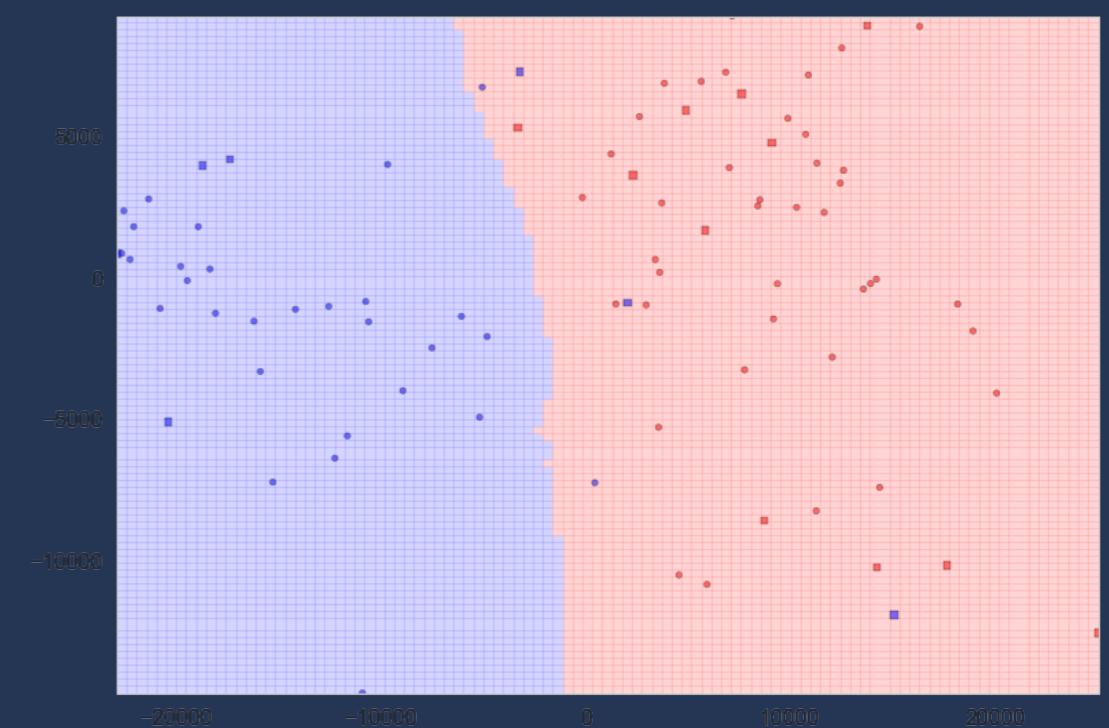
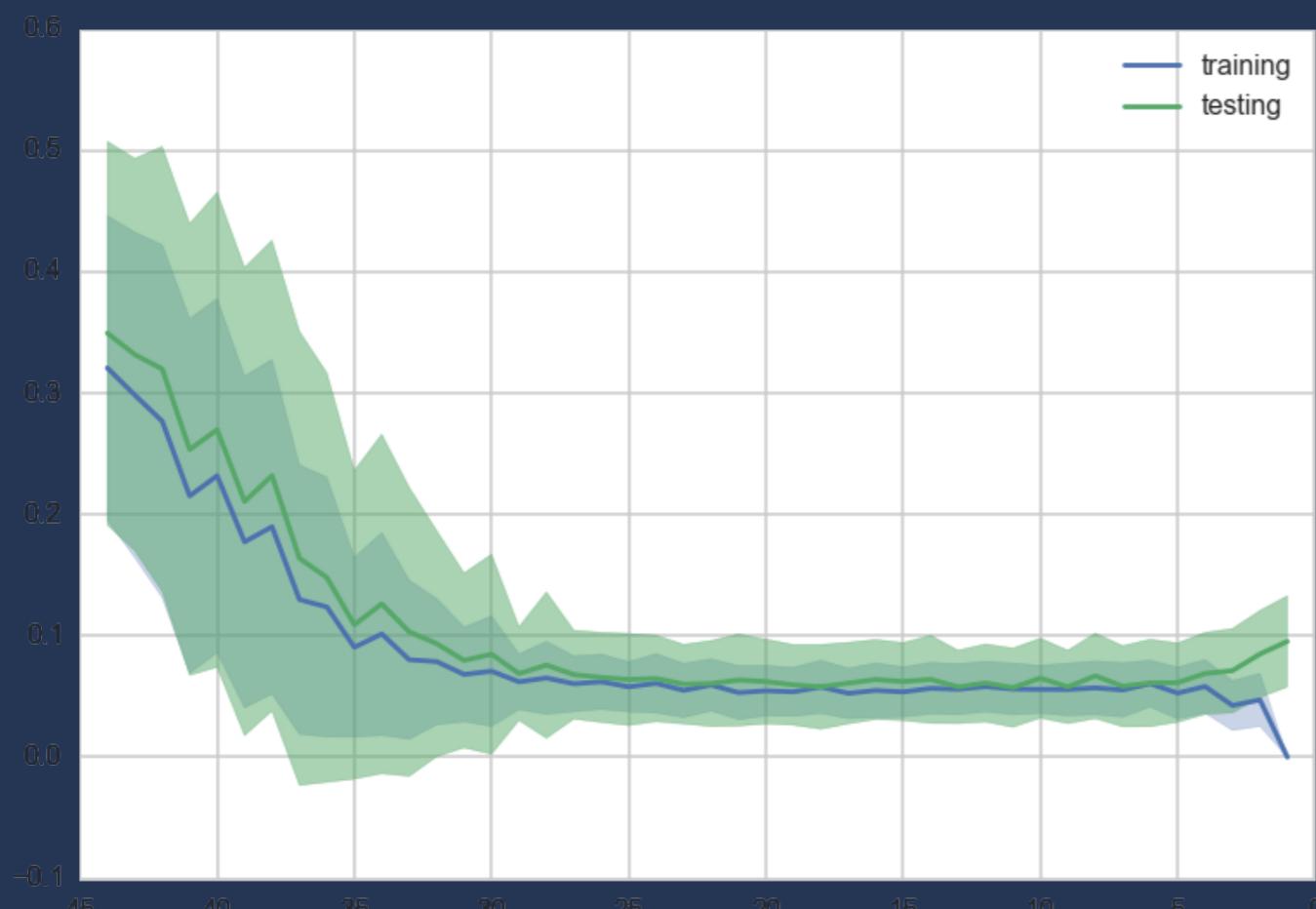
kNN



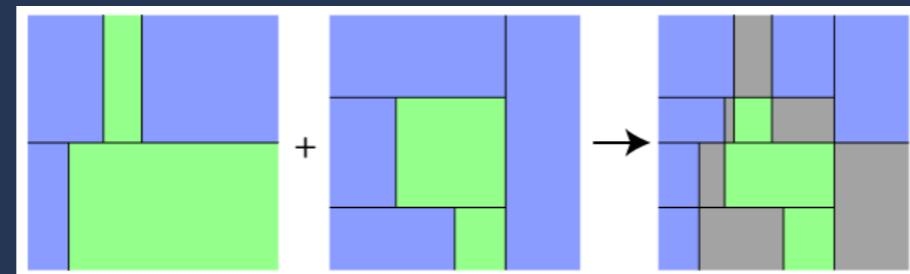
BIAS and VARIANCE in KNN



kNN CROSS- VALIDATED



ENSEMBLE LEARNING

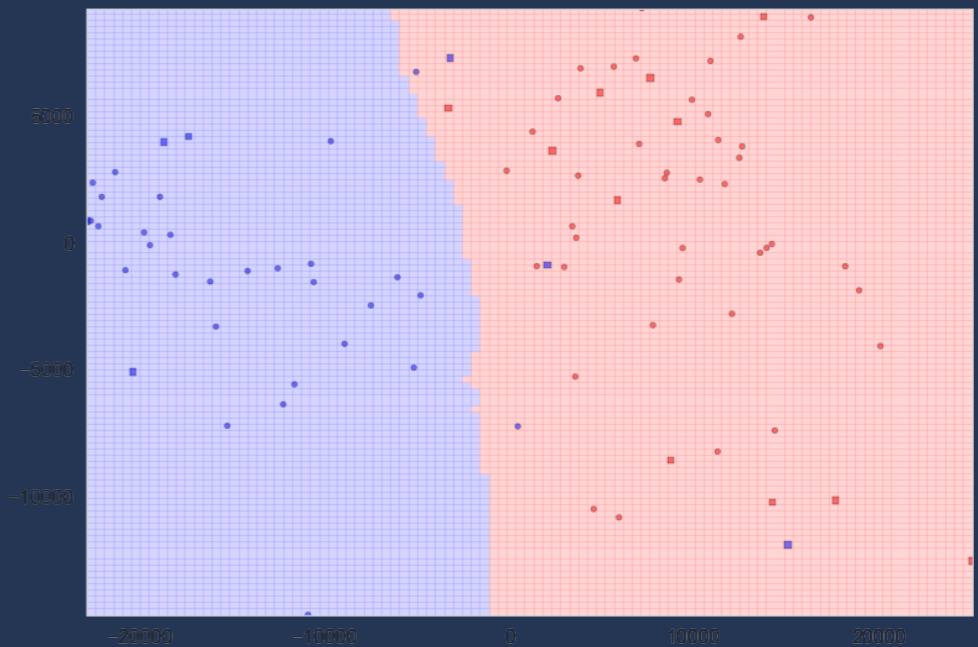


Combine multiple classifiers and vote.

Egs: Decision Trees → random forest, bagging, boosting

EVALUATING CLASSIFIERS

		Predicted	
		0	1
Observed	0	TN True Negative	FP False Positive
	1	FN False Negative	TP True Positive
		PN Predicted Negative	PP Predicted Positive

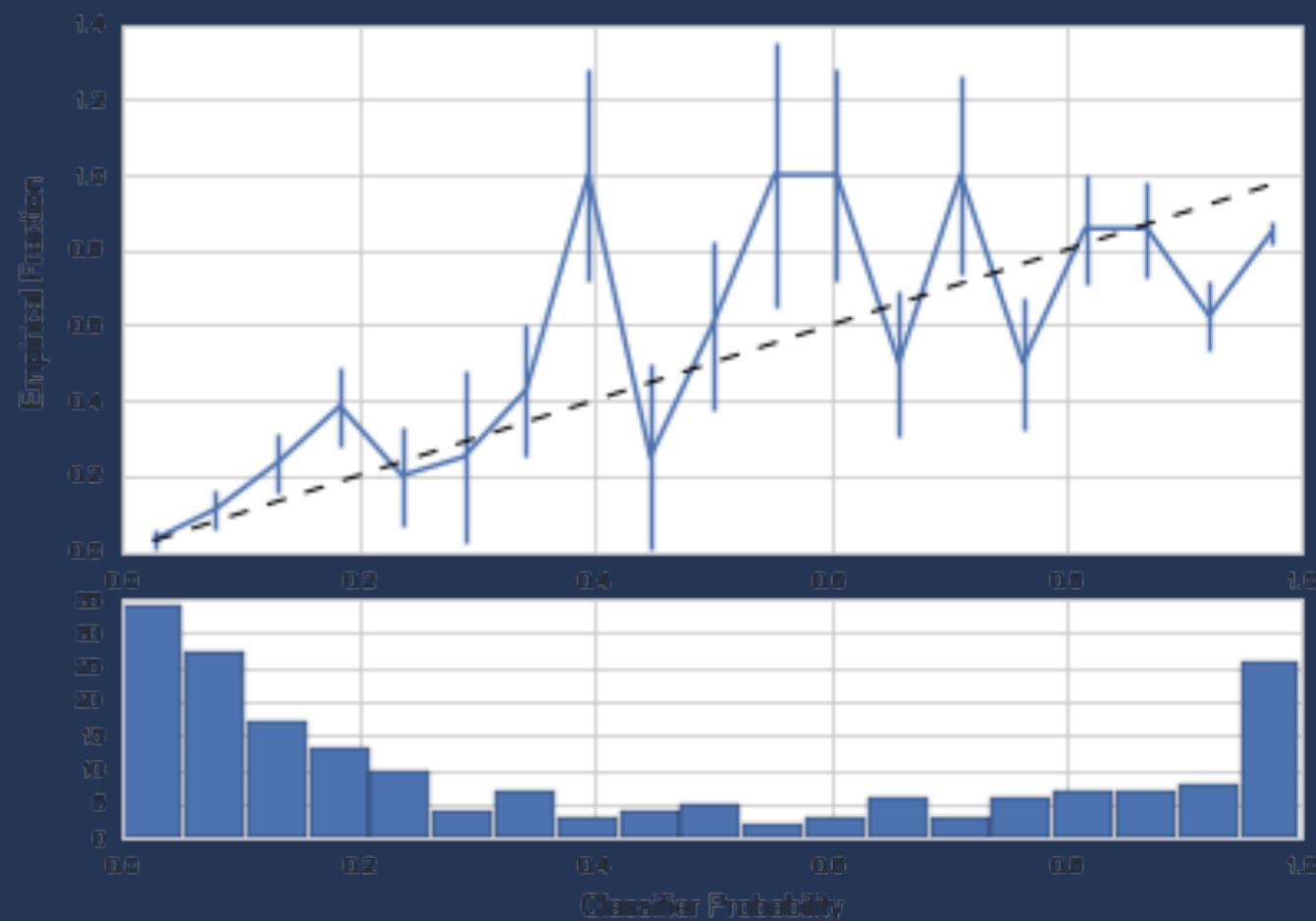
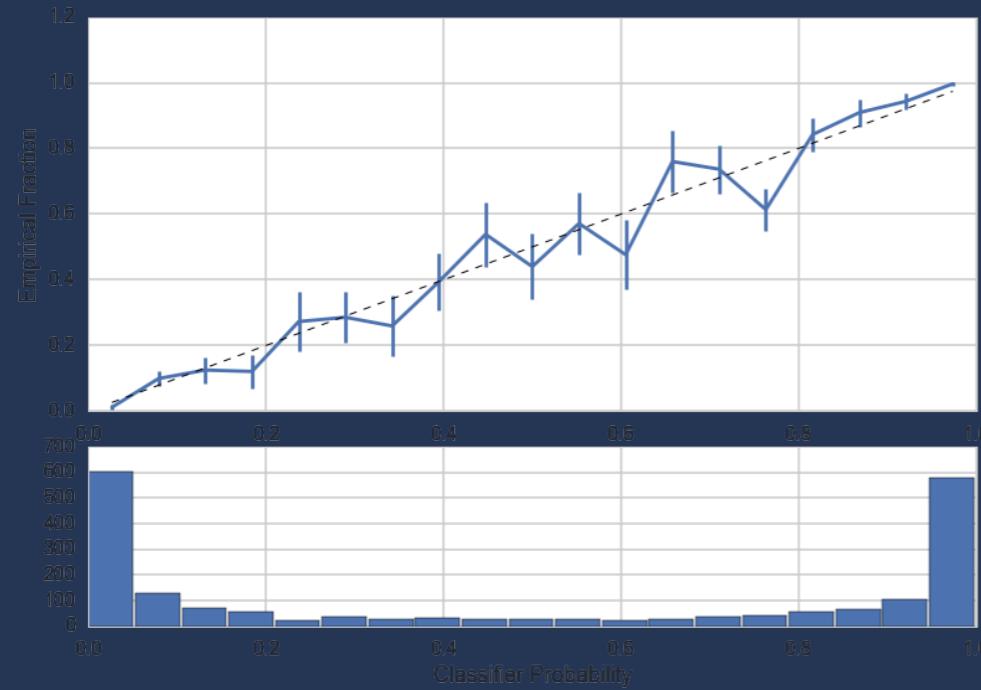


```
confusion_matrix(bestcv.predict(Xtest), ytest)  
[[11, 0],  
 [3, 4]]
```

checks=blue=1,
dollars=red=0

CLASSIFIER PROBABILITIES

- classifiers output rankings or probabilities
- ought to be well calibrated, or, atleast, similarly ordered



CLASSIFICATION RISK

- $R_{g,\mathcal{D}}(x) = P(y_1|x)\ell(g, y_1) + P(y_0|x)\ell(g, y_0)$
- The usual loss is the 1-0 loss $\ell = 1_{g \neq y}$.
- Thus, $R_{g=y_1}(x) = P(y_0|x)$ and $R_{g=y_0}(x) = P(y_1|x)$

CHOOSE CLASS WITH LOWEST RISK

$$1 \text{ if } R_1 \leq R_0 \implies 1 \text{ if } P(0|x) \leq P(1|x).$$

choose 1 if $P(1|x) \geq 0.5$! Intuitive!

ASYMMETRIC RISK[^]

want no checks misclassified as dollars, i.e., no false negatives: $\ell_{10} \neq \ell_{01}$.

Start with $R_g(x) = P(1|x)\ell(g, 1) + P(0|x)\ell(g, 0)$

$$\begin{aligned}R_1 &= l_{11}p_1 + l_{10}p_0 = l_{10}p_0 \\R_0 &= l_{01}p_1 + l_{00}p_0 = l_{01}p_1\end{aligned}$$

choose 1 if $R_1 < R_0$

[^] image of breast carcinoma from <http://bit.ly/1QgqhBw>

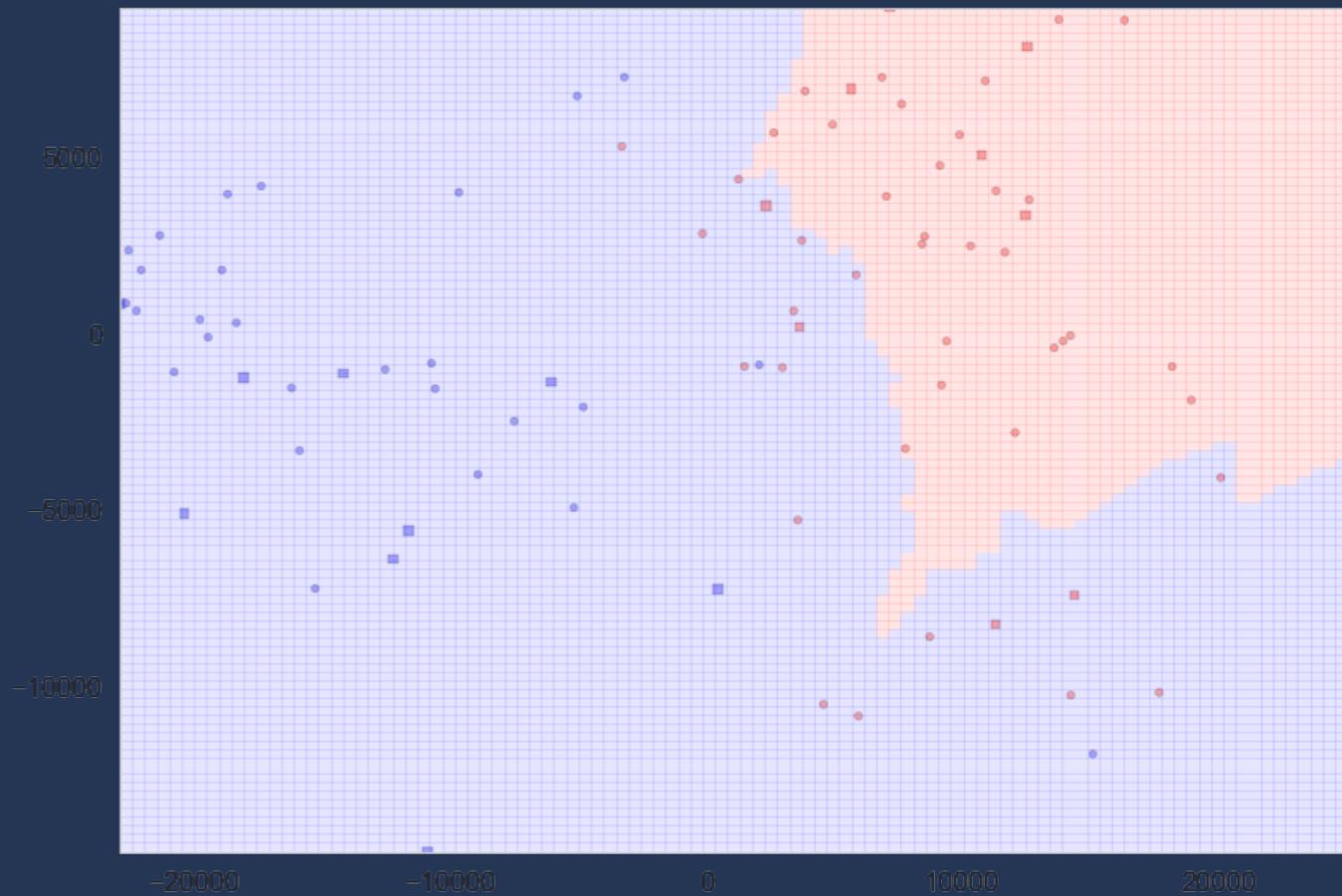
ASYMMETRIC RISK

i.e. $rp_0 < p_1$ where

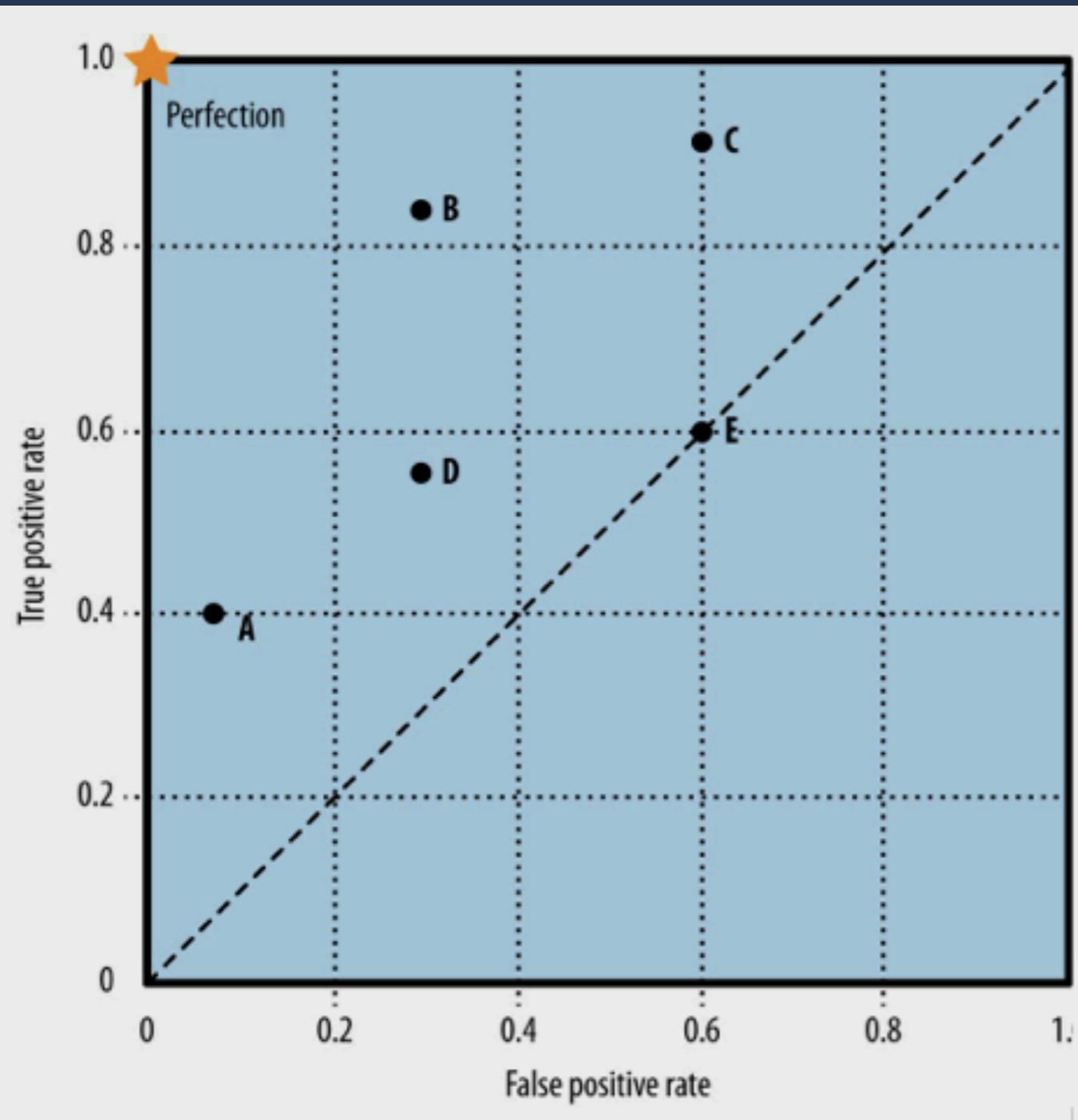
$$r = \frac{l_{10}}{l_{01}} = \frac{l_{FP}}{l_{FN}}$$

Or $P_1 > t$ where $t = \frac{r}{1+r}$

```
confusion_matrix(ytest, bestcv2, r=10, Xtest)
[[5, 4],
 [0, 9]]
```



ROC SPACE⁺



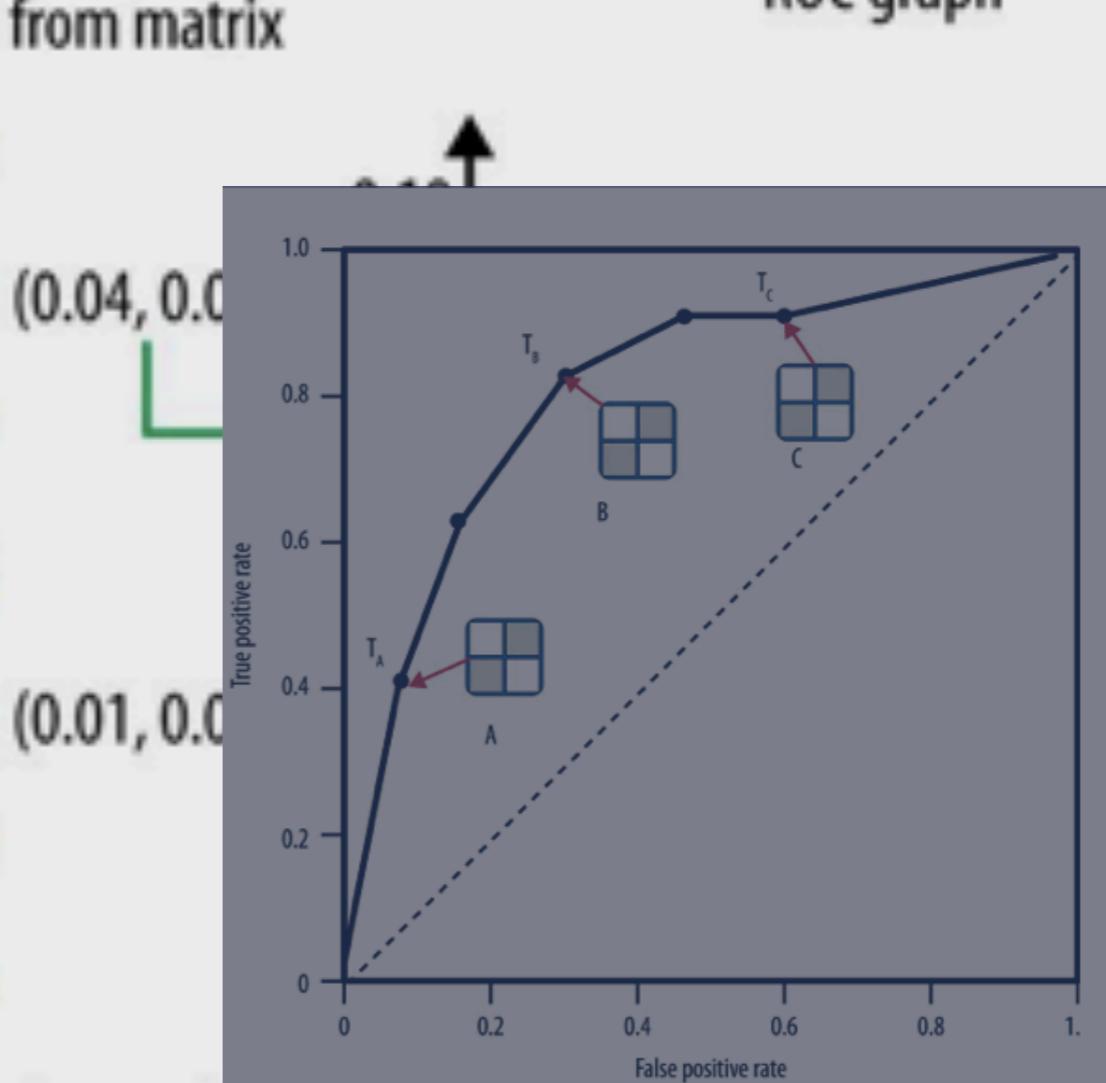
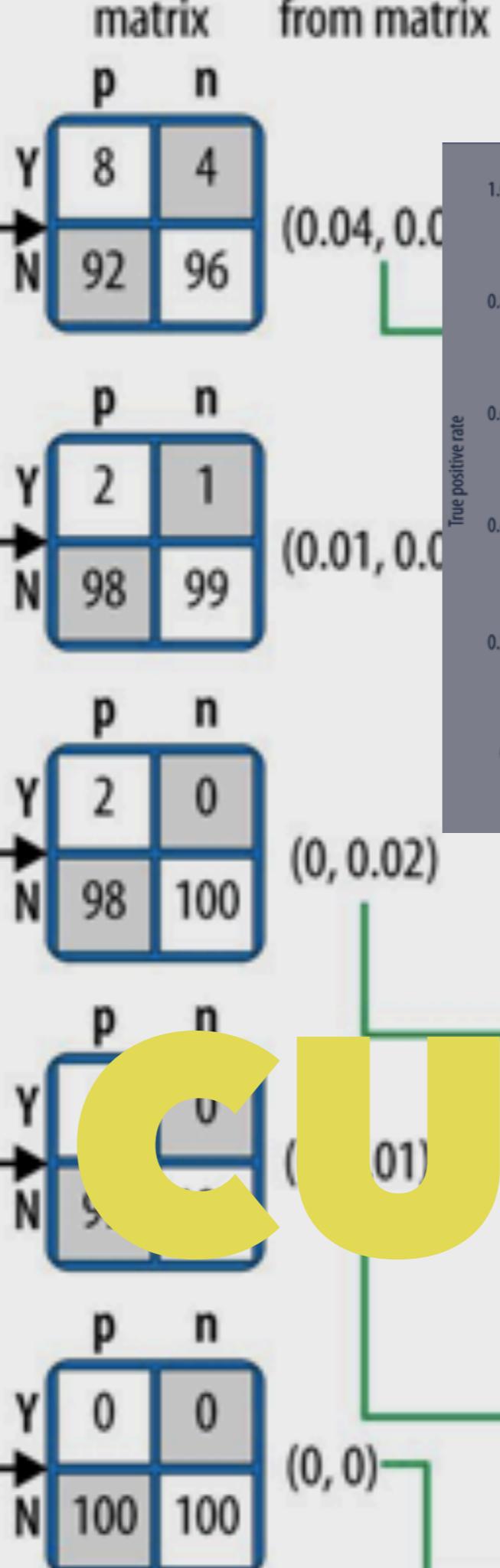
$$TPR = \frac{TP}{OP} = \frac{TP}{TP + FN}.$$

$$FPR = \frac{FP}{ON} = \frac{FP}{FP + TN}$$

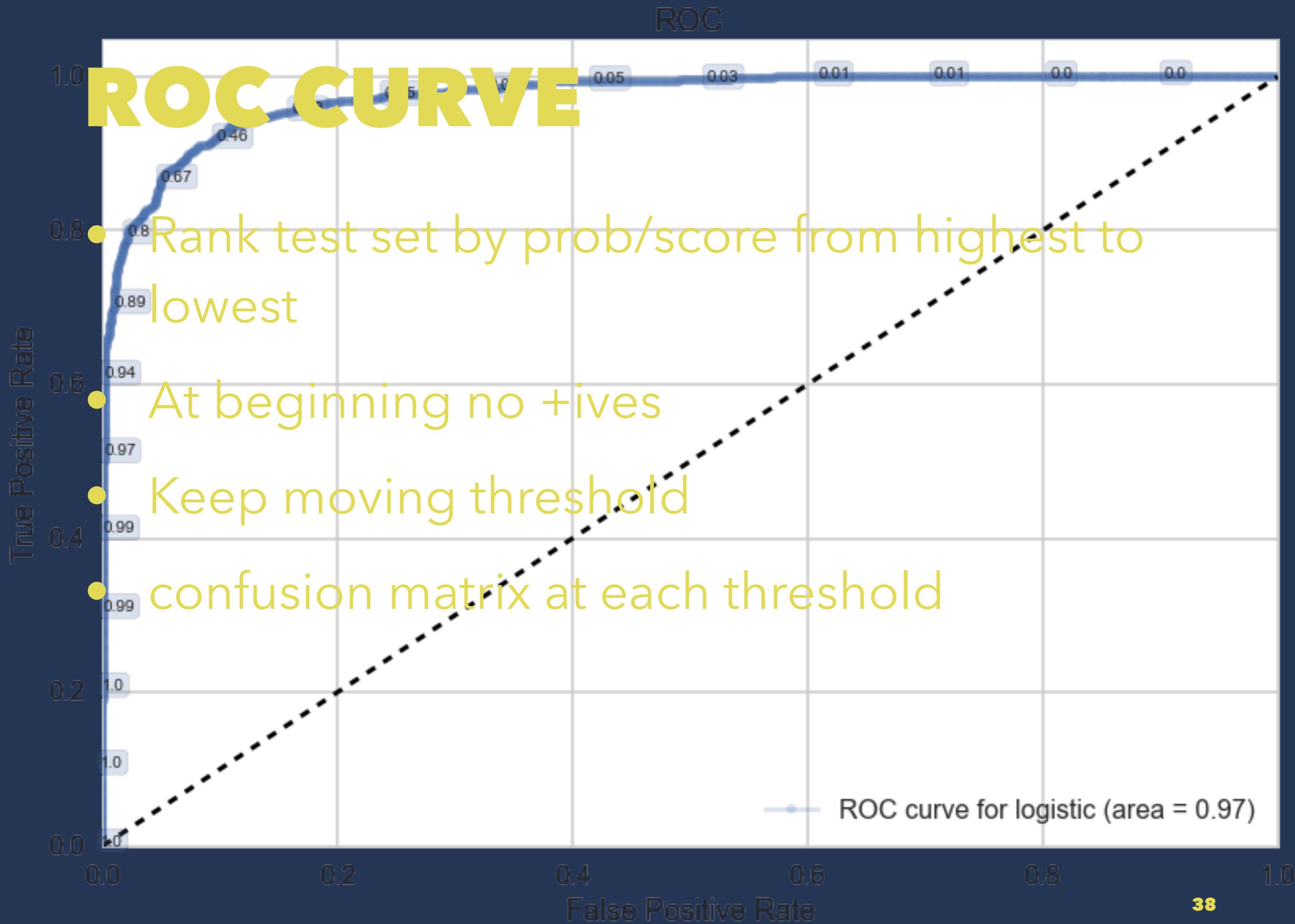
		Predicted		
		0	1	
Observed	0	TN True Negative	FP False Positive	ON Observed Negative
	1	FN False Negative	TP True Positive	OP Observed Positive
		PN Predicted Negative	PP Predicted Positive	

⁺ this+next fig: Data Science for Business, Foster et. al.

Instance	Class	Score
.	.	.
.	.	.
.	.	.
.....	p	0.65
.....	p	0.71
.....	n	0.74
.....	p	0.75
.....	n	0.80
.....	p	0.84
.....	p	0.85
.....	p	0.87
.....	p	0.88
.....	n	0.90
.....	n	0.96
.....	p	0.98
.....	p	0.99



ROC CURVE



COMPARING CLASSIFIERS

Telecom customer Churn data set from
@YhatHQ <

VMail Message	Day Mins	Day Calls	Day Charge	Eve Mins	Eve Calls	Eve Charge	Night Mins	Night Calls	Night Charge	Intl Mins	Intl Calls	Intl Charge	CustServ Calls	Churn?
25	265.1	110	45.07	197.4	99	16.78	244.7	91	11.01	10.0	3	2.70	1	False.
26	161.6	123	27.47	195.5	103	16.62	254.4	103	11.45	13.7	3	3.70	1	False.
0	243.4	114	41.38	121.2	110	10.30	162.6	104	7.32	12.2	5	3.29	0	False.
0	299.4	71	50.90	61.9	88	5.26	196.9	89	8.86	6.6	7	1.78	2	False.
0	166.7	113	28.34	148.3	122	12.61	186.9	121	8.41	10.1	3	2.73	3	False.

< <http://blog.yhathq.com/posts/predicting-customer-churn-with-sklearn.html>

ROC

True Positive Rate

1.0
0.8
0.6
0.4
0.2
0.0

ROC curves

0.0 0.2 0.4 0.6 0.8 1.0

False Positive Rate

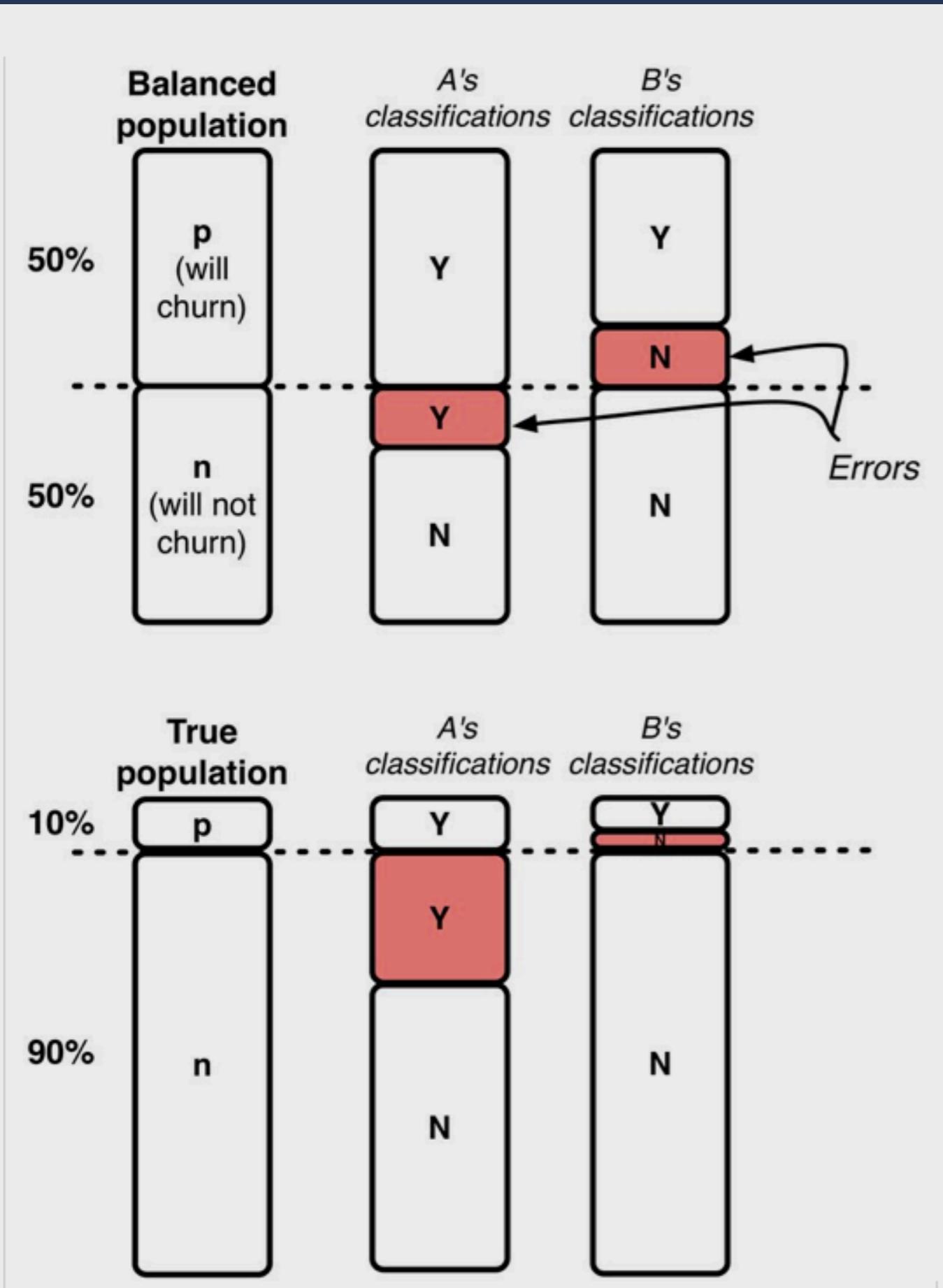
ROC curve for knn (area = 0.69)
ROC curve for rf (area = 0.91)

40

0.68
0.78
0.5
0.75
0.88
0.57
0.46
0.35
0.25
0.15
0.05
0.0

ASYMMETRIC CLASSES

- A has large FP[#]
- B has large FN.
- On asymmetric data sets, A will do very bad.
- But is it so?



[#] figure from Data Science for Business, Foster et. al.

EXPECTED VALUE FORMALISM

Can be used for risk or profit/utility (negative risk)

$$EP = p(1,1)\ell_{11} + p(0,1)\ell_{10} + p(0,0)\ell_{00} + p(1,0)\ell_{01}$$

$$\begin{aligned} EP &= p_a(1)[TPR\ell_{11} + (1 - TPR)\ell_{10}] \\ &\quad + p_a(0)[(1 - FPR)\ell_{00} + FPR\ell_{01}] \end{aligned}$$

Fraction of test set pred to be positive $x = PP/N$:

$$x = (TP + FP)/N = TPR p_o(1) + FPR p_o(0)$$

ASYMMETRIC CLASSES

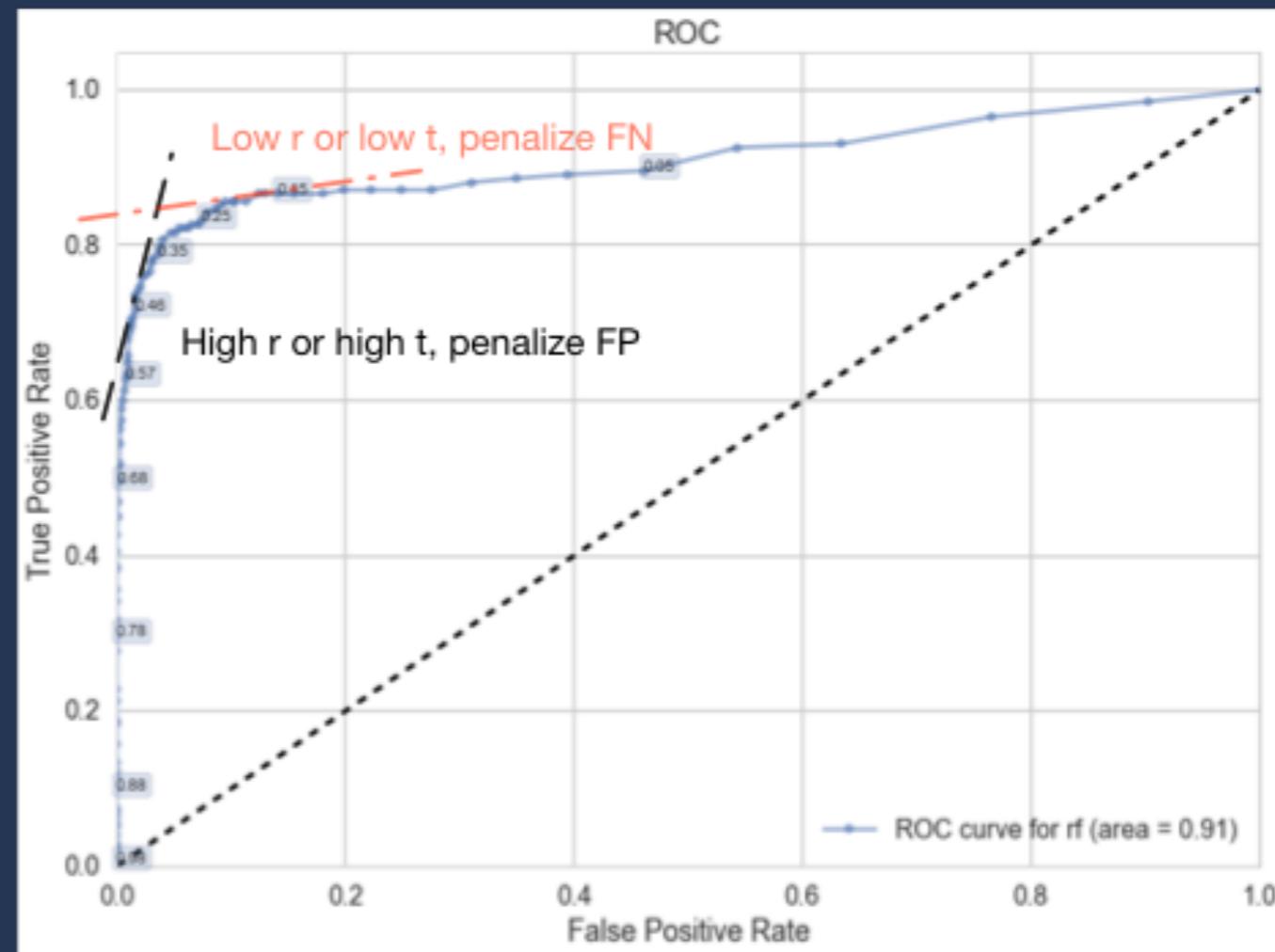
$$r = \frac{l_{FP}}{l_{FN}}$$

Equicost lines with slope

$$\frac{rp(0)}{p(1)} = \frac{rp(-)}{rp(+)}$$

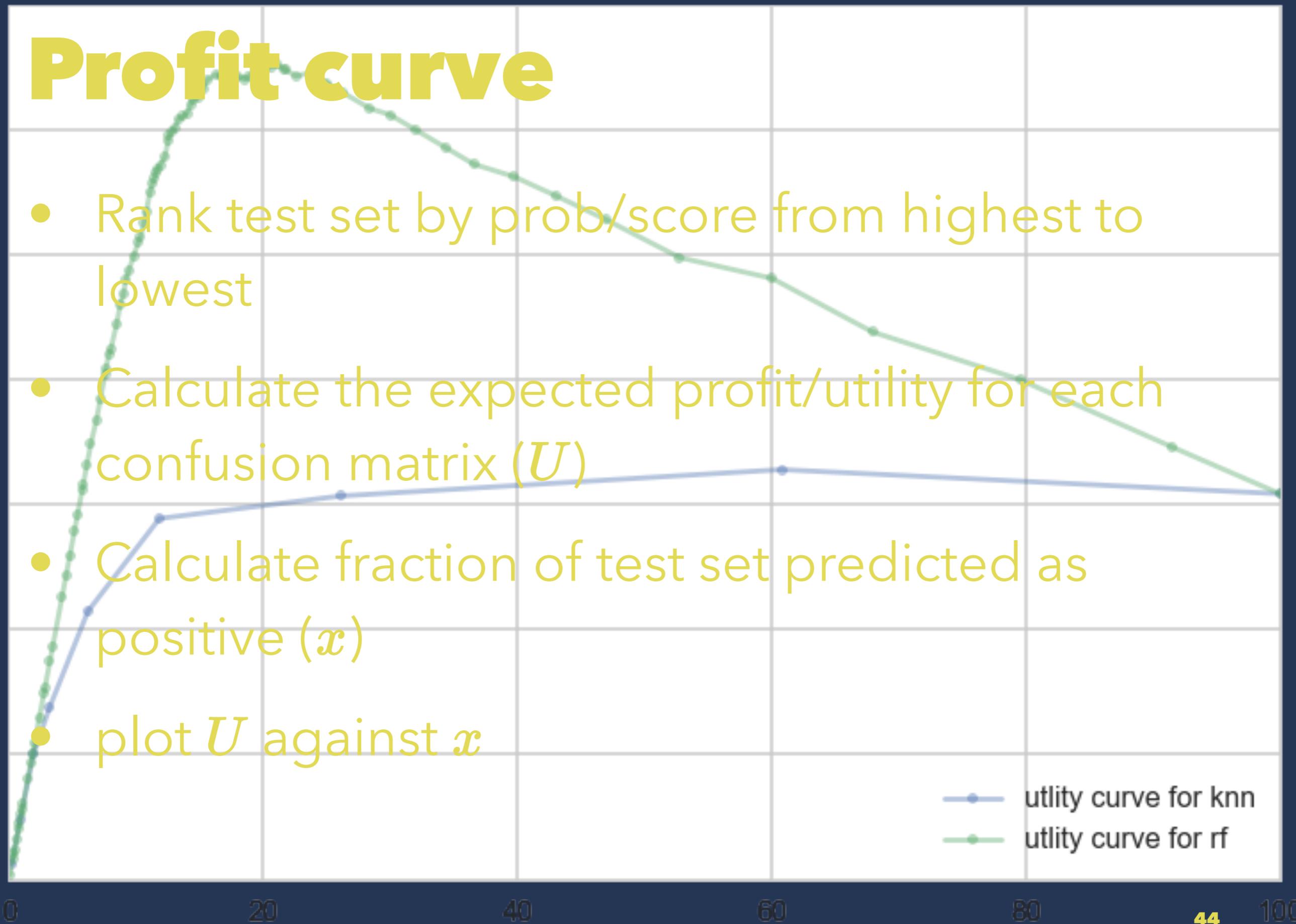
Small r penalizes FN.

Churn and Cancer u dont want FN: an uncaught chunner or cancer patient ($P=\text{churn}/\text{cancer}$)



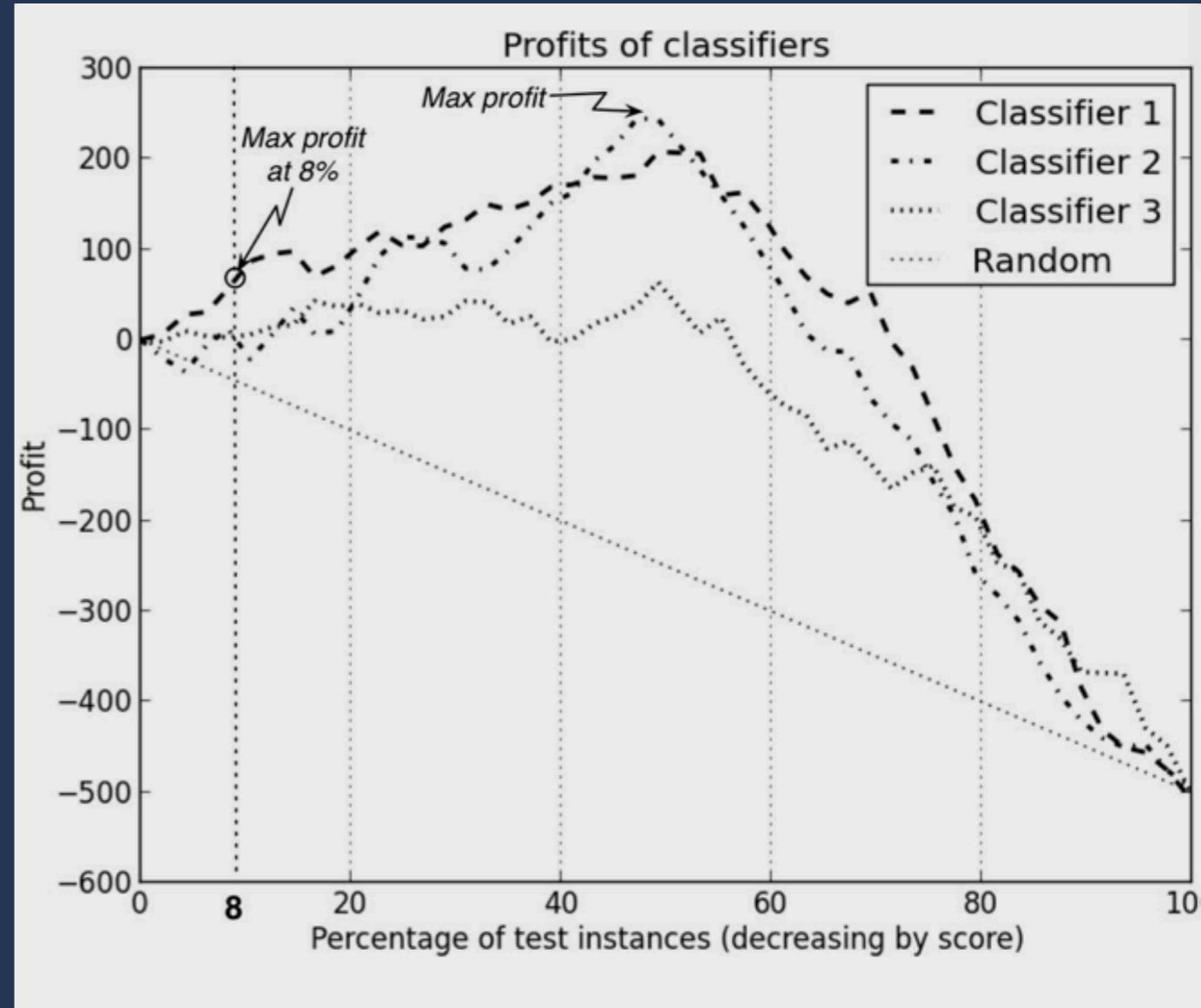
Profit curve

- Rank test set by prob/score from highest to lowest
- Calculate the expected profit/utility for each confusion matrix (U)
- Calculate fraction of test set predicted as positive (x)
- plot U against x



Finite budget[#]

- 100,000 customers, 40,000 budget, 5\$ per customer
- we can target 8000 customers
- thus target top 8%
- classifier 1 does better there, even though classifier 2 makes max profit



[#] figure from Data Science for Business, Foster et. al.

WHERE TO GO FROM HERE?

- Follow @YhatHQ, @kdnuggets etc
- Read Provost, Foster; Fawcett, Tom. Data Science for Business: What you need to know about data mining and data-analytic thinking. O'Reilly Media.
- check out Harvard's cs109: cs109.org
- Ask lots of questions of your data science team
- Follow your intuition

THANKS!!

@rahuldave