


A photograph of a snowy forest path. The path is covered in snow and leads into a dense forest of bare trees. A stone wall runs along the right side of the path. The text is overlaid on the image.

Measures of Diversity in Combining Classifiers

Part 2. Non-pairwise diversity measures

For fewer cartoons and more formulas:

<http://www.bangor.ac.uk/~mas00a/publications.html>

Random forest :  , x, θ_k (i.i.d, $k=1, \dots, L$), L is large

Strength and correlation:

$D(x)$: the class label of x suggested by D

Define margin function for a random forest to be

$$\text{mr}(x, \omega_j) = P_{\theta}(D(x)=\omega_j) - \max_{t \neq j} P_{\theta}(D(x)=\omega_t) ,$$

and the strength of the set of classifiers to be

$$s = \mathbf{E}_{x, \omega} [\text{mr}(x, \omega)]$$

Denote $\omega_s = \text{argmax}_{t \neq j} P_{\theta}(D(x)=\omega_t)$ and define raw margin function to be

$$\text{rmr}(x, \omega_j, \theta) = I(D(x)=\omega_j) - I(D(x)=\omega_s) ,$$

where $I(\cdot)$ is an indicator function.

The probability of error of the ensemble is bounded as follows

$$PE^* \leq \rho (1 - s^2) / s^2$$

(mean) correlation
between $\text{rnr}(D_i)$, $\text{rnr}(D_k)$
(averaged across all pairs
of classifiers)

strength

The diagram illustrates the components of the ensemble error bound formula. The formula is $PE^* \leq \rho (1 - s^2) / s^2$. The Greek letter ρ is enclosed in an orange box, and an arrow points from the text "(mean) correlation between $\text{rnr}(D_i)$, $\text{rnr}(D_k)$ (averaged across all pairs of classifiers)" to this box. The term s^2 appears twice in the denominator and once in the numerator. An arrow points from the word "strength" to both s^2 terms.

“Although the bound is likely to be loose, it fulfils the same suggestive function for random forests as VC-type bounds do for other types of classifiers.”

The 2-class case:

$$mr(x, \omega_i) = 2 P_{\theta}(D(x)=\omega_i) - 1, \quad i = 1,2$$

the strength of the set of classifiers is

$$s = \mathbf{E}_{x,\omega} [mr(x, \omega)]$$

$$\approx 2/N [\Sigma_1 P_{\theta}(D(x)=\omega_1) + \Sigma_2 P_{\theta}(D(x)=\omega_2)] - 1$$

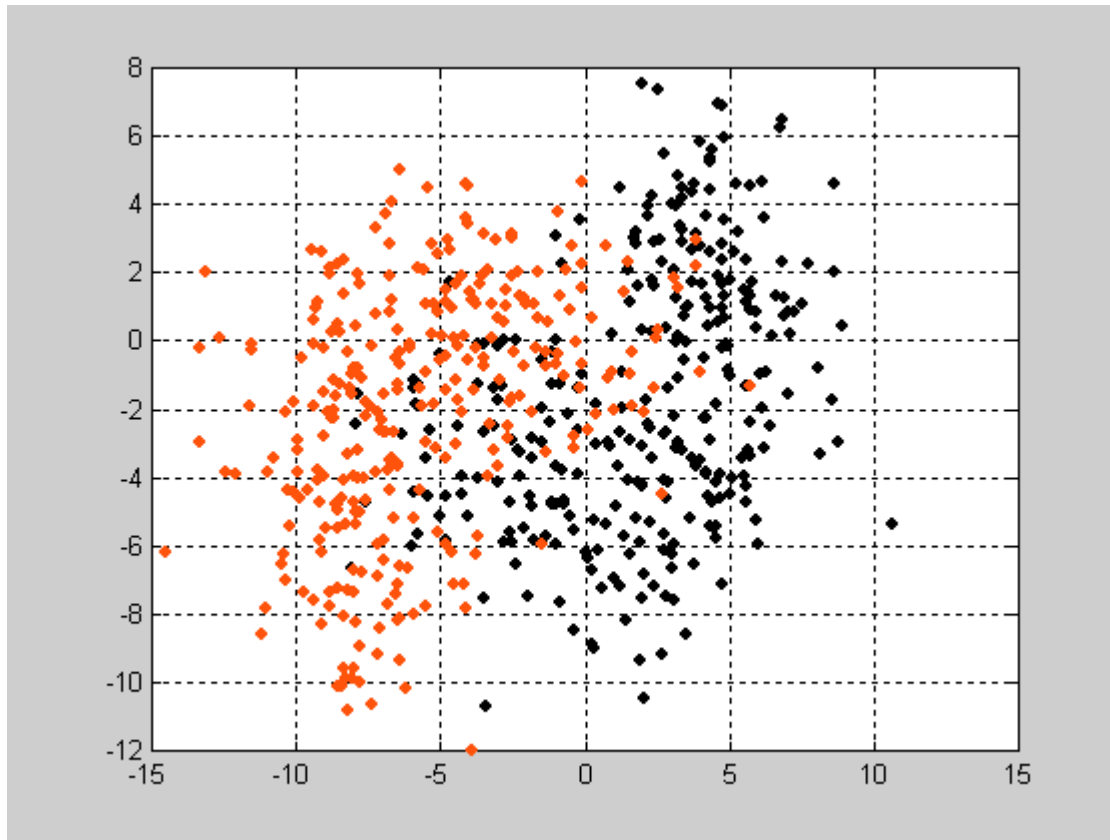
True label ω_1

True label ω_2

The correlation ρ can be calculated as the averaged pairwise correlation between the oracle outputs

An example:

banana-shaped data (gendatb routine from Matlab toolbox PRtools)



Training

N = 600 data points

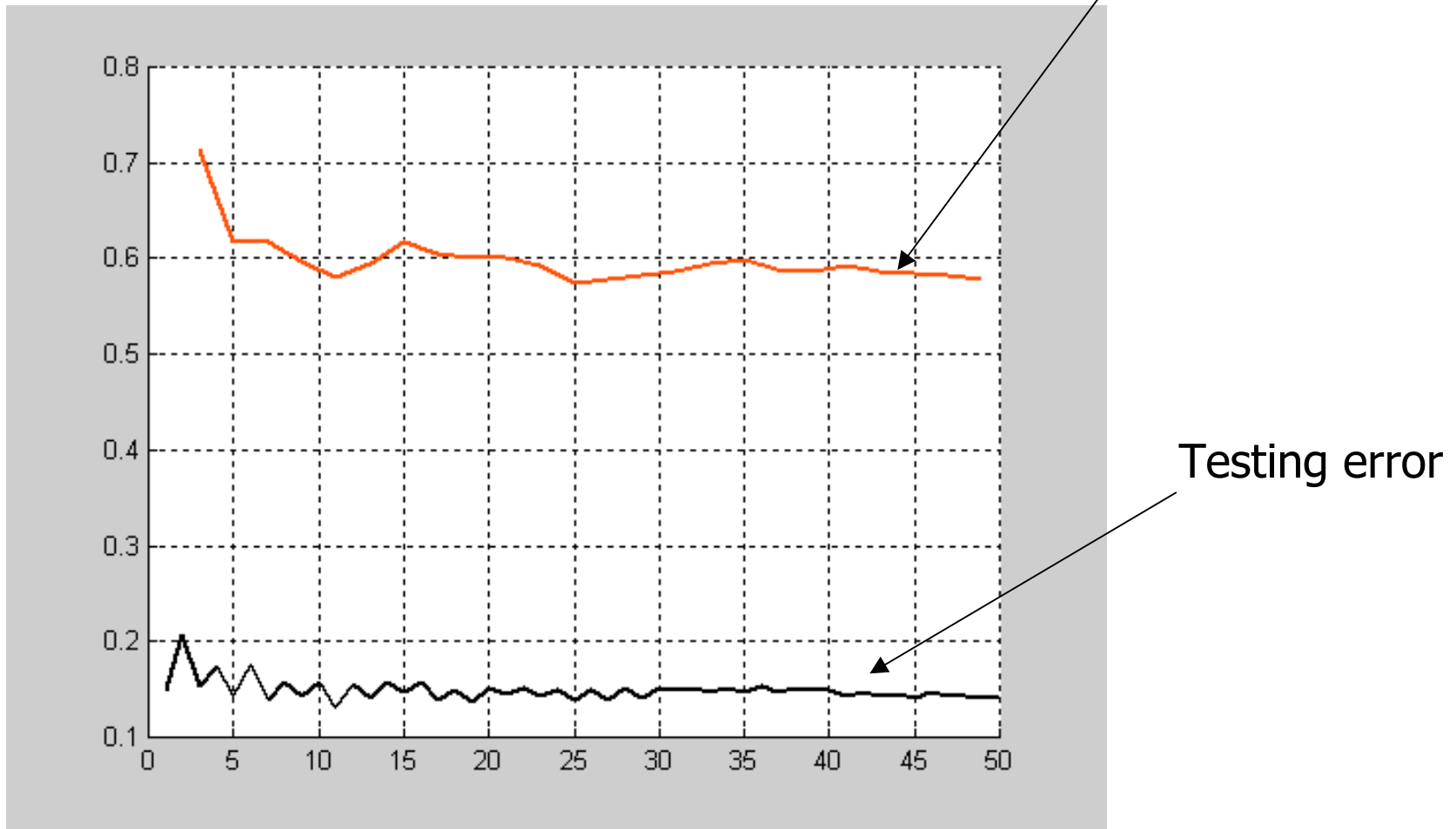
Testing

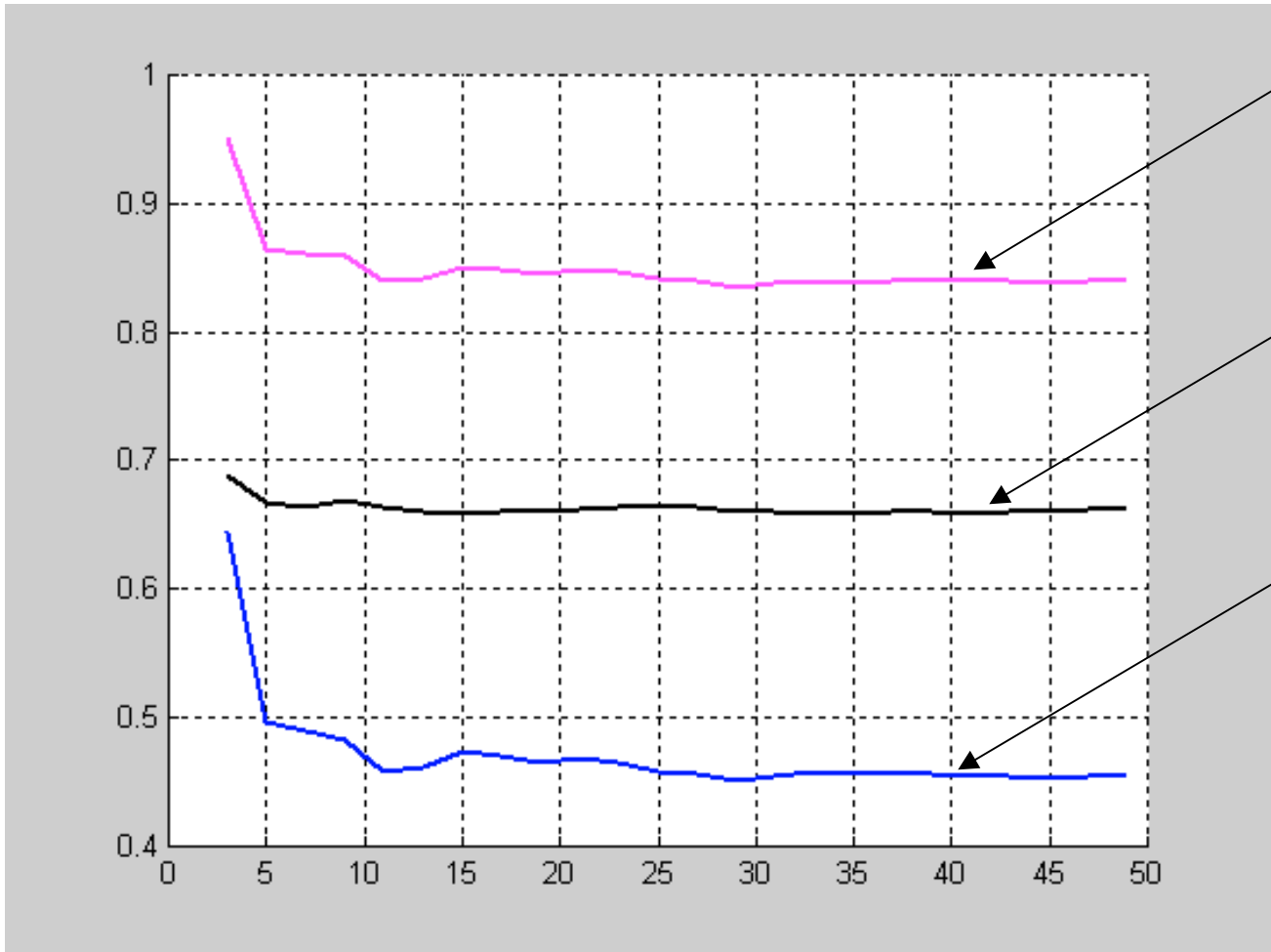
(a separate set)

N = 600 data points

The idea was to avoid using OB estimates which anyway simulate estimates on an independent testing set of the same size

Simple bagging, L = 50 classifiers



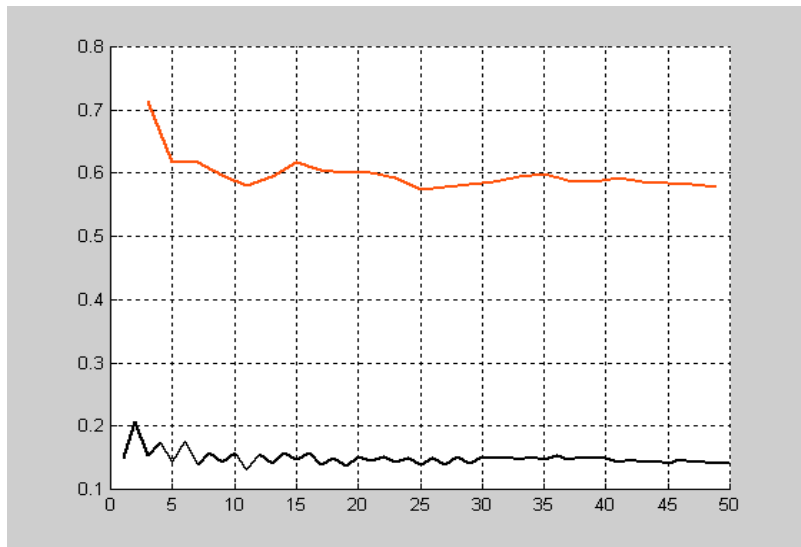


Q

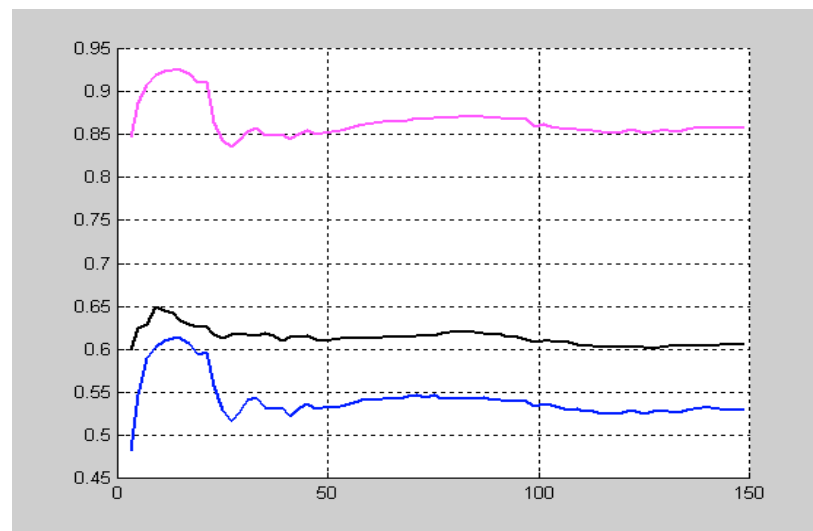
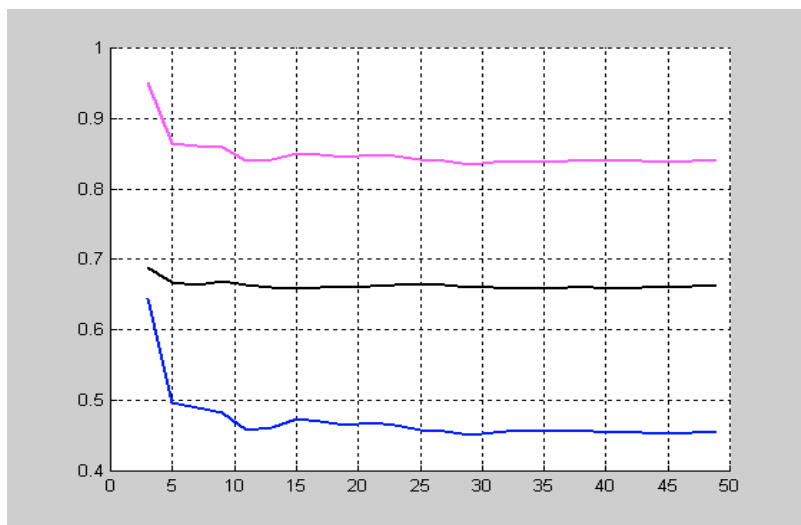
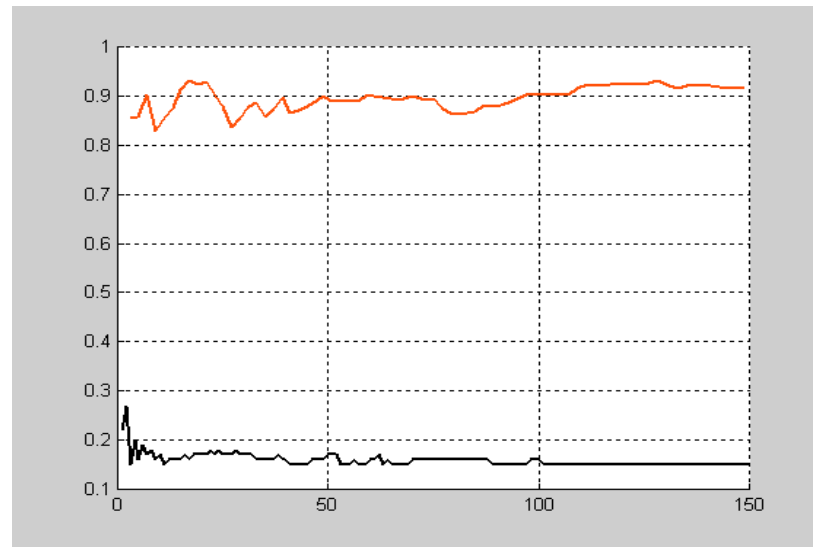
strength

correlation

L=50, N=600



L=150, N=100



Is strength related to accuracy?

true labels

guessed labels

	D_1	D_2	D_3
1	1	1	1
1	2	0	2
1	1	1	1
1	1	1	2
1	1	1	1
2	1	0	2
2	2	1	1
2	1	0	2
2	2	1	2
2	2	1	1
	7/10	7/10	6/10
	4/10	4/10	1/3

$P_{\theta}(D(x)=\omega_1)$

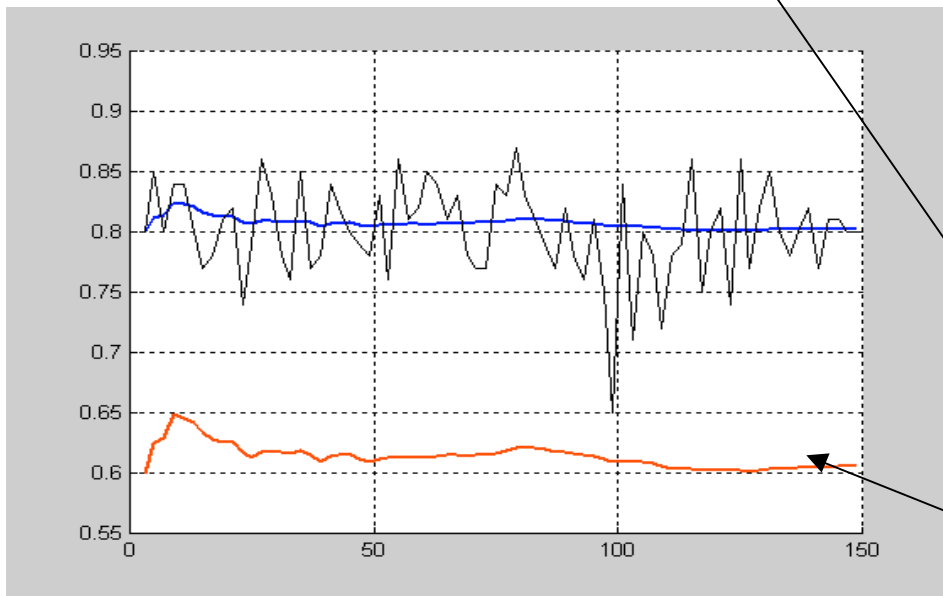
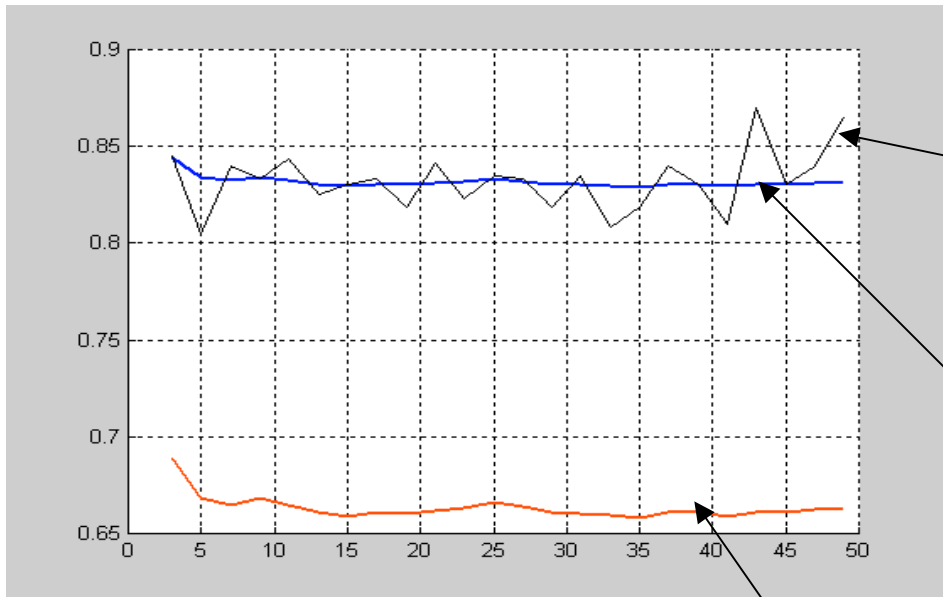
$P_{\theta}(D(x)=\omega_2)$

$2 \times (7/10) - 1$

$(2/10) \times (20/3) - 1$

accuracy

strength



individual testing error

averaged individual testing error

strength

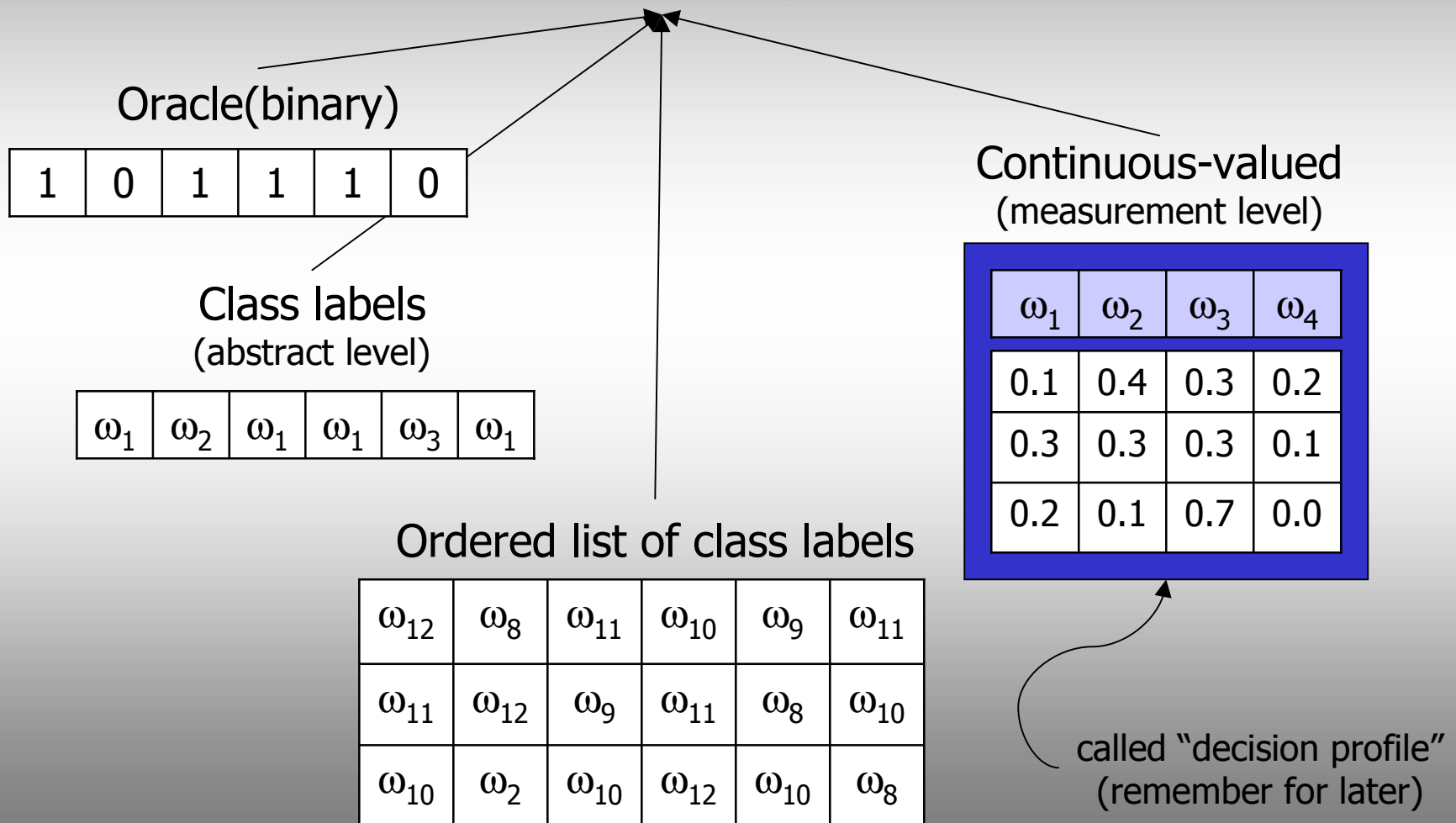
Vietri sul Mare, 27 09 02

Part 2: Non-pairwise diversity measures

0. A note on pairwise diversity (ρ) for random forests
 - Measures based on a single data point + averaging (entropy, spread, KW variance)
 - Interrater agreement (kappa for multiple raters)
 - Measures based on difficulties of the data points
 - Relationship with accuracy
 - Open problems

Now we look at the whole ensemble of classifiers.

Classifier outputs



- Measures based on a single data point (case, instance, example, object, whatever) and subsequently averaged over the whole data set.

2. Measures based on all data points.

For oracle outputs and $L = 8$ classifiers, are these diverse?

1	1	1	1	1	1	1	1
---	---	---	---	---	---	---	---

No-o-o-o-o-o-o-o!

0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---

Nope.

1	0	1	0	1	0	1	0
---	---	---	---	---	---	---	---

Yes.

1	1	1	1	0	0	0	0
---	---	---	---	---	---	---	---

Yes.

ENTROPY (oracle outputs)

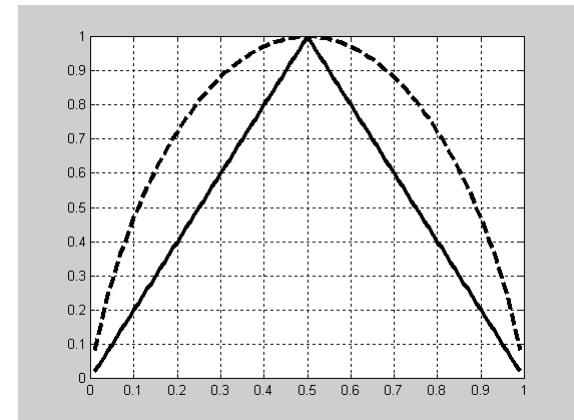
How do we measure how far we are from the desired pattern of $L/2$ 0's and $L/2$ 1s for N objects?

$$E = \frac{1}{\lceil L/2 \rceil N} \sum_k \left[\min \{ \Sigma 0's, \Sigma 1's \} \right]_k$$

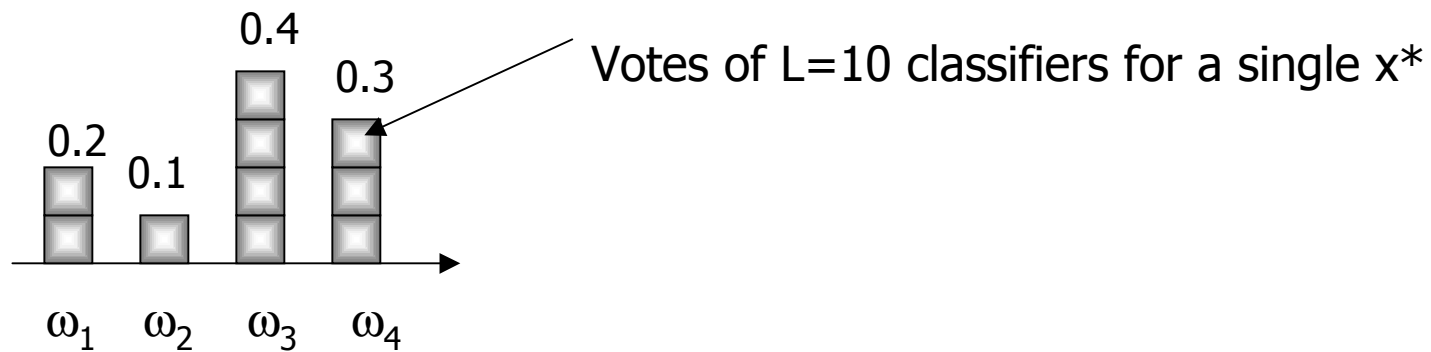
Consider the output 0 or 1 as a random variable with relative frequencies $p_0 = (\Sigma 0's) / L$ and $p_1 = (\Sigma 1's) / L$, respectively. Then the (proper) formula for the entropy of the distribution, averaged across the N data points will be

$$H = - \frac{1}{N} \sum_k \left[p_0 \log p_0 + p_1 \log p_1 \right]_k$$

[Cunningham Carney, 2000]



ENTROPY (label outputs)



$$H = -1/N \sum_k \left[\sum_i p_i \log p_i \right]_k$$

Breiman's Bias-Variance decomposition, 1996

Assume that classifier output for a given x^* is a random variable with p.m.f. $P(\omega_1|x^*,D), \dots, P(\omega_c|x^*,D)$. The classification error is

$$P(\text{error}|x^*) = 1 - \sum_j P(\omega_j | x^*) P(\omega_j | x^*, D)$$

$$= 1 - \{ P(\omega_B | x^*) - P(\omega_B | x^*) - \sum_j P(\omega_j | x^*) P(\omega_j | x^*, D) \}$$

$$= [1 - P(\omega_B | x^*)] + \sum_j [P(\omega_B | x^*) - P(\omega_j | x^*)] P(\omega_j | x^*, D)$$

$$= P_B(x^*) + \sum_j [P(\omega_B | x^*) - P(\omega_j | x^*)] P(\omega_j | x^*, D)$$

$$[P(\omega_B | x^*) - P(\omega_s | x^*)] P(\omega_s | x^*, D)$$

bias

$$+ \sum_{j \neq s} [P(\omega_B | x^*) - P(\omega_j | x^*)] P(\omega_j | x^*, D)$$

spread

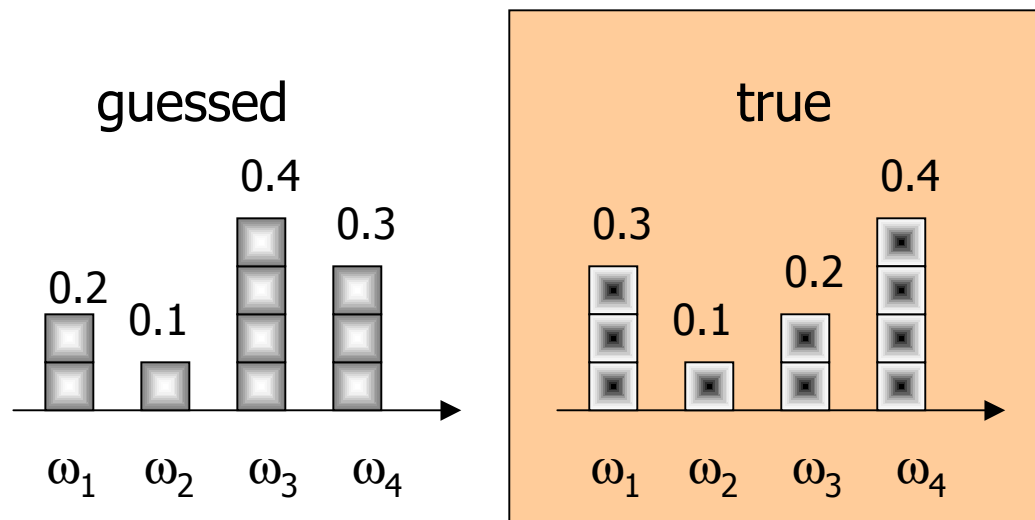
$$P_B(x^*)$$

$$+ [P(\omega_B | x^*) - P(\omega_s | x^*)] P(\omega_s | x^*, D) \quad \text{(bias)}$$

$$+ \sum_{j \neq s} [P(\omega_B | x^*) - P(\omega_j | x^*)] P(\omega_j | x^*, D) \quad \text{(spread)}$$

Is the spread related to diversity?

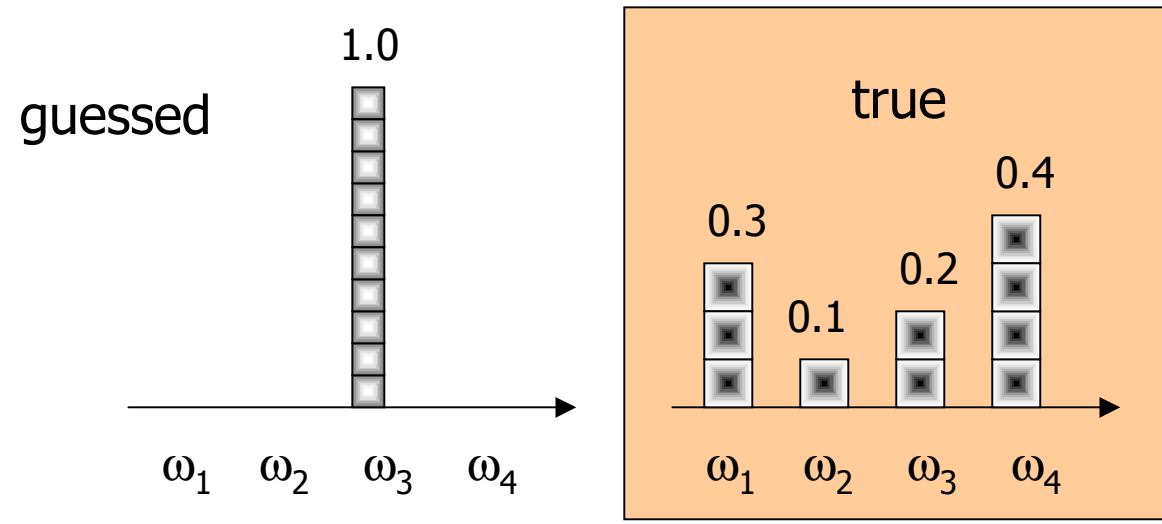
An example: If we drew a classifier at random from the distribution P_D ,



$$P(\text{error}|x) = 0.6 + [0.4-0.2] 0.4 + [0.1 \times 0.2 + 0.3 \times 0.1] = 0.73$$

$$P(\text{error}|x) = 0.6 + \frac{[0.4-0.2] 0.4}{0.08} + \frac{[0.1 \times 0.2 + 0.3 \times 0.1]}{0.05} = \underline{0.73}$$

Take majority vote. This means "decide always ω_s for x^* ".



$$P(\text{error}|x) = 0.6 + \frac{[0.4-0.2] 1.0}{0.2} = \underline{0.80}$$

KW variance (label outputs)

[Kohavi Wolpert, 1996, Bias plus variance decomposition for zero-one loss functions]

The c-class case:

$$P(\text{error}|x) = \text{bias}^2(x) + \text{variance}(x) + \text{noise}^2(x)$$

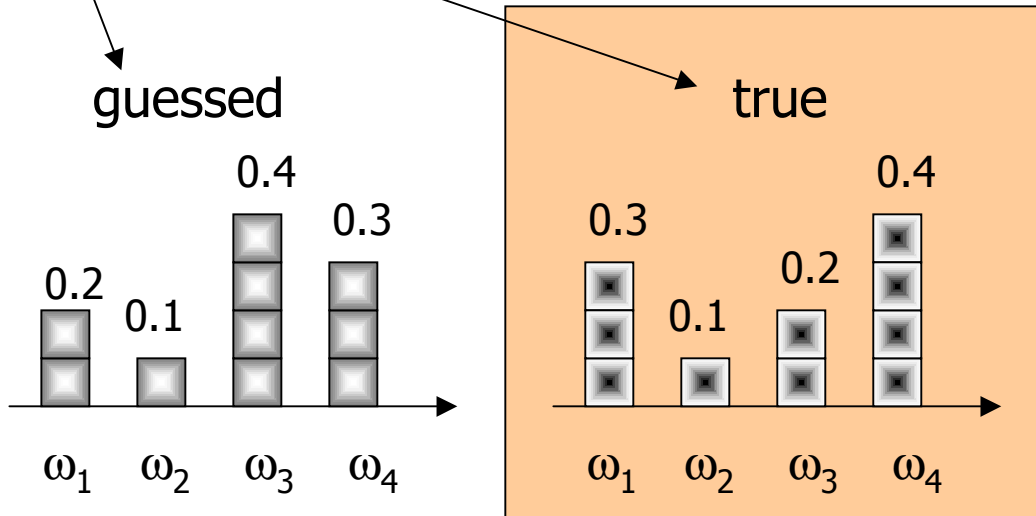
$$\underline{\text{bias}^2(x)} \quad \frac{1}{2} \sum_{\omega} (P_{\text{true}}(\omega|x) - P_{\text{guessed}}(\omega|x))^2$$

$$\underline{\text{variance}(x)} \quad \frac{1}{2} (1 - \sum_{\omega} (P_{\text{guessed}}(\omega|x))^2)$$

$$\underline{\text{noise}^2(x)} \quad \frac{1}{2} (1 - \sum_{\omega} (P_{\text{true}}(\omega|x))^2)$$

$$\text{bias}^2(x) = \frac{1}{2} \sum_{\omega} (P_{\text{true}}(\omega|x) - P_{\text{guessed}}(\omega|x))^2$$

$$\frac{1}{2} (0.3 - 0.2)^2 + (0.4 - 0.2)^2 + (0.3 - 0.4)^2 = \mathbf{0.03}$$



$$\text{variance}(x) = \frac{1}{2} (1 - \sum_{\omega} (P_{\text{guessed}}(\omega|x))^2)$$

$$\frac{1}{2} [1 - ((0.2)^2 + (0.1)^2 + (0.4)^2 + (0.3)^2)] = \mathbf{0.35}$$

$$\text{noise}^2(x) = \frac{1}{2} (1 - \sum_{\omega} (P_{\text{true}}(\omega|x))^2)$$

$$\frac{1}{2} [1 - ((0.3)^2 + (0.1)^2 + (0.2)^2 + (0.4)^2)] = \mathbf{0.35}$$

KW variance (oracle outputs)

Consider again the output 0 or 1 as a random variable with relative frequencies $p_0 = (\sum 0's) / L$ and $p_1 = (\sum 1's) / L$, respectively. Then the variance is

$$\text{variance}(x) = \frac{1}{2} (1 - (p_0)^2 - (p_1)^2)$$

Averaging across the whole data set,

$$KW = 1/(N \times L^2) \sum_k [(\sum 0's) \times (\sum 1's)]_k$$

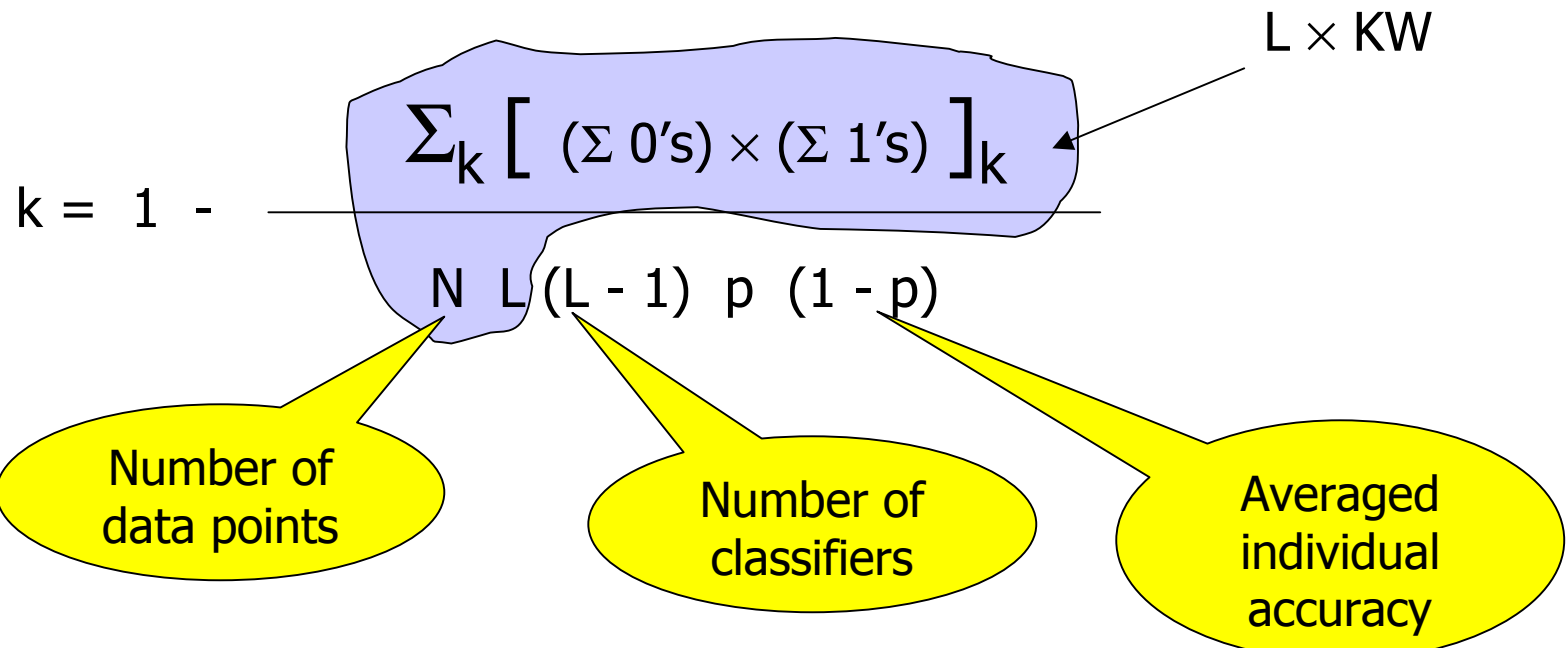
Curiously, KW and the averaged pairwise disagreement measure are related through

$$KW = (L-1)/(2L) D_{av}$$

- Measures based on a single data point (case, instance, example, object, whatever) and subsequently averaged over the whole data set.

2. Measures based on all data points.

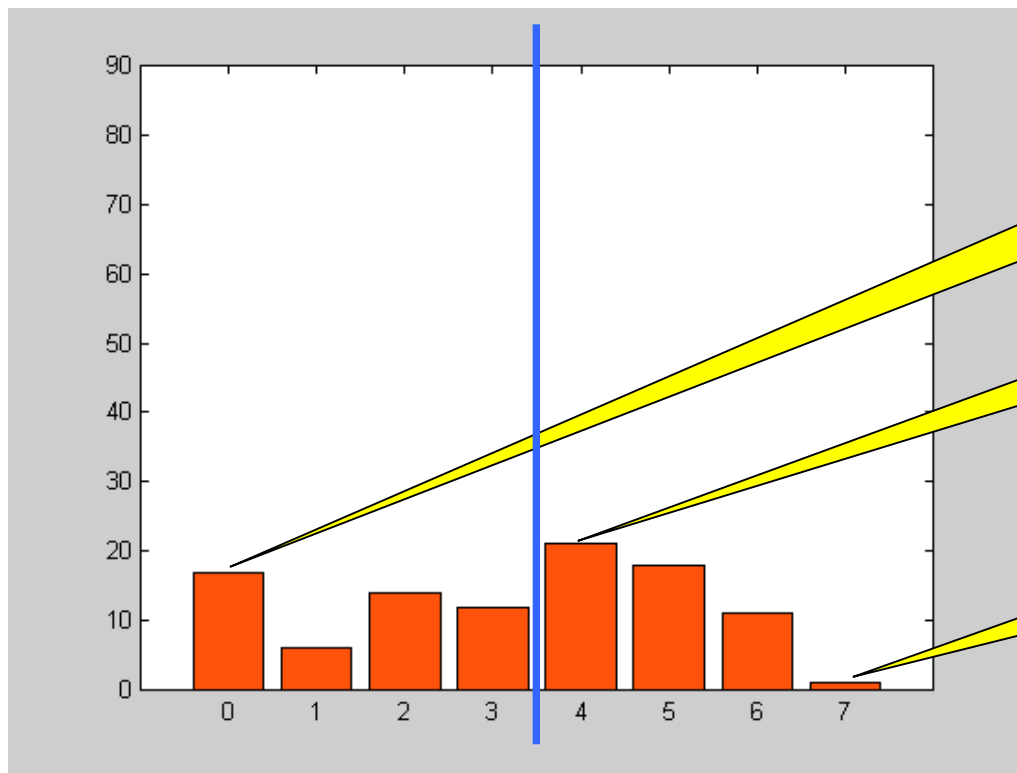
Interrater agreement, kappa, (oracle outputs)



Measure of difficulty θ

[Hansen Salamon, 1990]

Define a random variable X = proportion of classifiers which correctly classify a randomly drawn sample x . Let $L = 7$.

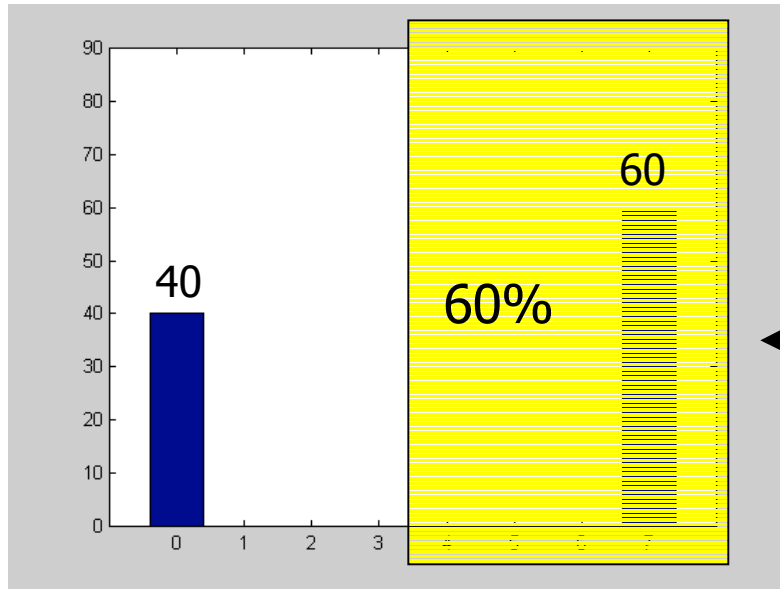


Number of points misclassified by all 7

Number of points recognized by any 4

Number of points recognized by all 7

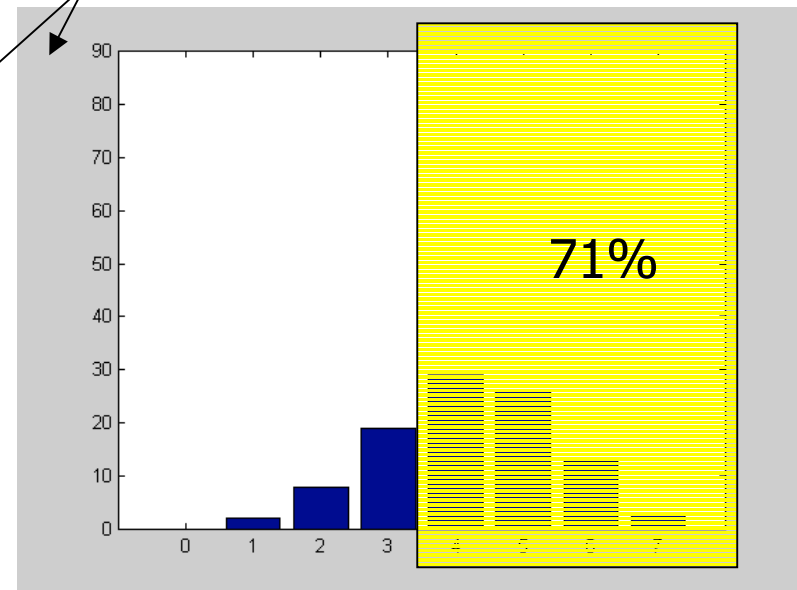
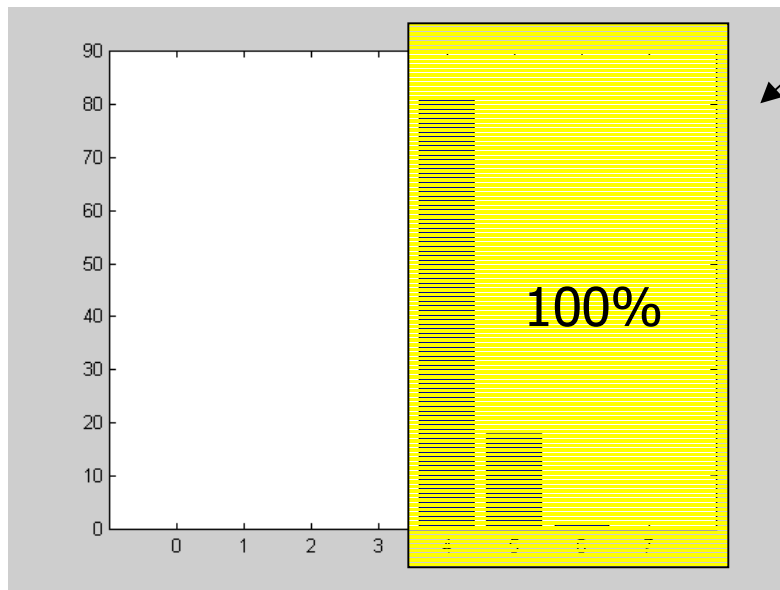
$L = 7, p = 0.6$



independent

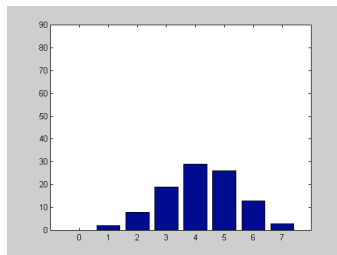
identical

negatively dependent



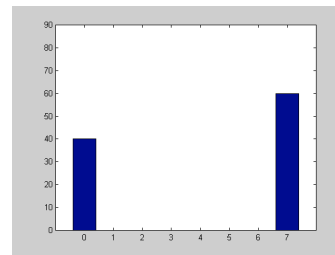
measure of diversity $\theta = \text{Var}(X)$

independent



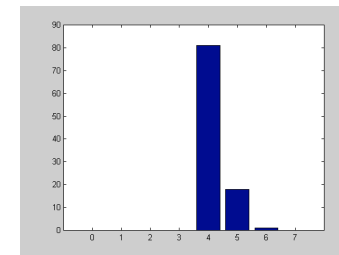
$\theta = 0.034$

identical



$\theta = 0.240$

diverse

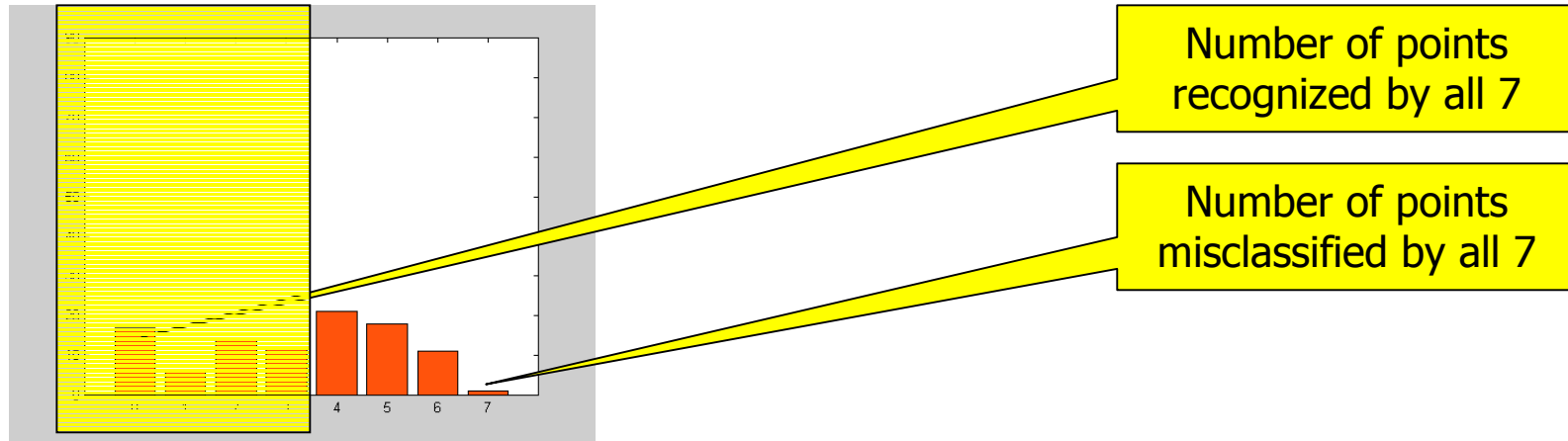


$\theta = 0.004$

Generalized diversity

[Partridge Krzanowski, 1997]

Define a random variable Y = proportion of classifiers which **misclassify** a randomly drawn sample x . ($Y = 1 - X$ defined before)



Denote by p_i the probability that $Y = i / L$, and by $p(k)$ the probability that k randomly chosen classifiers will fail on a randomly drawn x .

$$p(1) = \sum_i p_i \times i / L \text{ (the probability of single classifier failing)}$$

$$p(2) = \sum_i p_i \times i (i - 1) / (L (L - 1)) \text{ (the probability that two randomly chosen classifiers will fail together)}$$

$$GD = 1 - p(1)/p(2)$$

Coincidence failure diversity

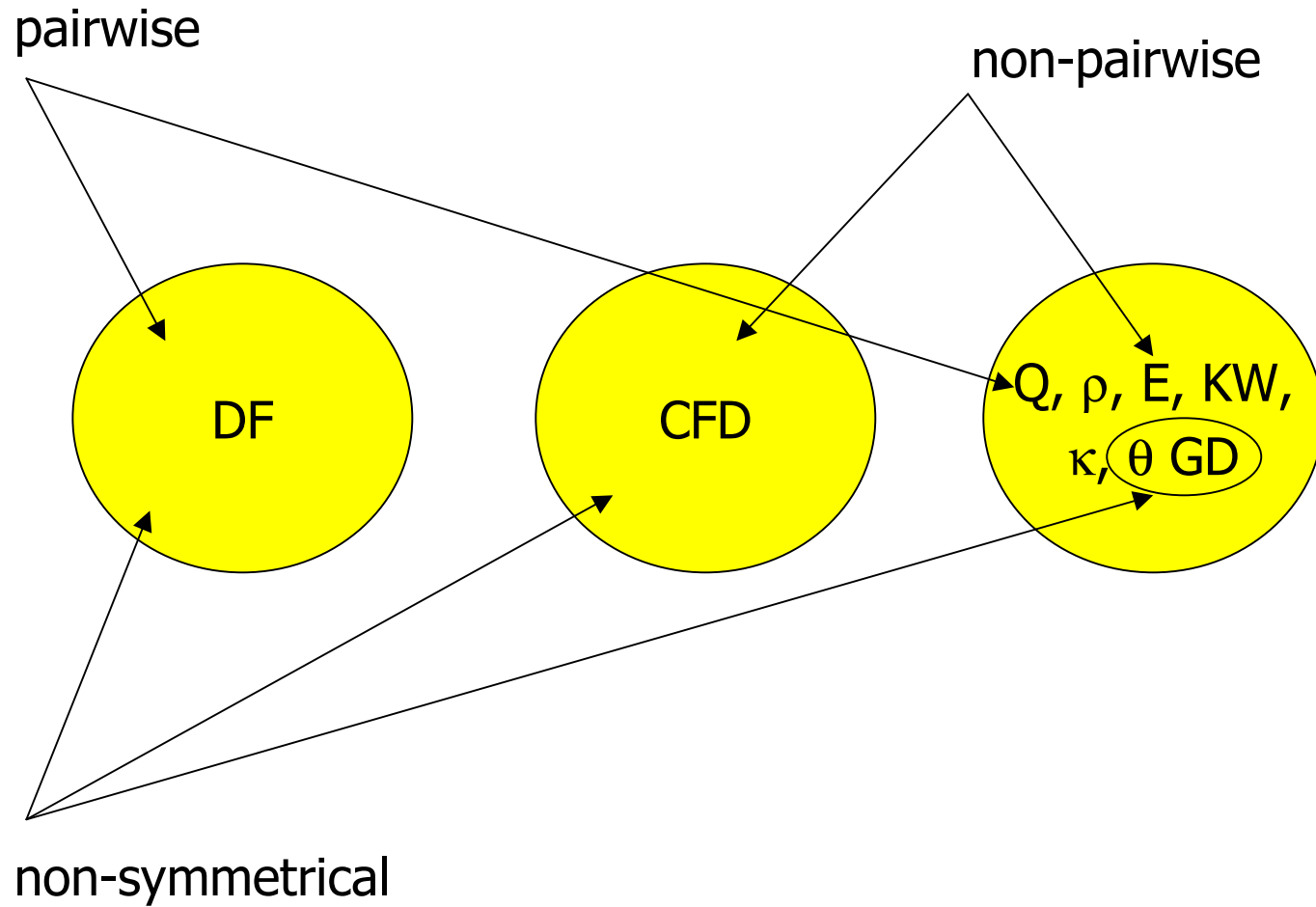
$$CFD = \begin{cases} 0, & \text{if } p_0 = 1, \\ 1/(1 - p_0) \sum_i p_i \times (L - i) / (L - 1), & \text{if } p_0 < 1 \end{cases}$$

Relationship between diversity and accuracy

Correlations between the improvement on the single best classifier and some diversity measures (WBC)

	Q	ρ	Dis	DF	κ	θ	GD	CFD
MAJ	-17	-21	33	18	-20	35	28	38
NB	-15	-20	32	20	-18	37	26	36
BKS	-15	-17	17	5	-15	16	18	18
WER	-15	-17	17	5	-16	17	19	18
MAX	-1	-0	20	38	0	45	7	11
AVR	-13	-15	34	33	-14	47	22	30
PRO	-11	-11	29	33	-11	44	18	24
DT	-12	-15	32	30	-14	44	22	29

Relationship between diversity measures





Open problems

- How to **narrow down** the study? (Use a specific methodology for building the ensemble)
- Some **theory** would not go amiss.
- Diversity for **label outputs** and **continuous-valued outputs** might lead somewhere.

The difficulty comes from the fact that the output of the classifiers are vectors

	ω_1	ω_2	ω_3	ω_4
D_1	0.1	0.4	0.3	0.2
D_2	0.3	0.3	0.3	0.1
D_3	0.2	0.1	0.7	0.0



similarity between distributions
(pairwise)

WAKE UP!!!

Vietri sul Mare, 27 09 02