# Minimum Redundancy and Maximum Relevance Feature Selection and Its Applications

#### Hanchuan Peng

Janelia Farm Research Campus, Howard Hughes Medical Institute

## Outline

- What is mRMR feature selection
- Applications in cancer classification
- Applications in image pattern recognition
- Theoretical basis of mRMR
- Combinations with other methods
- How to use mRMR programs











TYPE	ACRONYM	FULL NAME	Formula
RETE	MID	Mutual information difference	$\max_{i \in \Omega_{S}} [I(i,h) - \frac{1}{ S } \sum_{j \in S} I(i,j)]$
DISC	MIQ	Mutual information quotient	$\max_{i\in\Omega_{S}}\{I(i,h)/[\frac{1}{ S }\sum_{j\in S}I(i,j)]\}$
	FCD	F-test correlation difference	$\max_{i \in \Omega_{S}} [F(i,h) - \frac{1}{ S } \sum_{j \in S}  c(i,j) ]$
NUOUS	FCQ	F-test correlation quotient	$\max_{i\in\Omega_{S}}\{F(i,h)/[\frac{1}{ \mathcal{S} }\sum_{j\in\mathcal{S}} c(i,j) ]\}$
CONTI	FDM	F-test distance multiplicative	$\max_{i\in\Omega_{S}}[F(i,h)\cdot \frac{1}{ S }\sum_{j\in S}d(i,j)]$
	FSQ	F-test similarity quotient	$\max_{i\in\Omega_S} \{F(i,h)/[\frac{1}{ S }\sum_{j\in S}\frac{1}{d(i,j)}]\}$



 PCA and ICA: Features are orthogonal or independent, but not in the original feature space.

# Outline

- What is mRMR feature selection
- Applications in cancer classification
- Applications in image pattern recognition
- Theoretical basis of mRMR
- Combinations with other methods
- How to use mRMR programs









	Method	3	6	9	12	15	18	24	30	40	50	60	70	80	90	100
	Baseline	55	47	46	38	34	27	19	28	22	19	15	14	11	8	8
LDA	MID	50	43	32	29	30	29	22	15	13	10	10	- 9	7	8	- 9
	MIQ	43	43	34	27	23	21	18	16	11	11	6	4	6	6	4
	Baseline	56	55	49	37	33	33	27	35	29	30	23	20	18	14	13
SVM	MID	45	42	33	33	25	25	29	25	26	22	20	13	10	12	9
	MIQ	38	30	34	33	27	26	24	21	14	15	17	10	7	11	9

	M															
Data Sets	Method	3	6	- 9	12	15	18	21	24	27	30	36	42	48	54	6
NCL	Baseline	34	25	23	25	19	17	18	15	14	12	12	12	13	12	1
ner	MRMR (MIQ)	35	22	22	16	12	11	10	8	5	3	4	4	2	2	
Lymphoma	Baseline	58	52	44	39	44	17	17	14	16	13	11	10	13	10	1
2)	MRMR (MIQ)	24	17	7	8	4	2	1	2	4	3	2	2	2	2	
	)CV testing r	esult	ts (‡	#err	or)	for t	oina	rizeo	d NC	CI aı	nd					

Data	Method	NB	LDA	SVM	LR	Literature
NCI	Baseline	18.33	26.67	25.00		14.63 <sup>a</sup>
NCI	MRMR	1.67	13.33	11.67		5-class: 0 <sup>b</sup> , 0 <sup>b</sup>
Lymphomo	Baseline	17.71	11.46	5.21		$\frac{1}{2}$ aloss: $2.4^{\circ}$ 0
Lymphonia	MRMR	3.13	1.04	1.04		J-Class. 2.4 , 0
Lung	Baseline	10.96	10.96	10.96		
Lung	MRMR	2.74	5.48	5.48		
Thild Laukamia	Baseline	29.46	7.14	11.61		5 36 d
	MRMR	13.39	2.68	6.25		5.50
Laukamia	Baseline	0	1.39	1.39	1.39	0 e
Leukenna	MRMR	0	0	0	0	1.39 <sup>f</sup>
Calan	Baseline	11.29	11.29	11.29	11.29	9.68 <sup>e</sup>
Colon	MRMR	6.45	8.06	9.68	9.68	6.45 <sup>g</sup>

Ding & Peng, CSB2003, JBCB2005 17

# Outline What is mRMR feature selection Applications in cancer classification Applications in image pattern recognition Theoretical basis of mRMR Combinations with other methods How to use mRMR programs







13-16	11-12	9-10	7-8	4-6	1-3	Actual stage range/Predicted stage range
0	0	0	0	3	447	1-3
0	0	0	0	446	0	4-6
0	0	0	379	0	1	7-8
0	1	372	0	0	0	9-10
0	435	0	0	1	1	11-12
362	2	0	0	0	2	13-16
Bio, 2007	ig et al, BMC Cel	Per	%	/ >99	curacy	Overall ac

**Tier 1 Accuracy: Developmental Stage Prediction** 

Stage range	Set size	Set details	Score	Eigen	LeNet	LeNet +	Wavelet	Wavelet
				emoryo		mRMR	emoryo	emoryo + mKMP
4-6	98	78CB+20Subset	Rs	70	75	70	70	100
			RA	82	85	83	82	98
7-8	69	49TMA + 20PEA	Rs	60	50	55	60	100
			RA	67	59	80	67	100
9-10	69	49TMP + 20PEP	Rs	65	50	60	65	100
			RA	71	52	83	71	100
11-12	61	41 HPP + 20 PMP	R <sub>S</sub>	45	50	90	55	95
			RA	59	69	90	64	98
13-16	32	20EBNS+12VNC	Rs	83	58	67	83	100
			P.	84	50	81	84	100
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Reco	atures are used. le data set for ordium; PEP, p VNC, ventral gnition rates	Top 20 features are selected whe all classes. (Annotation terms: C reocephalic ectoderm primordia nerve cord.) (%) of feature extrction/sele	n using mRM B, cellular b n; HPP, hin ction metho	dR feature selecti lastoderm; TMA, dgut proper prim	on. R <sub>8</sub> : recogn trunk mesode ordium; PMP, : data sets w	tion rate on the s rm anlags; PEA, posterior midgu ith multi-object	maller class only posterior endod t primordium; E ive (M.O.) and	R <sub>A</sub> : overall recognition rm anlage; TMP, trani BNS, embryonic centra rotations
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recoj Stage range	atures are used. le data set for ordium; PEP, p VNC, ventral gnition rates Set size	Top 20 features are selected whe all classes. (Annotation terms: C receptualic ectoderm primordiu nerve cord.) (%) of feature extrction/sele Set details	n using mRM B, cellular b n; HPP, hin ction metho Sco	dR feature selecti lastoderm; TMA, dgut proper prim ods on synthetic ore Eigen	on. R <sub>8</sub> : recogn trunk mesode ordium; PMP, : data sets w LeNet	tion rate on the 1 rm anlage; PEA, posterior midgu ith multi-object LeNet +	maller class only posterior endod t primordium; E ive (M.O.) and Wavelet	RA: overall recognition rm anlage; TMP, truni BNS, embryonic centra notations Wavelet
Top 20 eigen fea rate on the who mesoderm prims nervous system; <b>Table 2.</b> Recop	atures are used. de data set for ordium; PEP, 5 VNC, ventral 1 gnition rates Set size	Top 20 features are selected whe all classes. (Annotation terms: C receptualic ectoderm prinzerdiu nerve cord.) (%) of feature extrction/sele Set details	n using mR3 B, cellular b n; HPP, hin ction metho Sco	dR feature selecti lastoderm; TMA, dgat proper prim ods on synthetic ore Eigen embryo	m. R <sub>5</sub> : recogn trunk mesode ordium; PMP, : data sets w LeNet	tion rate on the s rm anlage; PEA, posterior midgu th multi-object LeNet + mRMR	maller class only posterior ended t primordium; E ive (M.O.) and Wavelet embryo	RA: overall recognition rm anlage; TMP, trusi BNS, embryonic centra totations Wavelet embryo + mRMI
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recoy Stage range	tures are used. de data set for : ordium; PEP, g VNC, ventral i gnition rates : Set size 152	Top 20 features are selected whi all classes. (Annotation terms: C recorphalic ectoderm primordia terve cond.) (%) of feature extretion/sele Set details 78CB + 20Subset + 54Bot	n using mRM B, cellular b n; HPP, hin ction metho Sco h R <sub>S</sub>	dR feature selecti lastodern; TMA, dgat proper prim ods on synthetic ore Eigen embryo 49	m. R <sub>5</sub> : recogn trunk mesode ordium; PMP, : data sets w LeNet 42	tion rate on the s rm anlage; PEA, posterior midgu th multi-object LeNet + mRMR 55	maller class only posterior endod t primordium; E ive (M.O.) ann Wavelet embryo 49	R <sub>4</sub> , overall recognition R <sub>4</sub> , overall recognition BNS, embryonic centra notations Wavelet embryo + mRMI 76
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recop Stage range 4-6	tures are used. le data set for ordium; PEP, 5 VNC, ventral 1 gnition rates - Set size 152	Top 20 features are selected whe all classes. (Annotation terms: C recephalic excidents primordia nerve cord.) (%) of feature extrction/sele Set details 78CB + 20Subset + 54Bot	n using mRM B, cellular b n; HPP, hin ction metho Sco h R <sub>S</sub>	dR feature selecti lastoderm; TMA, dgut proper prim ods on synthetis ore Eigen embryo 49 , 64	on. R <sub>3</sub> : recogn trunk mesode ordium; PMP, : data sets w LeNet 42 63	tion rate on the s m anlags; PEA, posterior midgu th multi-object LeNet + mRMR 55 71	maller class only posterior endod i primordium; E ive (M.O.) ann Wavelet embryo 49 64	Not over all recognitions rm anlage; TMP, trani lNS, embryonic centra totations Wavelet embryo + mRM1 76 82
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recop Stage range 4-6 7-8	tures are used. le data set for - ordium; PEP, ; VNC, ventral : gnition rates - Set size 152 89	Top 20 features are selected whe all classes. (Annotation series of recorphale extoderm primordia preve cont.) (%) of feature extrection/sele Set details 78CB + 20Subset + 54Bot 49TMA + 20PEA + 20Bo	n using mR3 B, cellular b n; HPP, hin ction metho Sco h R <sub>S</sub> h R <sub>S</sub>	dR feature selecti lastoderm; TMA, ägat proper prim ods on synthetic ore Eigen embryo 49 64 45	on. R <sub>8</sub> : recogn trunk mesode ordium; PMP, e data sets w LeNet 42 63 50	tion rate on the r m anlage; PEA, posterior midgu ith multi-object LeNet + mRMR 55 71 58	maller class only posterior endod t primordium; E ive (M.O.) ann Wavelet embryo 49 64 45	No R <sub>A</sub> : overall recognition m anlage: TMP, truni BNS, embryonic centra notations Wavelet embryo + mRMI 76 82 83
Top 20 eigen fea rate on the who mesoderm prim nervous system; <b>Table 2.</b> Recoy Stage range 4-6 7-8	tures are used. le data set for ordium; PEP, 5 VNC, ventral 1 gnition rates - Set size 152 89	Top. 2D faitures are selected whe all classes, (Annotation terms: C- receptual: ecitodems primordiu terres cond.) (%) of feature extretion/sele Set details 78CB + 20Subset + 54Bot 49TMA + 20PEA + 20Bo	n using mRM B, cellular b n; HPP, him ction metho Sco h R <sub>S</sub> h R <sub>S</sub> R,	dR feature selecti lastoderm; TMA, ágat proper prins ods on synthetic ore Eigen embryo 49 64 45 40	on. R <sub>4</sub> : recogn trunk mesede ordium; PMP, data sets w LeNet 42 63 50 46	tion rate on the e m anlage; PEA, posterior midgu th multi-object LeNet + mRMR 55 71 58 62	maller class only posterior endod t primordium; E ive (M.O.) ann Wavelet embryo 49 64 45 40	R <sub>s</sub> : overall recognition R <sub>s</sub> : overall recognition BNS, embryonic centra totations Wavelet embryo + mRMI 76 82 83 83
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recop Stage range 4-6 7-8 9-10	tures are used. le data set for ordium; PEP, ; VNC, ventral : gnition rates : Set size 152 89 87	Top 20 fattures are selected whe all classes, (Annotation terms: C receptualic actioners primordia serve cost.) (%) of feature extruction/sele 78CB + 205ubset + 54Bot 49TMA + 20PEA + 20Bo 49TMP + 20PEA + 20Bo	n using mRM B, cellular b n; HPP, hin ztion metho Sco h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub>	dR feature selecti lastoderm; TMA, dgat proper prim ods on synthetic ore Eigen embryo 49 64 45 40 42	n. R <sub>4</sub> : recogn trunk mesode ordium; PMP, : data sets w LeNet 63 50 42 63 50 46 34	tion rate on the r rm anlage; PEA, posterior midgu th multi-object LeNet + mRMR 55 71 58 62 42	maller class only posterior endod t primordiam; E ive (M.O.) and Wavelet embryo 49 64 45 40 39	R <sub>4</sub> : overall recognition R <sub>4</sub> : overall recognition BNS, embryonic centra totations Wavelet embryo + mRM1 76 82 83 85 79
Top 20 eigen fea rate on the who mesoderm prim nervous system; <b>Table 2.</b> Recop Stage range 4-6 7-8 9-10	tures are used. le data set for ordium; PEP, ; VNC, ventral ; gnition rates ; Set size 152 89 87	Top 20 fattures are selected whe all classes, (Annotation terms: Co- copyhale catologne primordius terrer cord.) (%) of Cratture extretion/sele Set details 78CB + 20Subset + 54Bot 49TMA + 20PEA + 20Bo 49TMP + 20PEP + 18Bot	n using mR3 B, cellular b n; HPP, hin ction metho Sco h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> R <sub>S</sub>	dR feature selecti lastodern; TMA, dgut proper prim ods on synthetic ore Eigen embry 49 64 45 40 42 49	on. R <sub>6</sub> : recogn trank mesode ordium; PMP, c data sets w LeNet 6 42 63 50 46 34 33	tion rate on the r rm anlage; PEA, posterior midgu th multi-object LeNet + mRMR 55 71 58 62 42 55	maller class only posterior endod ; primordium; E ive (M.O.) and Wavelet embryo 49 64 45 40 39 47	Ros overall recognition Rac overall recognition RNS, embryonic centra notations Wavelet embryo + mRMI 76 82 83 85 79 80
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recop Stage range 4-6 7-8 9-10 11-12	atures are used. de data set for ordium; PEP, j VNC, ventral i gnition rates Set size 152 89 87 75	Top 20 fattures are selected whe all classes, (Annotation terms: C receptable schedures primordia serve cost). Set details 78CB + 205kubset + 54Bot 49TMA + 20PEA + 20Bo 49TMP + 20PEA + 20Bo 49TMP + 20PEP + 18Bot 41HPP + 20PMP + 14Bot	n using mR3 B, cellular b n; HPP, hin ction metho Sco h R <sub>S</sub> R <sub>A</sub> h R <sub>S</sub> R <sub>A</sub> h R <sub>S</sub> R <sub>A</sub> h R <sub>S</sub>	dR feature selecti lastodern; TMA, dgat proper prim ods on synthetic pre Eigen embryo 49 64 45 40 42 49 38	n. R <sub>4</sub> : recogn trank mosed erdium; PMP, : data sets w LeNets 42 63 50 46 34 33 22	tion rate on the translager PEA, posterior midgu th multi-object LeNet + mRMR 55 71 58 62 42 55 59	maller class only posterior ended i primordium; E ive (M.O.) ann Wavelet embryo 49 64 45 40 39 47 38	Ros overall recognition Ros overall recognition INSS, enhyponic centra notations Wavelet embryo + mRMI 76 82 83 85 79 98 80 88
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Reco Stage range 4-6 7-8 9-10 11-12	tures are used. le data set for ordium; PEP, p VNC, ventral i gnition rates : Set size 152 89 87 75	Top 20 fattures are selected wh all classes, (Annotation terms: Co- copyration expension of the selection (selection) (%) of Crature extrection/selection Set details 78CB + 20Subset + 54Bot 49TMA + 20PEA + 20Bo 49TMP + 20PEA + 18Bot 41HPP + 20PEP + 18Bot	n using mRN B, cellular b m; HPP, him tion metho Sco h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> h R <sub>S</sub> R <sub>A</sub> h R <sub>S</sub>	dR feature selecti lastodern; TMA, dgut proper prim ods on synthetic ore Eigen embryo 49 64 45 40 42 49 38 33	LeNet 42 63 50 42 63 50 44 33 32 36	tion rate on the t m anlage; PEA, posterior midgu th multi-object LeNet + mRMR 55 71 58 62 42 55 59 55	maller class only posterior ended t primordiam; E ive (M.O.) ann Wavelet embryo 49 64 45 40 39 47 38 37	Ros overall secognition sec, overall secognition notations Wavelet embryo + mRMI 76 82 83 85 79 80 88 89
Top 20 eigen fea rate on the who mesoderm prim nervous system; Table 2. Recop Stage range 4-6 7-8 9-10 11-12 13-16	tures are used. le data set for ordium; PEP, 1 VNC, ventral 1 gnition rates : 152 89 87 75 97	Top 20 faitures are selected whe all classes. (Annotation series: C scephiale sciences primordius rerve cord.) Set details 78CB + 20Subset + 54Bot 49TMP + 20PEP + 18Bot 40TMP + 20PEP + 18Bot 41HPP + 20PEP + 14Bot 20EBNS + 12VNC - 65Bs	n using mRX B, cellular b B, cellular b cellular b scilon method scilon metho	dR feature selecti lastoderm; TMA, ágat proper prim ods on synthetio pre Eigen embryu 49 , 64 , 45 , 40 , 42 , 49 , 38 , 33 , 57		LeNet + mRMR 55 71 58 62 42 55 59 55 70	maller class only posterior endock t primordiam; E ive (M.O.) and Wavelet embryco 49 64 45 40 39 47 38 37 57	R <sub>4</sub> : overall recognition R <sub>4</sub> : overall recognition INSS, enhyponic centra iotations Wavelet embryo + mRMI 76 82 83 85 85 85 85 85 85 88 88 88 87

# Outline

- What is mRMR feature selection
- Applications in cancer classification
- Applications in image pattern recognition
- Theoretical basis of mRMR
- Combinations with other methods
- How to use mRMR programs













Factorize 
$$I(S_m, h)$$
  
• Relevance of  $S = \{x_1, x_2, ...\}$  and h, or  $R_L(S,h)$   
• Redundancy among variables  $\{x_1, x_2, ...\}$  or  $R_D(S)$   
 $R_L = \frac{1}{|S|} \sum_{x \in S} I(x_i, h) \qquad R_D = \frac{1}{|S|^2} \sum_{x,x_i \in S} I(x_i, x_i)$   
 $I(S_m, h) = J(S_{m-1}, x_m, h) - J(S_{m-1}, x_m).$   
• For incremental search, max  $I(S,h)$  is "equivalent" to max  $[R_L(S,h) - R_D(S)]$ , i.e. combination of min-Redundancy-Max-Relevance (mRMR).









#### **Use Wrappers to Refine Features**

- mRMR is a filter approach
  - Fast
  - Features might be redundant
  - Independent of the classifier
- Wrappers seek to minimize the number of errors directly
  - Slow
  - Features are less robust Dependent on classifier
  - Better prediction accuracy
- Use mRMR first to generate a short feature pool and use wrappers to get a least redundant feature set with better accuracy

37

41



### Outline

- What is mRMR feature selection
- Applications in cancer classification
- Applications in image pattern recognition
- Theoretical basis of mRMR
- Combinations with other methods
- How to use mRMR programs



## **Available Versions**

Matlab version

- All source codes. Can be embedded in your programs easily.
- C/C++ version
  - For Linux/Unix/Mac, simple command line executable.
- Online version
  - Upload the data sets (csv format: comma separated values) and get the results right away.

## **Available Datasets**

- NCI data (9 cancers, discretized as 3-states)
- Lung Cancer data (7 subtypes, discretized as 3-states)
- Lymphoma data (9 cancer-subtypes, discretized as 3-states)
- Leukemia data (2 subtypes, discretized as 3states)
- Colon Cancer data (2 subtypes, discretized as 3-states)
- The continuous-value raw data should be obtained from the original sources.

#### Citations

- [TPAMI05] Hanchuan Peng, Fuhui Long, and Chris Ding, "Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 8, pp.1226-1238, 2005.
  - 1238, 2005. This paper presents a theory of mutual information based feature selection. Demonstrates the relationship of four selection schemes: maximum dependency, mRMR, maximum relevance, and minimal redundancy. Also gives the combination scheme of "mRMR + wrapper" selection and mutual information estimation of continuous/hybrid variables.
- [JBCB05] Chris Ding, and Hanchuan Peng, "Minimum redundancy feature selection from microarray gene expression data," Journal of Bioinformatics and Computational Biology, Vol. 3, No. 2, pp.185-205, 2005.
  - This paper presents a comprehensive suite of experimental results of mRMR for microarray gene selection on many different conditions. It is an extended version of the CSB03 paper.
- [CSB03] Chris Ding, and Hanchuan Peng, "Minimum redundancy feature selection from microarray gene expression data," Proc. 2nd IEEE Computational Systems Bioinformatics Conference (CSB 2003), pp.523-528, Stanford, CA, Aug, 2003.
   This paper presents the first set of mRMR results and different definitions of relevance/redundancy terms.
- [Bioinfo07] Jie Zhou, and Hanchuan Peng, "Automatic recognition and annotation of gene expression patterns of fly embryos," Bioinformatics, Vol. 23, No. 5, pp. 589-596, 2009.

43

45

One application of mRMR in selecting good wavelet image features.

#### Conclusions

- The Max-Dependency feature selection can be efficiently implemented as the mRMR algorithm.
- Significantly outperforms the widely used max-relevance selection method: mRMR features cover a broader feature space with less features.
- mRMR is very efficient and useful for gene selection and many other applications. The programs are ready!

mRMR website: http://research.janelia.org/peng/proj/mRMR

